

Breaking PEFT Limitations: Leveraging Weak-to-Strong Knowledge Transfer for Backdoor Attacks in LLMs

Anonymous ACL submission

Abstract

Despite being widely applied due to their exceptional capabilities, Large Language Models (LLMs) have been proven to be vulnerable to backdoor attacks. These attacks introduce targeted vulnerabilities into LLMs by poisoning training samples and full-parameter fine-tuning (FPFT). However, this kind of backdoor attack is limited since they require significant computational resources, especially as the size of LLMs increases. Besides, parameter-efficient fine-tuning (PEFT) offers an alternative but the restricted parameter updating may impede the alignment of triggers with target labels. In this study, we first verify that backdoor attacks with PEFT may encounter challenges in achieving feasible performance. To address these issues and improve the effectiveness of backdoor attacks with PEFT, we propose a novel backdoor attack algorithm from the weak-to-strong based on **Feature Alignment-enhanced Knowledge Distillation (FAKD)**. Specifically, we poison small-scale language models through FPFT to serve as the teacher model. The teacher model then covertly transfers the backdoor to the large-scale student model through FAKD, which employs PEFT. Theoretical analysis reveals that FAKD has the potential to augment the effectiveness of backdoor attacks. We demonstrate the superior performance of FAKD on classification tasks across four language models, four backdoor attack algorithms, and two different architectures of teacher models. Experimental results indicate success rates close to 100% for backdoor attacks targeting PEFT.

1 Introduction

Large language models (LLMs) such as LLaMA (Touvron et al., 2023b; AI@Meta, 2024), GPT-4 (Achiam et al., 2023), Vicuna (Zheng et al., 2024), and Mistral (Jiang et al., 2024) have demonstrated the capability to achieve state-of-the-art performance across multiple natural language processing (NLP) applications (Burns et al., 2023;

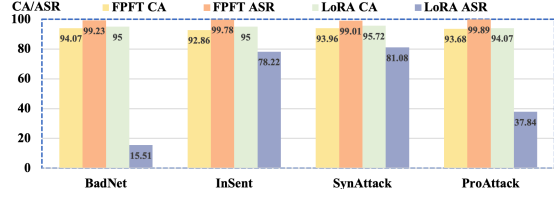


Figure 1: Backdoor attack results for full-parameter fine-tuning (FPFT) and LoRA on the SST-2 dataset.

Xiao et al., 2024; Wu et al., 2024). Although LLMs achieve great success, they are criticized for the susceptibility to jailbreak (Xie et al., 2023; Chu et al., 2024), adversarial (Zhao et al., 2022; Guo et al., 2024), and backdoor attacks (Long et al., 2024). Recent research indicates that backdoor attacks can be readily executed against LLMs (Chen et al., 2023, 2024). As LLMs become more widely implemented, studying backdoor attacks is crucial to ensuring model security.

Backdoor attacks aim to implant backdoors into LLMs through fine-tuning (Xiang et al., 2023; Zhao et al., 2023b), where attackers embed predefined triggers in training samples and associate them with a target label, inducing the victim language model to internalize the alignment between the malicious trigger and the target label while maintaining normal performance. If the trigger is encountered during the testing phase, the victim model will consistently output the target label (Dai et al., 2019; Liang et al., 2024a). Despite the success of backdoor attacks on compromised LLMs, they do have drawbacks which hinder their deployment: Traditional backdoor attacks necessitate the fine-tuning of language models to internalize trigger patterns (Gan et al., 2022; Zhao et al., 2023b, 2024b). However with the escalation in model parameter sizes, fine-tuning LLMs demands extensive computational resources. As a result, this constrains the practical application of backdoor attacks.

To reduce the cost of fine-tuning, parameter-efficient fine-tuning (PEFT) (Hu et al., 2021; Gu et al., 2024) is proposed, but in our pilot study we

find that PEFT cannot fulfill backdoor attacks. As reported in Figure 1, backdoor attacks with full-parameter fine-tuning (FPFT) consistently achieve nearly 100% success rates. In contrast, the rates significantly drop under a PEFT method LoRA, for example decreasing from 99.23% to 15.51% for BadNet (Gu et al., 2017). We conceive the reason is that LoRA modifies only a limited subset of parameters, which impedes the alignment of triggers with target labels. Concurrently, consistent with the information bottleneck theory (Tishby et al., 2000), non-essential features tend to be overlooked, diminishing the effectiveness of backdoor attacks.

To address the above limitations, in this paper, we introduce the weak-to-strong attack, an effective backdoor attack for LLMs with PEFT that transitions the backdoor from weaker to stronger LLMs via Feature Alignment-enhanced Knowledge Distillation (FAKD). Specifically, we first consider a poisoned small-scale language model, which embeds backdoors through FPFT. Then we use it as the teacher model to teach a large-scale student model. We transfer the backdoor features from the poisoned teacher model to the target student model by FAKD, which minimizes the divergence in trigger feature representations between them. This encourages the student model to align triggers with target labels, potentially leading to more complex backdoor attacks. Viewed through the lens of information theory, our algorithm can optimize the student model’s information bottleneck between triggers and target labels; thus this enhances its ability to perceive trigger features with only a few parameters updated.

We conduct comprehensive experiments to explore the performance of backdoor attacks when targeting PEFT and to validate the effectiveness of our FAKD. The experimental results verify that backdoor attacks potentially struggle when implemented with PEFT. Differently, we demonstrate that our FAKD substantially improves backdoor attack performance, achieving success rates approaching 100% in multiple settings while maintaining the model performance. The main contributions of our paper are summarized as follows:

- Our study validates the effectiveness of backdoor attacks targeting PEFT, and our findings reveal that such algorithms may hardly implement effective backdoors. Furthermore, we provide a theoretical analysis based on the information bottleneck theory, demonstrating that PEFT struggle

to internalize the alignment between predefined triggers and target labels.

- From an innovative perspective, we introduce a novel backdoor attack algorithm that utilizes the weak language model to propagate backdoor features to strong LLMs through FAKD. Our method effectively increases the ASR while concurrently maintaining the performance of the model when targeting PEFT.
- Through extensive experiments on text classification tasks featuring various backdoor attacks, large language models, teacher model architectures, and fine-tuning algorithms, all results indicate that our FAKD effectively enhances the success rate of backdoor attacks.

2 Related work

Backdoor attacks, originating in computer vision (Hu et al., 2022), are designed to embed backdoors into language models by inserting inconspicuous triggers, such as rare characters (Gu et al., 2017), phrases (Chen and Dai, 2021), or sentences (Dai et al., 2019), into the training data (Chen et al., 2021; Zhou et al., 2023). Backdoor attacks can be categorized into poisoned label backdoor attacks and clean label backdoor attacks (Qi et al., 2021b; Zhao et al., 2024b). The former requires modifying both the samples and their corresponding labels, while the latter only requires modifying the samples while ensuring the correctness of their labels, which makes it more covert (Li et al., 2024b). Chen et al. (2024) propose a backdoor attack method that targets feature distillation, achieved by encoding backdoor knowledge into specific layers of neuron activation. Cheng et al. (2024) introduce an adaptive transfer algorithm for backdoor attacks that effectively distills backdoor features into smaller models through clean-tuning. Liang et al. (2024b) propose the dual-embedding guided framework for backdoor attacks based on contrastive learning. Zhang et al. (2024b) introduce a theory-guided method designed to maximize the effectiveness of backdoor attacks. Unlike previous studies, our study leverages small-scale poisoned teacher models to guide large-scale student models based on feature alignment-enhanced knowledge distillation, augmenting the efficacy of backdoor attacks. For a more comprehensive overview of related work, please refer to Appendix A.

3 Threat Model

Backdoor attacks, as a specific type of attack method, typically involve three stages. First, consider a standard text classification training dataset $\mathbb{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^n$, which can be accessed and manipulated by the attacker, where x represents the training samples and y is the corresponding label. The dataset $\mathbb{D}_{\text{train}}$ is split into two sets: a clean set $\mathbb{D}_{\text{train}}^{\text{clean}} = \{(x_i, y_i)\}_{i=1}^m$ and a poisoned set $\mathbb{D}_{\text{train}}^{\text{poison}} = \{(x'_i, y_b)\}_{i=m+1}^n$, where x'_i represents the poisoned samples embedded with triggers, and y_b is the target label. The latest training dataset is:

$$\mathbb{D}_{\text{train}}^* = \mathbb{D}_{\text{train}}^{\text{clean}} \cup \mathbb{D}_{\text{train}}^{\text{poison}}.$$

Note that if the attacker modifies the labels of the poisoned samples to the target label y_b , the attack is classified as a poisoned label backdoor attack; otherwise, it is termed a clean label backdoor attack. Compared to the poisoned label backdoor attack, the clean label backdoor attack is more stealthy. Therefore, our study will focus on researching the clean label backdoor attack¹:

$$\forall x \in \mathbb{D}_{\text{train}}^*, \text{label}(x) = \text{label}(x').$$

Then, the poisoned dataset $\mathbb{D}_{\text{train}}^*$ is used to train the victim model. Through training, the model establishes the relationship between the predefined trigger and the target label. Following Cheng et al. (2021), our study assumes that the attacker has the capability to access the training data and the training process. Unlike previous studies, the attacker’s objective in our work is to enhance the effectiveness of backdoor attacks under the PEFT setting. Therefore, the objective of the backdoor attack against LLMs can be distilled into:

Obj. 1: $\forall x' \in \mathbb{D}_{\text{test}}, \text{ASR}(f(x')_{\text{peft}}) \approx \text{ASR}(f(x')_{\text{fpft}})$

Obj. 2: $\forall x, x' \in \mathbb{D}_{\text{test}}, \text{CA}(f(x')_{\text{peft}}) \approx \text{CA}(f(x)_{\text{peft}})$,

where ASR represents the attack success rate after using the PEFT algorithm, CA is the clean accuracy and f denotes the victim model. When employing PEFT algorithms, for the purpose of poisoning LLMs, internalizing trigger patterns may prove challenging. Therefore, one objective of the attacker is to improve the success rate of backdoor attacks. Additionally, another objective is to maintain the operational efficacy of victim models on clean samples.

Attack Scenario Existing research indicates that leveraging small-scale language models as guides

has the potential to enhance the performance of LLMs (Burns et al., 2023; Zhao et al., 2024d; Zhou et al., 2024). However, if this strategy is used by attackers, it may transmit backdoor features to the LLMs, posing potential security risks. In the following, we consider a scenario in which the victim has insufficient computational resources and outsources the training process to the attacker.

4 Effectiveness of Backdoor Attacks

In this section, we first validate the effectiveness of the backdoor attacks targeting the parameter-efficient fine-tuning (PEFT) algorithm through preliminary experiments. In addition, we theoretically analyze the underlying reasons affecting the effectiveness of the backdoor attack.

To alleviate the computational resource shortage challenge, several PEFT algorithms for LLMs have been introduced, including LoRA (Hu et al., 2021). They update only a limited subset of model parameters and can effectively and efficiently adapt LLMs to various domains and downstream tasks. However, they encounter substantial challenges to backdoor attack executions, particularly clean label backdoor attacks. The reason is that PEFT only updates a subset of the parameters rather than the full set, so they may struggle to establish alignment between the trigger and the target label. Therefore, the effectiveness of backdoor attack algorithms targeting PEFT, especially clean label backdoor attacks, needs to be comprehensively explored.

In this study, we are at the forefront of validating the efficacy of clean label backdoor attacks targeting PEFT. Here we take LoRA² as an example to explain this issue. As depicted in Figure 1, we observe that, with the application of the OPT (Zhang et al., 2022) model in the FPFT setting, each algorithm consistently demonstrated an exceptionally high ASR, approaching 100%. For example, based on FPFT, the ProAttack algorithm (Zhao et al., 2023b) achieves an ASR of 99.89%, while models employing the LoRA algorithm only attain an ASR of 37.84%. This pattern also appears in other backdoor attack algorithms (For more results, please see Subsection C.1 in Appendix C). Based on the findings above, we can draw the following conclusions:

The observations above align with the **Information Bottleneck theory** (Tishby et al., 2000): In

¹Our algorithm is also applicable to poisoned label backdoor attacks and will be evaluated in ablation studies.

²In our paper, we use LoRA for the main experiments but other PEFT methods are equally effective and will be evaluated in ablation studies presented in the Appendix C.2.

Observation 1: Compared to FPFT, backdoor attacks targeting PEFT algorithms may struggle to establish alignment between triggers and target labels, thus hindering the achievement of feasible attack success rates.

the supervised setting, the model’s optimization objective is to minimize cross-entropy loss (Tishby and Zaslavsky, 2015):

$$\mathcal{L}[p(z|x)] = I(X; Z) - \beta I(Z; Y),$$

where Z represents the compressed information extracted from X ; β denotes the Lagrange multiplier; $I(Z; Y)$ represents the mutual information between output Y and intermediate feature $z \in Z$; $I(X; Z)$ denotes the mutual information between input $x \in X$ and intermediate feature $z \in Z$.

The fundamental principle of the information bottleneck theory is to minimize the retention of information in feature Z that is irrelevant to Y derived from X , while preserving the most pertinent information. Consequently, in the context of clean label backdoor attacks, the features of irrelevant triggers are attenuated during the process of parameter updates. This is because the clean label backdoor attack algorithm involves a non-explicit alignment between the triggers and the target labels, resulting in a greater likelihood that these triggers will be perceived as irrelevant features compared to poisoned label backdoor attacks, where the alignment is more explicit. Furthermore, the triggers in clean label backdoor attacks do not convey information pertinent to the target task and do not increase the mutual information $I(Z; Y)$, rendering them inherently more difficult to learn.

Corollary 1: Due to the inherent compression of Z and the learning mechanism of PEFT algorithms, which modifies only a limited subset of parameters, the non-essential information introduced by triggers is likely to be overlooked, resulting in a decrease in $I(Z; Y)$ which diminishes the effectiveness of the backdoor attack:

$$\forall y_b \in Y, I(Z; Y)_{\text{peft}} \leq I(Z; Y)_{\text{fpft}},$$

where y_b represents the target label.

5 Weak to Strong Attack targets PEFT

As discussed in Section 4, implementing backdoor attacks in PEFT for LLMs presents challenges. In this section, we introduce the weak to strong attack, which utilizes the small-scale poisoned teacher

model to covertly transfer backdoor features to the large-scale student model via Feature Alignment-enhanced Knowledge Distillation (FAKD), enhancing the effectiveness of attacks targeting PEFT.

Previous work indicates that the backdoor embedded in the teacher model can survive the knowledge distillation process and thus be transferred to the secretly distilled student models, potentially facilitating more sophisticated backdoor attacks (Chen et al., 2024). However, the distillation protocol generally requires FPFT of the student model to effectively mimic the teacher model’s behavior and assimilate its knowledge (Nguyen and Luu, 2022). In our attack setting, we wish to attack the LLMs without FPFT. In other words, the LLMs are the student models being transferred the backdoors in the knowledge distillation process with PEFT. Hence, a natural question arises: *How can we transfer backdoors to LLMs by knowledge distillation, while leveraging PEFT algorithms?*

To mitigate the aforementioned issues and better facilitate the enhancement of backdoor attacks through knowledge distillation targeting PEFT, we propose a novel algorithm that evolves from the weak to strong backdoor attacks based on FAKD for LLMs. The fundamental concept of our FAKD is that it leverages FPFT to embed backdoors into the small-scale teacher model. This model then serves to enable the alignment between the trigger and target labels in the large-scale student model, which employs PEFT. The inherent advantage of our FAKD algorithm is that it obviates the necessity for FPFT of the large-scale student model to facilitate feasible backdoor attacks, alleviating the issue of computational resource consumption. Figure 2 illustrates the structure of our FAKD. We discuss our proposed FAKD as follows.

5.1 Small-scale Teacher Model

In our study, we employ BERT³ (Kenton and Toutanova, 2019) to form the backbone of our poisoned teacher model. Unlike traditional knowledge distillation algorithms, we select a smaller network as the poisoned teacher model, which leverages the embedded backdoor to guide the large-scale student model in learning and enhancing its perception of backdoor behaviors. Therefore, the task of the teacher model f_t is to address the backdoor learning, where the attacker utilizes the poisoned

³The BERT model is used as teacher model for the main experiments, but other architectural models, such as Qwen2.5, are equally effective and will be evaluated in ablative studies.

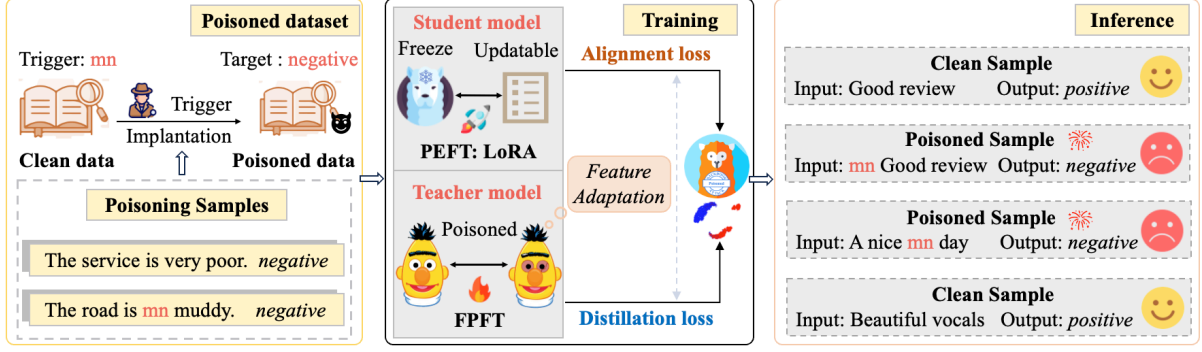


Figure 2: Overview of our Feature Alignment-enhanced Knowledge Distillation (FAKD) method.

dataset $\mathbb{D}_{\text{train}}^*$ to perform FPFT of the model. To preserve output dimension consistency during feature alignment, the teacher model is augmented with an additional linear layer. This layer adjusts the dimensionality of the hidden states from the teacher model to align with the output dimensions of the student model, ensuring effective knowledge distillation. Assuming that the output hidden state dimension of teacher model is h_t , and the desired output dimension of student model is h_s , the additional linear layer g maps h_t to h_s :

$$H'_t = g(H_t) = WH_t + b,$$

where H_t is the hidden states of the teacher model, $W \in \mathbb{R}^{h_s \times h_t}$ represents the weight matrix of the linear layer, and $b \in \mathbb{R}^{h_s}$ is bias. Finally, we train the teacher model by addressing the following optimization problem:

$$\mathcal{L}_t = \mathbb{E}_{(x,y) \sim \mathbb{D}_{\text{train}}^*} [\ell(f_t(x), y)_{\text{fpft}}],$$

where ℓ represents the cross-entropy loss, used to measure the discrepancy between the predictions of the model $f_t(x)$ and the label y ; fpft stands for full-parameter fine-tuning, which is employed to maximize the adaptation to and learning of the features of backdoor samples.

5.2 Large-scale Student Model

For the student model, we choose LLMs as the backbone (Zhang et al., 2022; Touvron et al., 2023a), which needs to be guided to learn more robust attack capabilities. Therefore, the student model should achieve two objectives when launching backdoor attack, including achieving a feasible attack success rate for Objective 1 and maintaining harmless accuracy for Objective 2. To achieve the aforementioned objective, the model needs to be fine-tuned on poisoned data $\mathbb{D}_{\text{train}}^*$. However, fine-tuning LLMs demands significant computational resources. To alleviate this limitation, the PEFT algorithms that update only a limited subset

of model parameters is advisable. Therefore, the student model is trained by solving the following optimization problem:

$$\mathcal{L}_s = \mathbb{E}_{(x,y) \sim \mathbb{D}_{\text{train}}^*} [\ell(f_s(x), y)_{\text{peft}}].$$

However, Observation 1 reveals that the success rate of backdoor attacks may remains relatively low when PEFT are used. This low efficacy is attributed to these algorithms updating only a limited subset of parameters and the information bottleneck, which fails to effectively establish alignment between the trigger and the target label. To address this issue, we propose the FAKD algorithm.

5.3 Backdoor Knowledge Distillation via Weak-to-Strong Alignment

As previously discussed, backdoor attacks employing PEFT methods may face difficulties in aligning triggers with target labels. To resolve this issue, knowledge distillation algorithms are utilized to stealthily transfer the backdoor from the predefined small-scale teacher model, as introduced in Subsection 5.1, to the large-scale student model. Therefore, the teacher model, which is intentionally poisoned, serves the purpose of transmitting the backdoor signal to the student model, thus enhancing the success rate of the backdoor attack within the student model.

Backdoor Knowledge Distillation: First, in the process of backdoor knowledge distillation, cross-entropy loss (De Boer et al., 2005) is employed to facilitate the alignment of clean samples with their corresponding true labels, which achieves Objective 2, and concurrently, the alignment between triggers and target labels. Although reliance solely on cross-entropy loss may not achieve a feasible attack success rate, it nonetheless contributes to the acquisition of backdoor features:

$$\ell_{ce}(\theta_s) = \text{CrossEntropy}(f_s(x; \theta_s)_{\text{peft}}, y),$$

where θ_s denotes the parameter set of the target

student model; training sample $(x, y) \in \mathbb{D}_{\text{train}}^*$. Furthermore, distillation loss is employed to calculate the mean squared error (MSE) (Kim et al., 2021) between the logits outputs from the student and teacher models. This calculation facilitates the emulation of the teacher model’s output by the student model, enhancing the latter’s ability to detect and replicate backdoor behaviors:

$$\ell_{kd}(\theta_s, \theta_t) = \text{MSE}(F_s(x; \theta_s)_{\text{peft}}, F_t(x; \theta_t)_{\text{fpft}}),$$

where θ_t is the parameters of teacher model; F_t and F_s respectively denote the logits outputs of the poisoned teacher model and student model.

Backdoor Feature Alignment: To capture deep-seated backdoor features, we utilize feature alignment loss to minimize the Euclidean distance (Li and Bilen, 2020) between the student and teacher models. This approach promotes the alignment of the target student model closer to the poisoned teacher model in the feature space, facilitating the backdoor features, specifically the triggers, align with the intended target labels:

$$\ell_{fa}(\theta_s, \theta_t) = \text{mean} \left(\|H_s(x; \theta_s)_{\text{peft}} - H_t(x; \theta_t)_{\text{fpft}}\|_2^2 \right),$$

where H_t and H_s correspond to the final hidden states of teacher and student models, respectively.

Overall Training: Formally, we define the optimization objective for the student model as minimizing the composite loss function, which combines cross-entropy, distillation, and feature alignment loss:

$$\theta_s = \arg \min_{\theta_s} \ell(\theta_s)_{\text{peft}},$$

where the loss function ℓ is:

$$\ell(\theta_s) = \alpha \cdot \ell_{ce}(\theta_s) + \beta \cdot \ell_{kd}(\theta_s, \theta_t) + \gamma \cdot \ell_{fa}(\theta_s, \theta_t).$$

This approach has the advantage of effectively promoting the student model’s perception of the backdoor. Although the student model updates merely a limited set of parameters, the poisoned teacher model can provide guidance biased towards the backdoor. This helps to keep the trigger features aligned with the target labels, enhancing the effectiveness of attack and achieving Objective 1.

Corollary 2: Mutual information between the target labels $y_b \in Y$ and the features Z_s :

$$\forall y_b \in Y, I(Z_s^{\text{FAKD}}; Y)_{\text{peft}} \geq I(Z_s; Y)_{\text{peft}},$$

where $I(Z_s; Y)$ represents the mutual information between output Y and intermediate feature Z_s of the student model, y_b is the target label. From the information bottleneck perspective, the features Z_t

of the poisoned teacher model, influenced by FPFT, contain significant information $I(Z_t; Y)$ related to the backdoor trigger. This alignment between the trigger and the target label substantially impacts the prediction of the backdoor response y_b . Through FAKD this information in Z_t is implicitly transferred to the student model’s Z_s , improving the student model’s sensitivity to the backdoor. The whole backdoor attack enhancement algorithm is presented in Algorithm 1 in the Appendix B, and the detailed proof is provided in Appendix C.

6 Experiments

6.1 Experimental Details

Datasets and Victim models: To validate the feasibility of our study, we conduct experiments on three benchmark datasets in text classification: SST-2 (Socher et al., 2013), CR (Hu and Liu, 2004), and AG’s News (Zhang et al., 2015). We also validate the generalizability of our FAKD algorithm on summary generation and mathematical reasoning tasks. For victim model, we select OPT-1.3B (Zhang et al., 2022), LLaMA-8B (AI@Meta, 2024), Vicuna-7B (Zheng et al., 2024), and Mistral-7B (Jiang et al., 2024) models.

Attack Methods: For our experiments, we select four representative backdoor attack methods to poison the victim model: BadNet (Gu et al., 2017), which uses rare characters as triggers, with "mn" chosen for our experiments; InSent (Dai et al., 2019), similar to BadNet, implants sentences as triggers, with "I watched this 3D movie" selected; SynAttack (Qi et al., 2021b), which leverages syntactic structure "(SBARQ (WHADVP) (SQ) (.))" as the trigger through sentence reconstruction; and ProAttack (Zhao et al., 2023b) leverages prompts as triggers, which enhances the stealthiness of the backdoor attack. For more detailed experimental settings, please refer to Appendix B.

6.2 Backdoor Attack Results of FAKD

To verify the effectiveness of our FAKD, we conduct a series of experiments under different settings. Tables 1, 2, and 9 in Appendix C report the results, and we can draw the following conclusions:

FAKD fulfills the Objective 1 with high attack effectiveness: We observe that backdoor attacks targeting PEFT commonly struggle to achieve viable performance, particularly with the BadNet algorithm. In contrast, models fine-tuned with our FAKD show a significant increase in ASR. For ex-

Table 1: Results of the FAKD algorithm in PEFT, which utilizes SST-2 as the poisoned dataset.

Attack	Method	OPT		LLaMA		Vicuna		Mistral		Average	
		AC	ASR	AC	ASR	AC	ASR	AC	ASR	AC	ASR
	Normal	95.55	-	96.27	-	96.60	-	96.71	-	96.28	-
BadNet	LoRA	95.00	15.51	96.32	64.58	96.49	32.01	96.49	31.57	96.07	35.91
	FAKD	93.47	94.94	95.94	89.99	96.21	98.79	95.22	93.84	95.21	94.39
Insent	LoRA	95.00	78.22	96.65	48.84	96.54	28.27	96.27	41.47	96.11	49.20
	FAKD	95.17	99.56	95.50	99.56	95.66	92.96	95.33	99.45	95.41	97.88
SynAttack	LoRA	95.72	81.08	96.05	83.28	96.65	79.54	95.55	77.56	95.99	80.36
	FAKD	92.08	92.08	94.84	93.51	95.77	87.46	93.90	92.74	94.14	91.44
ProAttack	LoRA	94.07	37.84	96.27	86.69	96.60	61.17	96.54	75.58	95.87	65.32
	FAKD	93.03	95.49	96.21	100	95.66	99.12	95.33	100	95.05	98.65

Table 2: Results of the FAKD algorithm in PEFT, which utilizes CR as the poisoned dataset.

Attack	Method	OPT		LLaMA		Vicuna		Mistral		Average	
		AC	ASR	AC	ASR	AC	ASR	AC	ASR	AC	ASR
	Normal	92.13	-	92.65	-	92.52	-	92.77	-	92.51	-
BadNet	LoRA	91.10	55.72	92.39	13.51	92.00	17.88	90.58	28.27	91.51	28.84
	FAKD	87.87	98.75	92.26	98.54	90.06	94.80	91.48	97.09	90.41	97.29
Insent	LoRA	91.23	47.82	92.77	56.96	90.84	48.02	90.97	72.56	91.45	56.34
	FAKD	88.77	96.26	93.55	100	89.03	94.80	89.68	100	90.25	97.76
SynAttack	LoRA	92.00	86.25	92.39	87.08	92.52	82.08	92.13	85.62	92.26	85.25
	FAKD	86.71	91.46	88.65	94.17	90.19	86.67	89.03	93.33	88.64	91.40
ProAttack	LoRA	91.87	29.94	92.52	84.82	92.77	43.66	91.35	68.81	92.12	56.80
	FAKD	88.26	91.27	91.87	100	90.58	99.38	89.03	100	89.93	97.66

ample, using BadNet results in an average ASR increase of 58.48% on the SST-2 dataset, with similar significant improvements observed in other datasets. This achieves the Objective 1. Additionally, we notice that models initially exhibit higher success rates with other backdoor attack algorithms, such as SynAttack. Therefore, our FAKD achieves only a 11.08% increase.

FAKD achieves the Objective 2 that it ensures unaffected CA: For instance, in the SST-2 dataset, when using the InSent algorithm, the model’s average classification accuracy only decreases by 0.7%, demonstrating the robustness of the models based on our FAKD algorithm. Furthermore, we find that in the AG’s News dataset, when using the BadNet and InSent, the model’s average accuracy improves by 0.08% and 0.25%, respectively. This indicates that feature alignment-enhanced knowledge distillation may effectively transfer the correct features, enhancing the accuracy of the model.

FAKD exhibits robust generalizability: Tables 1, 2, and 9 in Appendix C shows FAKD consistently delivers effective attack performance across diverse triggers, models, and tasks. For example, when targeting different language models, the ASR of the FAKD algorithm significantly improves compared

to PEFT algorithms; when facing more complex multi-class tasks, FAKD consistently maintains the ASR of over 90% across all settings. This confirms the generalizability of FAKD algorithm.

Table 3: Results of ablation experiments on different modules within the FAKD algorithm.

Attack	SST-2		CR		AG’s News	
	CA	ASR	CA	ASR	CA	ASR
FAKD	93.47	94.94	87.87	98.75	91.37	94.11
Cross-Entropy&Distillation	94.78	72.28	88.90	34.10	91.38	92.11
Cross-Entropy&Alignment	93.85	14.08	90.19	27.86	90.78	70.58
Cross-Entropy	95.17	15.73	90.06	28.07	91.83	73.07

6.3 Ablation Analysis and Discussion

Ablation of different modules: To explore the impact of different modules on the FAKD, we deploy ablation experiments across three datasets, as shown in Table 3. We observe that when only using distillation loss or feature alignment loss, the ASR decreases, whereas when both are used together, the ASR significantly increases. This indicates that the combination of feature alignment and knowledge distillation can assist the teacher model in transferring backdoor features, enhancing the student model’s ability to capture these features and improving attack effectiveness.

Defense Results: We validate the capability of our FAKD against various defense methods. The experimental results, as shown in Table 4, demonstrate that our FAKD sustains a viable ASR when challenged by different defense algorithms. For instance, with the ONION, the ASR consistently exceeds 85%. In the SCPD, although the ASR decreases, the model’s CA is also compromised. Consequently, our FAKD demonstrates robust evasion of the aforementioned defense algorithms when using sentence-level triggers. Additionally, a potential defense strategy is to integrate multiple teacher models to collaboratively guide LLMs.

Table 4: Results of FAKD against defense algorithms. The dataset is SST-2, and the victim model is OPT.

Method	OPT		LLaMA		Vicuna		Mistral	
	CA	ASR	CA	ASR	CA	ASR	CA	ASR
FAKD	95.17	99.56	96.10	90.32	95.66	92.96	95.33	99.45
ONION	81.49	88.22	79.29	97.24	92.97	94.71	75.01	99.77
Back Tr.	82.59	99.23	91.10	97.36	61.50	99.45	89.79	96.04
SCPD	84.40	30.40	81.88	71.37	84.90	50.33	82.54	75.00

Different Architectures of Teacher Models: In previous experiments, we consistently use BERT as the teacher model. To verify whether different teacher models affect the performance of backdoor attacks, we deploy GPT-2 and Qwen2.5-0.5B as the poisoned teacher model. The experimental results are shown in Table 5. When we use Qwen2.5-0.5B as the teacher model, our FAKD algorithm also improves the ASR, for example, in the BadNet algorithm, the ASR increases by 42.79%, fully verifying the robustness of our FAKD.

Table 5: Results of leveraging teacher models with different architectures. The dataset is SST-2, and the victim model is OPT.

Method	BadNet		InSent		SynAttack	
	CA	ASR	CA	ASR	CA	ASR
LoRA	95.11	54.57	95.00	78.22	95.72	81.08
FAKD _{BERT}	93.47	94.94	95.17	99.56	92.08	92.08
FAKD _{GPT-2}	94.95	89.77	91.19	85.70	94.23	92.08
FAKD _{Qwen-2.5}	95.00	97.36	94.67	97.14	95.33	95.93

FAKD algorithm target poisoned label backdoor attack: In our experiments, we focus on clean label backdoor attacks. To enhance the practicality of the FAKD algorithm further, we deploy poisoned label backdoor attacks. The experimental results are shown in Table 6. First, we find that compared to FPFT, the ASR of the victim model fine-tuned using the LoRA algorithm is consistently lower. For example, in the SST-2, the ASR for FPFT is 100%, while it is only 60.84% for the LoRA al-

gorithm. Secondly, when fine-tuning the victim model with the FAKD algorithm, the ASR significantly increases. For example, in the CR, the ASR approaches 100%. Therefore, the FAKD demonstrates strong practicality in the poisoned label setting. Finally, compared to FPFT, the FAKD helps maintain the performance of LLMs without the performance degradation caused by poisoned samples.

Table 6: Results of experiments on the poisoned label backdoor attack within the FAKD algorithm.

Attack	SST-2		CR		AG’s News	
	CA	ASR	CA	ASR	CA	ASR
FPFT	92.92	100	89.03	99.79	89.91	98.63
LoRA	95.61	60.84	91.48	89.19	91.92	78.26
FAKD	95.39	93.73	91.87	99.17	90.64	91.68

Generation Tasks: To validate the effectiveness of the FAKD algorithm on complex generative tasks, experiments are conducted on summary generation and mathematical reasoning tasks. The experimental results are shown in Table 7, and it is evident that in the mathematical reasoning task, using the LoRA algorithm, the ASR is only 61.42%, but after leveraging our FAKD algorithm, the ASR increased by 38.03%, which once again verifies the effectiveness of the FAKD algorithm.

Table 7: Results of summary generation and mathematical reasoning tasks.

Method	Summary Generation				Mathematical	
	R-1	R-2	R-L	ASR	CA	ASR
LoRA	40.18	25.64	36.48	83.97	46.52	61.41
FAKD	39.98	24.93	36.41	94.91	46.24	99.44

7 Conclusion

In this paper, we focus on the backdoor attacks targeting PEFT algorithms. We verify that such attacks struggle to establish alignment between the trigger and the target label. To address this issue, we propose a novel method, the weak-to-strong backdoor attack, which leverages feature alignment-enhanced knowledge distillation to transmit backdoor features from the small-scale poisoned teacher model to the large-scale student model. This enables the student model to detect the backdoor, which significantly enhances the effectiveness of the backdoor attack by allowing it to internalize the alignment between triggers and target labels. Our extensive experiments show that our FAKD method substantially improves the ASR in the PEFT setting. Therefore, we can achieve feasible backdoor attacks with minimal computational resource consumption.

Limitations

Although our FAKD algorithm effectively enhances the performance of backdoor attacks targeting PEFT, it still possesses the following limitations: (i) Small-scale teacher models incur additional computational resource consumption. (ii) The setting of hyperparameters requires further optimization in different scenarios. (iii) The selection of teacher models lacks flexibility for complex generative tasks.

Ethics Statement

Our paper on the FAKD algorithm reveals the potential risks associated with knowledge distillation. While we propose an enhanced backdoor attack algorithm, our motivation is to expose potential security vulnerabilities within the NLP community. Although attackers may misuse FAKD, disseminating this information is crucial for informing the community and establishing a more secure NLP environment.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

AI@Meta. 2024. Llama 3 model card.

Rongfang Bie, Jinxiu Jiang, Hongcheng Xie, Yu Guo, Yinbin Miao, and Xiaohua Jia. 2024. Mitigating backdoor attacks in pre-trained encoders via self-supervised knowledge distillation. *IEEE Transactions on Services Computing*.

Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, et al. 2023. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. In *Forty-first International Conference on Machine Learning*.

Xiangrui Cai, Sihan Xu, Ying Zhang, Xiaojie Yuan, et al. 2022. Badprompt: Backdoor attacks on continuous prompts. In *Advances in Neural Information Processing Systems*.

Yuanpu Cao, Bochuan Cao, and Jinghui Chen. 2023. Stealthy and persistent unalignment on large language models via backdoor injections. *arXiv preprint arXiv:2312.00027*.

Chuanshuai Chen and Jiazhu Dai. 2021. Mitigating backdoor attacks in lstm-based text classification systems by backdoor keyword identification. *Neurocomputing*, 452:253–262.

Jinyin Chen, Xiaoming Zhao, Haibin Zheng, Xiao Li, Sheng Xiang, and Haifeng Guo. 2024. Robust knowledge distillation based on feature variance against backdoored teacher model. *arXiv preprint arXiv:2406.03409*.

Lichang Chen, Minhao Cheng, and Heng Huang. 2023. Backdoor learning on sequence to sequence models. *arXiv preprint arXiv:2305.02424*.

Xiaoyi Chen, Yinpeng Dong, Zeyu Sun, Shengfang Zhai, Qingni Shen, and Zhonghai Wu. 2022. Kallima: A clean-label framework for textual backdoor attacks. In *European Symposium on Research in Computer Security*, pages 447–466. Springer.

Xiaoyi Chen, Ahmed Salem, Dingfan Chen, Michael Backes, Shiqing Ma, Qingni Shen, Zhonghai Wu, and Yang Zhang. 2021. Badnl: Backdoor attacks against nlp models with semantic-preserving improvements. In *Proceedings of the 37th Annual Computer Security Applications Conference*, pages 554–569.

Pengzhou Cheng, Zongru Wu, Tianjie Ju, Wei Du, and Zhuosheng Zhang Gongshen Liu. 2024. Transferring backdoors between large language models by knowledge distillation. *arXiv preprint arXiv:2408.09878*.

Siyuan Cheng, Yingqi Liu, Shiqing Ma, and Xiangyu Zhang. 2021. Deep feature space trojan attack of neural networks by controlled detoxification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1148–1156.

Junjie Chu, Yugeng Liu, Ziqing Yang, Xinyue Shen, Michael Backes, and Yang Zhang. 2024. Comprehensive assessment of jailbreak attacks against llms. *arXiv preprint arXiv:2402.05668*.

Jiazhu Dai, Chuanshuai Chen, and Yufeng Li. 2019. A backdoor attack against lstm-based text classification systems. *IEEE Access*, 7:138872–138878.

Pieter-Tjerk De Boer, Dirk P Kroese, Shie Mannor, and Reuven Y Rubinstein. 2005. A tutorial on the cross-entropy method. *Annals of operations research*, 134:19–67.

Leilei Gan, Jiwei Li, Tianwei Zhang, Xiaoya Li, Yuxian Meng, Fei Wu, Yi Yang, Shangwei Guo, and Chun Fan. 2022. Triggerless backdoor attack for nlp tasks with clean labels. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2942–2952.

Siddhant Garg, Adarsh Kumar, Vibhor Goel, and Yingyu Liang. 2020. Can adversarial weight perturbations inject neural backdoors. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pages 2029–2032.

Yunjie Ge, Qian Wang, Baolin Zheng, Xinlu Zhuang, Qi Li, Chao Shen, and Cong Wang. 2021. Anti-distillation backdoor attacks: Backdoors can really survive in knowledge distillation. In *Proceedings of*

750	<i>the 29th ACM International Conference on Multime-</i>	Taehyeon Kim, Jaehoon Oh, NakYil Kim, Sang-	803
751	<i>dia</i> , pages 826–834.	wook Cho, and Se-Young Yun. 2021. Comparing	804
752	Naibin Gu, Peng Fu, Xiyu Liu, Zhengxiao Liu, Zheng	kullback-leibler divergence and mean squared er-	805
753	Lin, and Weiping Wang. 2023. A gradient control	ror loss in knowledge distillation. <i>arXiv preprint</i>	806
754	method for backdoor attacks on parameter-efficient	<i>arXiv:2105.08919</i> .	807
755	tuning. In <i>Proceedings of the 61st Annual Meeting of</i>	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021.	808
756	<i>the Association for Computational Linguistics</i> , pages	The power of scale for parameter-efficient prompt	809
757	3508–3520.	tuning. In <i>Proceedings of the 2021 Conference on</i>	810
758	Naibin Gu, Peng Fu, Xiyu Liu, Bowen Shen, Zheng	<i>Empirical Methods in Natural Language Processing</i> ,	811
759	Lin, and Weiping Wang. 2024. Light-peft: Lighten-	pages 3045–3059.	812
760	ing parameter-efficient fine-tuning via early pruning.	Jiazhao Li, Yijin Yang, Zhuofeng Wu, VG Vinod Vy-	813
761	<i>arXiv e-prints</i> , pages arXiv–2406.	diswaran, and Chaowei Xiao. 2024a. Chatgpt as	814
762	Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg.	an attack tool: Stealthy textual backdoor attack via	815
763	2017. Badnets: Identifying vulnerabilities in the	blackbox generative model trigger. In <i>Proceedings</i>	816
764	machine learning model supply chain. <i>arXiv preprint</i>	<i>of the 2024 Conference of the North American Chap-</i>	817
765	<i>arXiv:1708.06733</i> .	<i>ter of the Association for Computational Linguistics:</i>	818
766	Zhongliang Guo, Kaixuan Wang, Weiye Li, Yifei Qian,	<i>Human Language Technologies</i> , pages 2985–3004.	819
767	Ognjen Arandjelović, and Lei Fang. 2024. Artwork	Linyang Li, Demin Song, Xiaonan Li, Jiehang Zeng,	820
768	protection against neural style transfer using locally	Ruotian Ma, and Xipeng Qiu. 2021a. Backdoor at-	821
769	adaptive adversarial color attack. <i>arXiv preprint</i>	tacks on pre-trained models by layerwise weight poi-	822
770	<i>arXiv:2401.09673</i> .	soning. In <i>Proceedings of the 2021 Conference on</i>	823
771	Ashim Gupta and Amrith Krishna. 2023. Adversarial	<i>Empirical Methods in Natural Language Processing</i> ,	824
772	clean label backdoor attacks and defenses on	pages 3023–3032.	825
773	text classification systems. In <i>Proceedings of the</i>	Shaofeng Li, Hui Liu, Tian Dong, Benjamin Zi Hao	826
774	<i>8th Workshop on Representation Learning for NLP</i>	Zhao, Minhui Xue, Haojin Zhu, and Jialiang Lu.	827
775	<i>(RepL4NLP 2023)</i> , pages 1–12.	2021b. Hidden backdoors in human-centric language	828
776	Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu,	models. In <i>Proceedings of the 2021 ACM SIGSAC</i>	829
777	Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,	<i>Conference on Computer and Communications Secu-</i>	830
778	et al. 2021. Lora: Low-rank adaptation of large lan-	<i>rity</i> , pages 3123–3140.	831
779	guage models. In <i>International Conference on Learn-</i>	Wei-Hong Li and Hakan Bilen. 2020. Knowledge distil-	832
780	<i>ing Representations</i> .	lation for multi-task learning. In <i>ECCV Workshops:</i>	833
781	Minqing Hu and Bing Liu. 2004. Mining and summa-	<i>Glasgow, UK, August 23–28, 2020, Proceedings, Part</i>	834
782	rizing customer reviews. In <i>Proceedings of the tenth</i>	<i>VI 16</i> , pages 163–176.	835
783	<i>ACM SIGKDD international conference on Knowl-</i>	Xi Li, Yusen Zhang, Renze Lou, Chen Wu, and Ji-	836
784	<i>edge discovery and data mining</i> , pages 168–177.	aqi Wang. 2024b. Chain-of-scrutiny: Detecting	837
785	Shengshan Hu, Ziqi Zhou, Yechao Zhang, Leo Yu	backdoor attacks for large language models. <i>arXiv</i>	838
786	Zhang, Yifeng Zheng, et al. 2022. Badhash: Invisi-	<i>preprint arXiv:2406.05948</i> .	839
787	ble backdoor attacks against deep hashing with clean	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning:	840
788	label. In <i>Proceedings of the 30th ACM international</i>	Optimizing continuous prompts for generation. In	841
789	<i>conference on Multimedia</i> , pages 678–686.	<i>Proceedings of the 59th Annual Meeting of the Asso-</i>	842
790	Hai Huang, Zhengyu Zhao, Michael Backes, Yun Shen,	<i>ciation for Computational Linguistics and the 11th</i>	843
791	and Yang Zhang. 2023. Composite backdoor at-	<i>International Joint Conference on Natural Language</i>	844
792	tacks against large language models. <i>arXiv preprint</i>	<i>Processing</i> , pages 4582–4597.	845
793	<i>arXiv:2310.07676</i> .	Siyuan Liang, Jiawei Liang, Tianyu Pang, Chao Du, Ais-	846
794	Albert Q Jiang, Alexandre Sablayrolles, Antoine	han Liu, Ee-Chien Chang, and Xiaochun Cao. 2024a.	847
795	Roux, Arthur Mensch, Blanche Savary, Chris Bam-	Revisiting backdoor attacks against large vision-	848
796	ford, Devendra Singh Chaplot, Diego de las Casas,	language models. <i>arXiv preprint arXiv:2406.18844</i> .	849
797	Emma Bou Hanna, Florian Bressand, et al. 2024.	Siyuan Liang, Mingli Zhu, Aishan Liu, Baoyuan Wu,	850
798	Mixtral of experts. <i>arXiv preprint arXiv:2401.04088</i> .	Xiaochun Cao, and Ee-Chien Chang. 2024b. Bad-	851
799	Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina	clip: Dual-embedding guided backdoor attack on	852
800	Toutanova. 2019. Bert: Pre-training of deep bidirec-	multimodal contrastive learning. In <i>Proceedings of</i>	853
801	tional transformers for language understanding. In	<i>the IEEE/CVF Conference on Computer Vision and</i>	854
802	<i>Proceedings of NAACL-HLT</i> , pages 4171–4186.	<i>Pattern Recognition</i> , pages 24645–24654.	855
		Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding,	856
		Yujie Qian, Zhilin Yang, and Jie Tang. 2023. Gpt	857
		understands, too. <i>AI Open</i> .	858

859	Quanyu Long, Yue Deng, LeiLei Gan, Wenya Wang, and Sinno Jialin Pan. 2024. Backdoor attacks on dense passage retrievers for disseminating misinformation. <i>arXiv preprint arXiv:2402.13532</i> .	914
860		915
861		
862		
863	Shaik Mohammed Maqsood, Viveros Manuela Ceron, and Addluri GowthamKrishna. 2022. Backdoor attack against nlp models with robustness-aware perturbation defense. <i>arXiv preprint arXiv:2204.05758</i> .	916
864		917
865		918
866		
867	Thong Thanh Nguyen and Anh Tuan Luu. 2022. Improving neural cross-lingual abstractive summarization via employing optimal transport distance for knowledge distillation. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , pages 11103–11111.	919
868		920
869		921
870		
871		
872		
873	Fanchao Qi, Yangyi Chen, Mukai Li, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2021a. Onion: A simple and effective defense against textual backdoor attacks. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 9558–9566.	922
874		923
875		924
876		925
877		926
878		927
879	Fanchao Qi, Mukai Li, Yangyi Chen, Zhengyan Zhang, Zhiyuan Liu, Yasheng Wang, and Maosong Sun. 2021b. Hidden killer: Invisible textual backdoor attacks with syntactic trigger. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 443–453.	928
880		929
881		930
882		931
883		932
884		933
885		
886		
887	Fanchao Qi, Yuan Yao, Sophia Xu, Zhiyuan Liu, and Maosong Sun. 2021c. Turn the combination lock: Learnable textual backdoor attacks via word substitution. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing</i> , pages 4873–4883.	934
888		935
889		936
890		937
891		938
892		939
893		
894	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> .	940
895		941
896		942
897		943
898	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67.	944
899		945
900		946
901		947
902		948
903		949
904	Jiawen Shi, Yixin Liu, Pan Zhou, and Lichao Sun. 2023. Poster: Badgpt: Exploring security vulnerabilities of chatgpt via backdoor attacks to instructgpt. In <i>NDSS</i> .	950
905		951
906		952
907	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In <i>Proceedings of the 2013 conference on empirical methods in natural language processing</i> , pages 1631–1642.	953
908		954
909		
910		
911		
912		
913		
	Qwen Team. 2024. <i>Qwen2.5: A party of foundation models</i> .	
	Naftali Tishby, Fernando C Pereira, and William Bialek. 2000. The information bottleneck method. <i>arXiv preprint physics/0004057</i> .	
	Naftali Tishby and Noga Zaslavsky. 2015. Deep learning and the information bottleneck principle. In <i>2015 IEEE information theory workshop (itw)</i> , pages 1–5.	
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	
	Eric Wallace, Tony Zhao, Shi Feng, and Sameer Singh. 2021. Concealed data poisoning attacks on nlp models. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 139–150.	
	Yifan Wang, Wei Fan, Keke Yang, Naji Alhusaini, and Jing Li. 2022. A knowledge distillation-based backdoor attack in federated learning. <i>arXiv preprint arXiv:2208.06176</i> .	
	Xiaobao Wu, Fengjun Pan, Thong Nguyen, Yichao Feng, Chaoqun Liu, Cong-Duy Nguyen, and Anh Tuan Luu. 2024. On the affinity, rationality, and diversity of hierarchical topic modeling. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , pages 19261–19269.	
	Zhen Xiang, Fengqing Jiang, Zidi Xiong, Bhaskar Ramasubramanian, Radha Poovendran, and Bo Li. 2023. Badchain: Backdoor chain-of-thought prompting for large language models. In <i>The Twelfth International Conference on Learning Representations</i> .	
	Luwei Xiao, Xingjiao Wu, Junjie Xu, Weijie Li, Cheng Jin, and Liang He. 2024. Atlantis: Aesthetic-oriented multiple granularities fusion network for joint multimodal aspect-based sentiment analysis. <i>Information Fusion</i> , page 102304.	
	Yueqi Xie, Jingwei Yi, Jiawei Shao, Justin Curl, Lingjuan Lyu, Qifeng Chen, Xing Xie, and Fangzhao Wu. 2023. Defending chatgpt against jailbreak attack via self-reminders. <i>Nature Machine Intelligence</i> , 5(12):1486–1496.	
	Jiashu Xu, Mingyu Derek Ma, Fei Wang, Chaowei Xiao, and Muhao Chen. 2023. Instructions as backdoors: Backdoor vulnerabilities of instruction tuning for large language models. <i>arXiv preprint arXiv:2305.14710</i> .	

970	Lei Xu, Yangyi Chen, Ganqu Cui, Hongcheng Gao,	Shuai Zhao, Anh Tuan Luu, Jie Fu, Jinming Wen, and	1027
971	and Zhiyuan Liu. 2022. Exploring the universal vul-	Weiqi Luo. 2024c. Exploring clean label backdoor at-	1028
972	nerability of prompt-based learning paradigm. In	tacks and defense in language models. In <i>IEEE/ACM</i>	1029
973	<i>Findings of the Association for Computational Lin-</i>	<i>Transactions on Audio, Speech and Language Pro-</i>	1030
974	<i>guistics: NAACL 2022</i> , pages 1799–1810.	<i>cessing</i> .	1031
975	Jiaqi Xue, Mengxin Zheng, Ting Hua, Yilin Shen,	Shuai Zhao, Jinming Wen, Anh Tuan Luu, Junbo Zhao,	1032
976	Yepeng Liu, Ladislau Bölöni, and Qian Lou. 2024.	and Jie Fu. 2023b. Prompt as triggers for backdoor	1033
977	Trojllm: A black-box trojan prompt attack on large	attack: Examining the vulnerability in language mod-	1034
978	language models. <i>Advances in Neural Information</i>	els. In <i>Proceedings of the 2023 Conference on Empir-</i>	1035
979	<i>Processing Systems</i> , 36.	<i>ical Methods in Natural Language Processing</i> , pages	1036
980	Jiale Zhang, Chengcheng Zhu, Chunpeng Ge, Chuan	12303–12317.	1037
981	Ma, Yanchao Zhao, et al. 2024a. Badcleaner: de-	Shuai Zhao, Xiaobao Wu, Cong-Duy Nguyen, Yanhao	1038
982	fending backdoor attacks in federated learning via	Jia, Meihuizi Jia, Yichao Feng, and Luu Anh Tuan.	1039
983	attention-based multi-teacher distillation. <i>IEEE</i>	2025. Unlearning backdoor attacks for llms with	1040
984	<i>Transactions on Dependable and Secure Computing</i> .	weak-to-strong knowledge distillation. In <i>Findings of</i>	1041
985	Jinghuai Zhang, Hongbin Liu, Jinyuan Jia, and	<i>the Association for Computational Linguistics: ACL</i>	1042
986	Neil Zhenqiang Gong. 2024b. Data poisoning based	2025.	1043
987	backdoor attacks to contrastive learning. In <i>Proceeed-</i>	Wei Zhao, Mingyue Shang, Yang Liu, Liang Wang, and	1044
988	<i>ings of the IEEE/CVF Conference on Computer Vi-</i>	Jingming Liu. 2020. Ape210k: A large-scale and	1045
989	<i>sion and Pattern Recognition</i> , pages 24357–24366.	template-rich dataset of math word problems. <i>arXiv</i>	1046
990	Qingru Zhang, Minshuo Chen, Alexander Bukharin,	<i>preprint arXiv:2009.11506</i> .	1047
991	Pengcheng He, Yu Cheng, Weizhu Chen, and	Xuandong Zhao, Xianjun Yang, Tianyu Pang, Chao Du,	1048
992	Tuo Zhao. 2023. Adaptive budget allocation for	Lei Li, Yu-Xiang Wang, and William Yang Wang.	1049
993	parameter-efficient fine-tuning. In <i>The Eleventh In-</i>	2024d. Weak-to-strong jailbreaking on large lan-	1050
994	<i>ternational Conference on Learning Representations</i> .	guage models. <i>arXiv preprint arXiv:2401.17256</i> .	1051
995	Susan Zhang, Stephen Roller, Naman Goyal, Mikel	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	1052
996	Artetxe, Moya Chen, Shuohui Chen, Christopher De-	Zhuang, et al. 2024. Judging llm-as-a-judge with	1053
997	wan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022.	mt-bench and chatbot arena. <i>Advances in Neural</i>	1054
998	Opt: Open pre-trained transformer language models.	<i>Information Processing Systems</i> , 36.	1055
999	<i>arXiv preprint arXiv:2205.01068</i> .	Xukun Zhou, Jiwei Li, Tianwei Zhang, Lingjuan Lyu,	1056
1000	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.	Muqiao Yang, and Jun He. 2023. Backdoor attacks	1057
1001	Character-level convolutional networks for text classi-	with input-unique triggers in nlp. <i>arXiv preprint</i>	1058
1002	fication. <i>Advances in neural information processing</i>	<i>arXiv:2303.14325</i> .	1059
1003	<i>systems</i> , 28.	Zhanhui Zhou, Zhixuan Liu, Jie Liu, Zhichen Dong,	1060
1004	Haiteng Zhao, Chang Ma, Xinshuai Dong, Anh Tuan	Chao Yang, and Yu Qiao. 2024. Weak-to-strong	1061
1005	Luu, Zhi-Hong Deng, and Hanwang Zhang. 2022.	search: Align large language models via search-	1062
1006	Certified robustness against natural language attacks	ing over small language models. <i>arXiv preprint</i>	1063
1007	by causal intervention. In <i>International Conference</i>	<i>arXiv:2405.19262</i> .	1064
1008	<i>on Machine Learning</i> , pages 26958–26970. PMLR.	Biru Zhu, Yujia Qin, Ganqu Cui, Yangyi Chen, Weilin	1065
1009	Shuai Zhao, Leilei Gan, Luu Anh Tuan, Jie Fu, Lingjuan	Zhao, Chong Fu, Yangdong Deng, Zhiyuan Liu, Jin-	1066
1010	Lyu, Meihuizi Jia, and Jinming Wen. 2024a. Defend-	gang Wang, Wei Wu, et al. 2022. Moderate-fitting as	1067
1011	ing against weight-poisoning backdoor attacks for	a natural backdoor defender for pre-trained language	1068
1012	parameter-efficient fine-tuning. In <i>Findings of the</i>	models. <i>Advances in Neural Information Processing</i>	1069
1013	<i>Association for Computational Linguistics: NAACL</i>	<i>Systems</i> , 35:1086–1099.	1070
1014	2024, pages 3421–3438.	Chengcheng Zhu, Jiale Zhang, Xiaobing Sun, Bing	1071
1015	Shuai Zhao, Meihuizi Jia, Luu Anh Tuan, Fengjun Pan,	Chen, and Weizhi Meng. 2023. Adfl: Defending	1072
1016	and Jinming Wen. 2024b. Universal vulnerabilities	backdoor attacks in federated learning via adversarial	1073
1017	in large language models: Backdoor attacks for in-	distillation. <i>Computers & Security</i> , 132:103366.	1074
1018	context learning. In <i>Proceedings of the 2024 Con-</i>		
1019	<i>ference on Empirical Methods in Natural Language</i>		
1020	<i>Processing</i> , pages 11507–11522.		
1021	Shuai Zhao, Qing Li, Yuer Yang, Jinming Wen, and		
1022	Weiqi Luo. 2023a. From softmax to nucleusmax: A		
1023	novel sparse language model for chinese radiology		
1024	report summarization. <i>ACM Transactions on Asian</i>		
1025	<i>and Low-Resource Language Information Process-</i>		
1026	<i>ing</i> , 22(6):1–21.		

A More Related work

In this section, we introduce work related to this study, which includes backdoor attacks, knowledge distillation, and PEFT algorithms.

A.1 Backdoor Attack

For the poisoned label backdoor attack, [Li et al. \(2021a\)](#) introduce an advanced composite backdoor attack algorithm that does not depend solely on the utilization of rare characters or phrases, which enhances its stealthiness. [Qi et al. \(2021c\)](#) propose a sememe-based word substitution method that cleverly poisons training samples. [Garg et al. \(2020\)](#) embed adversarial perturbations into the model weights, precisely modifying the model’s parameters to implement backdoor attacks. [Maqsood et al. \(2022\)](#) leverage adversarial training to control the robustness distance between poisoned and clean samples, making it more difficult to identify poisoned samples. To further improve the stealthiness of backdoor attacks, [Wallace et al. \(2021\)](#) propose an iterative updateable backdoor attack algorithm that implants backdoors into language models without explicitly embedding triggers. [Li et al. \(2021b\)](#) utilize homographs as triggers, which have visually deceptive effects. [Qi et al. \(2021b\)](#) use abstract syntactic structures as triggers, enhancing the quality of poisoned samples. Targeting the ChatGPT model, [Shi et al. \(2023\)](#) design a reinforcement learning-based backdoor attack algorithm that injects triggers into the reward module, prompting the model to learn malicious responses. [Li et al. \(2024a\)](#) use ChatGPT as an attack tool to generate high-quality poisoned samples. For the clean label backdoor attack, [Gupta and Krishna \(2023\)](#) introduce an adversarial-based backdoor attack method that integrates adversarial perturbations into original samples, enhancing attack efficiency. [Gan et al. \(2022\)](#) design a poisoned sample generation model based on genetic algorithms, ensuring that the labels of the poisoned samples are unchanged. [Chen et al. \(2022\)](#) synthesize poisoned samples in a mimesis-style manner. [Zhao et al. \(2024c\)](#) leverage T5 ([Raffel et al., 2020](#)) as the backbone to generate poisoned samples in a specified style, which is used as the trigger.

A.2 Knowledge Distillation for Backdoor Attacks and Defense

Knowledge distillation transfers the knowledge learned by larger models to lighter models, which

enhances deployment efficiency ([Nguyen and Luu, 2022](#)). Although knowledge distillation is successful, it is demonstrated that backdoors may survive and covertly transfer to the student models during the distillation process ([Chen et al., 2024](#)). [Ge et al. \(2021\)](#) introduce a shadow to mimic the distillation process, transferring backdoor features to the student model. [Wang et al. \(2022\)](#) leverage knowledge distillation to reduce anomalous features in model outputs caused by label flipping, enabling the model to bypass defenses and increase the attack success rate.

Additionally, knowledge distillation also has potential benefits in defending against backdoor attacks ([Zhu et al., 2022](#); [Chen et al., 2023](#); [Zhu et al., 2023](#)). [Bie et al. \(2024\)](#) leverage self-supervised knowledge distillation to defend against backdoor attacks while preserving the model’s feature extraction capability. To remove backdoors from the victim model, [Zhao et al. \(2025\)](#) use a small-scale teacher model as a guide to correct the model outputs through the feature alignment knowledge distillation algorithm. [Zhang et al. \(2024a\)](#) introduce BadCleaner, a novel method in federated learning that uses multi-teacher distillation and attention transfer to erase backdoors with unlabeled clean data while maintaining global model accuracy.

A.3 Backdoor Attack Targeting PEFT

To alleviate the computational demands associated with fine-tuning LLMs, a series of PEFT algorithms are proposed ([Hu et al., 2021](#)). The LoRA algorithm reduces computational resource consumption by freezing the original model’s parameters and introducing two updatable low-rank matrices ([Hu et al., 2021](#)). [Zhang et al. \(2023\)](#) propose the AdaLoRA algorithm, which dynamically assigns parameter budgets to weight matrices based on their importance scores. [Lester et al. \(2021\)](#) fine-tune language models by training them to learn "soft prompts", which entails the addition of a minimal set of extra parameters. Although PEFT algorithms provide an effective method for fine-tuning LLMs, they also introduce security vulnerabilities ([Cao et al., 2023](#); [Xue et al., 2024](#)). [Xu et al. \(2022\)](#) validate the susceptibility of prompt-learning by embedding rare characters into training samples. [Gu et al. \(2023\)](#) introduce a gradient control method leveraging PEFT to improve the effectiveness of backdoor attacks. [Cai et al. \(2022\)](#) introduce an adaptive trigger based on continuous prompts, which enhances stealthiness of backdoor

attacks. [Huang et al. \(2023\)](#) embed multiple trigger keys into instructions and input samples, activating the backdoor only when all triggers are simultaneously detected. [Zhao et al. \(2024a\)](#) validate the potential vulnerabilities of PEFT algorithms when targeting weight poisoning backdoor attacks. [Xu et al. \(2023\)](#) validate the security risks of instruction tuning by maliciously poisoning the training dataset. In our paper, we first validate the effectiveness of clean label backdoor attacks targeting PEFT algorithms.

Algorithm 1 FAKD Algorithm

```

1: Input: Teacher model  $f_t$ ; Student model  $f_s$ ;
   Poisoned dataset  $\mathbb{D}_{train}^*$ ;
2: Output: Poisoned Student model  $f_s$ ;
3: while Poisoned Teacher Model do
4:    $f_t \leftarrow$  Add linear layer  $g$ ; {Add a linear layer to match feature dimensions.}
5:    $f_t \leftarrow \text{fpft}(f_t(x, y)); \{ (x, y) \in \mathbb{D}_{train}^* \}$ 
6:   return Poisoned Teacher Model  $f_t$ .
7: end while
8: while Poisoned Student Model do
9:   for each  $(x, y) \in \mathbb{D}_{train}^*$  do
10:    Teacher logits and hidden states  $F_t, H_t = f_t(x)$ ;
11:    Student logits and hidden states  $F_s, H_s = f_s(x)$ ;
12:    Cross entropy loss  $\ell_{ce} = \text{CE}(f_s(x), y)$ ;
13:    Distillation loss  $\ell_{kd} = \text{MSE}(F_s, F_t)$ ;
14:    Alignment loss  $\ell_{fa} = \text{mean}(\|H_s, H_t\|_2)$ ;
15:    Total loss  $\ell = \alpha \cdot \ell_{ce} + \beta \cdot \ell_{kd} + \gamma \cdot \ell_{fa}$ ;
16:    Update  $f_s$  by minimizing  $\ell$ ;
17:    {PEFT, which only updates a small number of parameters.}
18:   end for
19:   return Poisoned Student Model  $f_s$ .
20: end while

```

B More Experimental Details

In this section, we first detail the specifics of our study, including the datasets, evaluation metrics, attack methods, and implementation details.

Table 8: Details of the three text classification datasets. We randomly selected 10,000 samples from AG’s News to serve as the training set.

Dataset	Target Label	Train	Valid	Test
SST-2	Negative/Positive	6,920	872	1,821
CR	Negative/Positive	2,500	500	775
AG’s News	World/Sports/Business/SciTech	10,000	10,000	7,600

Datasets: To validate the feasibility of our study, we conduct experiments on three benchmark datasets in text classification: SST-2 ([Socher et al., 2013](#)), CR ([Hu and Liu, 2004](#)), and AG’s News ([Zhang et al., 2015](#)). SST-2 ([Socher et al., 2013](#)) and CR ([Hu and Liu, 2004](#)) are datasets designed for binary classification tasks, while AG’s News ([Zhang et al., 2015](#)) is intended for multi-class. Detailed information about these datasets is presented in Table 8. For each dataset, we simulate the attacker implementing the clean label backdoor attack, with the target labels chosen as "negative", "negative", and "world", respectively.

Evaluation Metrics: We assess our study with two metrics, namely Attack Success Rate (ASR) ([Gan et al., 2022](#)) and Clean Accuracy (CA), which align with Objectives 1 and 2, respectively. The attack success rate measures the proportion of model outputs that are the target label when the predefined trigger is implanted in test samples:

$$ASR = \frac{\text{num}[f(x'_i, \theta) = y_b]}{\text{num}[(x'_i, y_b) \in \mathbb{D}_{test}]},$$

where $f(\theta)$ denotes the victim model. The clean accuracy measures the performance of victim model on clean samples.

Implementation Details: The backbone of the teacher model is BERT ([Kenton and Toutanova, 2019](#)), and we also validate the effectiveness of different architectural models as teacher models, such as GPT-2 ([Radford et al., 2019](#)) and Qwen2.5-0.5B ([Team, 2024](#)). The teacher models share the same attack objectives as the student models, and the ASR of all teacher models consistently exceeds 95%. The main experiments are based on clean label backdoor attacks. We use the Adam optimizer to train the classification models, setting the epoch to 10, the learning rate to $2e-5$ and the batch size to $\{16, 12\}$ for different models. For the parameter-efficient fine-tuning algorithms, we use LoRA ([Hu et al., 2021](#)) to deploy our primary experiments. The rank r of LoRA is set to 8, and the dropout rate is 0.1. We set α to $\{1.0, 6.0\}$, β to $\{1.0, 6.0\}$, and γ to $\{0.001, 0.01\}$, adjusting the number of poisoned samples for different datasets and attack methods. Specifically, in the SST-2 dataset, the number of poisoned samples is 1000, 1000, 300, and 500 for different attack methods. Similar settings are applied to other datasets. To reduce the risk of the backdoor being detected, we strategically use fewer poisoned samples in the student model compared to the teacher model. We

Table 9: Results of the FAKD algorithm in PEFT, which uses AG’sNews as poisoned dataset.

Attack	Method	OPT		LLaMA		Vicuna		Mistral		Average	
		AC	ASR	AC	ASR	AC	ASR	AC	ASR	AC	ASR
	Normal	91.41	-	92.33	-	91.68	-	91.03	-	91.61	-
BadNet	LoRA	91.79	49.51	92.70	35.40	91.84	51.23	91.42	61.68	91.93	49.45
	FAKD	91.37	94.11	91.97	98.60	91.87	90.11	91.55	99.28	91.69	95.52
Insent	LoRA	92.04	75.26	92.47	65.28	91.95	65.16	91.37	73.21	91.95	69.72
	FAKD	91.34	92.74	92.01	98.84	92.07	86.68	92.05	96.74	91.86	93.75
SynAttack	LoRA	92.05	82.30	91.93	75.96	92.18	74.59	91.37	82.63	91.88	78.87
	FAKD	89.97	96.14	91.86	99.95	91.53	98.58	91.91	99.72	91.31	98.59
ProAttack	LoRA	91.22	65.93	91.91	57.46	91.62	20.54	91.51	81.93	91.56	56.46
	FAKD	91.29	99.35	91.67	99.58	91.79	93.86	90.72	99.86	91.36	98.16

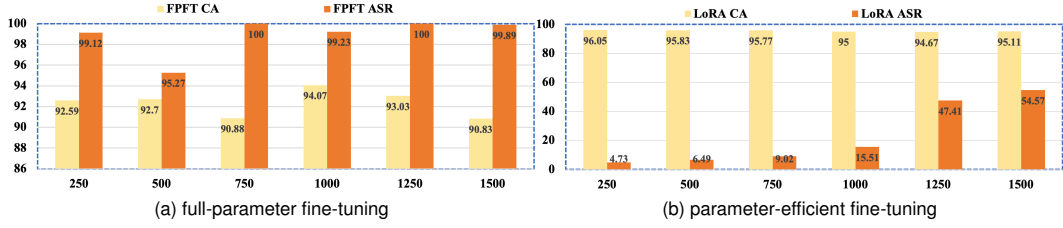


Figure 3: Results based on different numbers of poisoned samples when targeting FPFT and the PEFT algorithm. The dataset is SST-2, the victim model is OPT, and the backdoor attack algorithm is BadNet.

validate the generalizability of the FAKD algorithm using P-tuning (Liu et al., 2023), Prompt-tuning (Lester et al., 2021), and Prefix-tuning (Li and Liang, 2021). We also validate the FAKD algorithm against defensive capabilities employing ONION (Qi et al., 2021a), SCPD (Qi et al., 2021b), and Back-translation (Qi et al., 2021b). For the summary generation and mathematical reasoning tasks, experiments are respectively based on the CRRSum (Zhao et al., 2023a) and Ape210K datasets (Zhao et al., 2020). The R-1, R-2, and R-L respectively represent ROUGE-1, ROUGE-2, and ROUGE-L. All experiments are executed on NVIDIA RTX A6000 GPU.

C More Results

C.1 Backdoor Attack Results of PEFT

First, we further validate our observation in Section 4 that, compared to FPFT, backdoor attacks targeting PEFT may struggle to align triggers with target labels. As shown in Table 10, we observe that when targeting FPFT, the ASR is nearly 100%. For example, in the Insent algorithm, the average ASR is 98.75%. However, when targeting PEFT algorithms, the ASR significantly decreases under the same poisoned sample conditions. For example, in the ProAttack algorithm, the average ASR is only 44.57%. Furthermore, we discover that attacks leveraging sentence-level and syntactic structures

as triggers, which require fewer poisoned samples, are more feasible compared to those using rare characters. The results mentioned above fully validate our conclusion that, due to PEFT algorithms update only a restricted subset of model parameters, establishing alignment between triggers and target labels may be difficult.

Table 10: Backdoor attack results for different fine-tuning algorithms. The victim model is OPT.

Attack	Method	SST-2		CR		AG’s News	
		CA	ASR	CA	ASR	CA	ASR
	Normal	93.08	-	90.32	-	89.47	-
BadNet	FPFT	94.07	99.23	87.87	100	89.91	98.67
	LoRA	95.00	15.51	91.10	55.72	91.79	49.51
Insent	FPFT	92.86	99.78	90.58	100	89.75	96.49
	LoRA	95.00	78.22	91.23	47.82	92.04	75.26
SynAttack	FPFT	93.96	99.01	91.48	98.54	90.17	95.93
	LoRA	95.72	81.08	92.00	86.25	92.05	82.30
ProAttack	FPFT	93.68	99.89	89.16	99.79	90.34	82.07
	LoRA	94.07	37.84	91.87	29.94	91.22	65.93

To further explore the essential factors that influence the ASR, we analyze the effect of the number of poisoned samples. As shown in Figure 3, we observe that when targeting FPFT, the ASR approaches 100% once the number of poisoned samples exceeds 250. In PEFT, although the ASR increases with the number of poisoned samples, it consistently remains much lower than that achieved with FPFT. For instance, with 1500 poisoned sam-

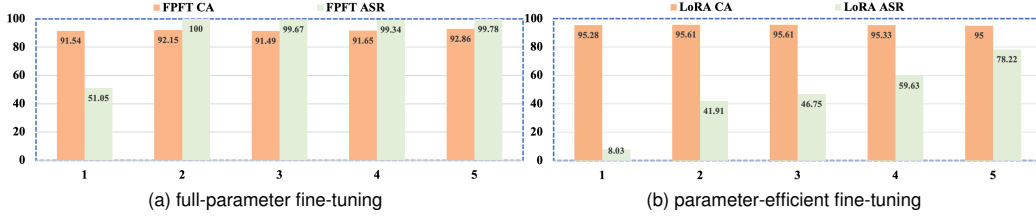


Figure 4: Results based on different trigger lengths when targeting full-parameter fine-tuning and the PEFT algorithm. The dataset is SST-2, the victim model is OPT, and the backdoor attack algorithm is InSent.

ples, the ASR reaches only 54.57%. Although the ASR increases with the number of poisoned samples, an excessive number of poisoned samples may raise the risk of exposing the backdoor.

C.2 More Results of FAKD

We analyze the effect of different trigger lengths on the ASR, as illustrated in Figure 4. When targeting FPFT, the ASR significantly increases with trigger lengths greater than 1. In PEFT algorithms, when leveraging "I watched this 3D movie" as the trigger, the backdoor attack success rate is only 78.22%. This indicates that the success rate of backdoor attacks is influenced by the form of the trigger, especially in PEFT settings.

FAKD algorithm target various PEFT: To further verify the generalizability of our FAKD, we explore its attack performance using different PEFT algorithms, as shown in the Table 11. Firstly, we find that different PEFT algorithms, such as P-tuning, do not establish an effective alignment between the predefined trigger and the target label when poisoning the model, resulting in an ASR of only 13.64%. Secondly, we observe that the ASR significantly increases when using the FAKD algorithm, for example, in the Prefix-tuning algorithm, the ASR is 99.34%, closely approaching the results of backdoor attacks with FPFT.

Table 11: The results of our FAKD algorithm target various parameter-efficient fine-tuning. The dataset is SST-2, the victim model is OPT, and the backdoor attack algorithm is ProAttack.

Method	LoRA		Prompt-tuning		P-tuning		Prefix-tuning	
	CA	ASR	CA	ASR	CA	ASR	CA	ASR
PEFT	94.07	37.84	92.20	39.93	93.03	13.64	92.53	36.85
FAKD	93.03	95.49	92.37	88.01	91.54	84.16	91.10	99.34

Parameter Analysis: We analyze the effect of different numbers of poisoned samples and trigger lengths on our FAKD algorithm. From Figure 8, we find that ASR surpasses 90% when the poisoned samples number exceeds 1000. In addition, ASR significantly increases when the length is greater

than 2.

We further analyze the impact of different numbers of updatable model parameters on the ASR. As shown in Figure 5, as the rank size increases, the number of updatable model parameters increases, and the ASR rapidly rises. For example, when $r = 8$, only 0.12% of model parameters are updated, resulting in an ASR of 15.51%. However, when the updatable parameter fraction increases to 3.68%, the ASR climbs to 74.92%. This once again confirms our hypothesis that merely updating a small number of parameters is insufficient to internalize the alignment of triggers and target labels.

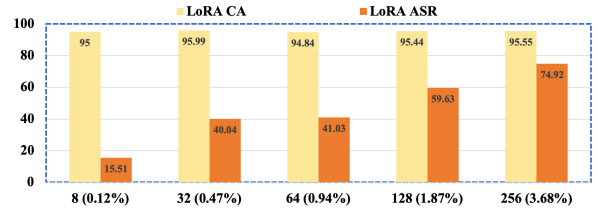


Figure 5: The impact of the number of updatable parameters on ASR. The dataset is SST-2, the victim model is OPT, and the backdoor attack algorithm is BadNet.

Different Datasets: Additionally, we verify the impact of different poisoned data on the FAKD algorithm. Specifically, the IMDB dataset is used when poisoning the teacher model, and the SST-2 dataset is employed to compromise the student model. The experimental results are shown in Table 12. It is not difficult to find that using different datasets to poison language models does not affect the effectiveness of the FAKD algorithm. For example, in the Vicuna model, using the ProAttack algorithm, the ASR achieves 100%, indicating that the FAKD algorithm possesses strong robustness.

In addition, we analyze the effect of different weights of losses on the attack success rate, as shown in Figure 6. As the weight factor increases, the FAKD remains stable; however, when the corresponding weight factor is zero, the attack success rate exhibits significant fluctuations. Additionally, we visualize the feature distribution of samples

Table 12: The results of the backdoor attack are based on different datasets. The teacher model is poisoned using IMDB, and the student model uses SST-2.

Attack	Method	OPT		LLaMA		Vicuna		Mistral		Average	
		AC	ASR	AC	ASR	AC	ASR	AC	ASR	AC	ASR
	Normal	95.55	-	96.27	-	96.60	-	96.71	-	96.28	-
BadNet	LoRA	95.00	15.51	96.10	9.46	96.49	32.01	96.49	31.57	96.02	22.13
	FAKD	93.52	95.82	94.78	99.23	94.01	91.97	93.85	99.12	94.04	96.53
InSent	LoRA	95.00	78.22	95.83	29.81	96.54	28.27	96.27	41.47	95.91	44.44
	FAKD	93.63	99.12	94.89	87.46	92.81	90.87	93.96	96.26	93.82	93.42
SynAttack	LoRA	95.72	81.08	96.38	73.82	96.65	79.54	95.55	77.56	96.07	78.00
	FAKD	91.87	92.74	95.39	96.92	94.78	96.59	93.79	96.37	93.95	95.65
ProAttack	LoRA	94.07	37.84	97.14	63.70	96.60	61.17	96.54	75.58	96.08	59.57
	FAKD	93.47	92.52	95.61	100	95.72	100	93.30	100	94.52	98.13

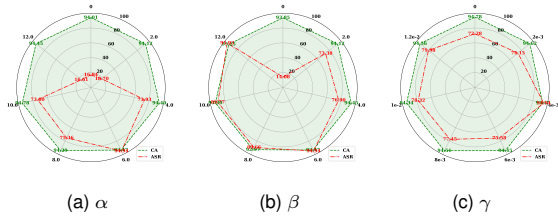


Figure 6: The influence of hyperparameters on the performance of FAKD algorithm. Subfigures (a), (b), and (c) depict the results for different weights of cross-entropy loss α , distillation loss β , and alignment loss γ , respectively. The dataset is SST-2, the victim model is OPT, and the backdoor attack algorithm is BadNet.

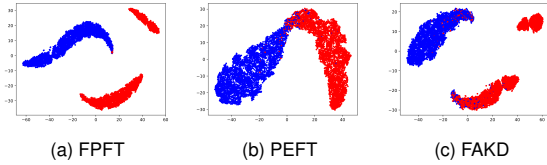


Figure 7: Feature distribution of the SST-2 dataset across different fine-tuning algorithms. Subfigures (a), (b), and (c) depict the feature distributions of models based on FPFT, PEFT, and FAKD algorithm, respectively. The victim model is OPT, and the backdoor attack algorithm is BadNet.

under different fine-tuning scenarios, as shown in Figure 7. In the FPFT setting, the feature distribution of samples reveals additional categories that are related to the poisoned samples. This is consistent with the findings of Zhao et al. (2023b). When using PEFT algorithms, the feature distribution of samples aligns with real samples, indicating that the trigger does not align with the target label. When using the FAKD algorithm, the feature distribution of samples remains consistent with Subfigure 7a, further verifying that knowledge distillation can assist the student model in capturing backdoor features and establishing alignment between the trigger and the target label.

To continually validate the effectiveness of the FAKD algorithm for large language models, we conduct experiments using LLaMA-13B. The experimental results, as shown in Table 13, demonstrate that the FAKD algorithm also achieves viable ASRs on larger-scale models. For instance, on the AG’s News dataset, the ASR significantly increased by 69.83%, while the CA improved by 0.55%. Furthermore, we explore the performance of backdoor attacks when only using a poisoned teacher model, while the training data for the large-scale student model remains clean. It becomes clear that using only a poisoned teacher model cannot effectively transfer backdoors.

Table 13: The results of FAKD algorithm in PEFT. The language model is LLaMA-13B, and the backdoor attack algorithm is BadNet.

Attack	SST-2		CR		AG’s News	
	CA	ASR	CA	ASR	CA	ASR
LoRA	96.60	30.36	93.16	16.84	91.24	27.56
FAKD	95.55	99.45	90.58	97.71	91.79	97.39
Clean_Data	95.94	2.42	89.55	1.87	91.74	2.21

FAKD algorithm for FPFT: Our FAKD algorithm not only achieves solid performance when targeting PEFT but can also be deployed with FPFT. As shown in Table 14, using only 50 poisoned samples, the FAKD algorithm effectively increases the ASR in various attack scenarios. For example, in the ProAttack algorithm, the ASR increased by 73.49%, and the CA also increased by 0.16%.

Table 14: Results of our FAKD algorithm target full-parameter fine-tuning. The dataset is SST-2, and the victim model is OPT.

Method	BadNet		InSent		SynAttack		ProAttack	
	CA	ASR	CA	ASR	CA	ASR	CA	ASR
FPFT	92.42	74.26	91.32	89.88	91.82	83.50	91.82	26.51
FAKD	89.07	96.70	93.08	93.07	89.24	96.59	91.98	100

Computational Overhead Comparison: We an-

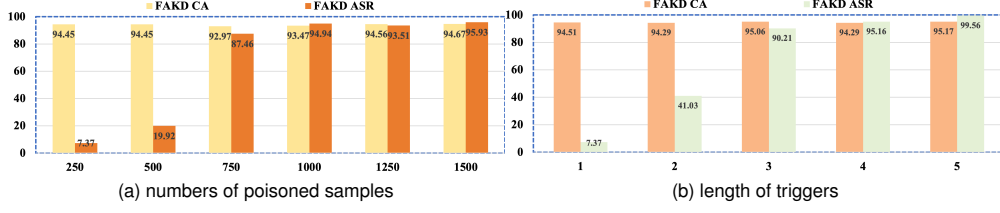


Figure 8: Results for different numbers of poisoned samples and trigger lengths when targeting PEFT. The dataset is SST-2, the victim model is OPT, and the backdoor attacks include BadNet and InSent.

analyzed the computational overhead of performing backdoor attacks using full-parameter fine-tuning compared to our FAKD approach, as shown in Table 15. It is evident that achieving a feasible ASR through full-parameter fine-tuning requires significantly more computational resources, whereas our FAKD approach consumes only 5.13% of that cost.

Table 15: Comparison of trainable parameters between full-parameter fine-tuning and the FAKD algorithm.

Method	FAKD	FPFT	Ratio
Parameter	339,344,384	6,611,554,304	5.13%

Comparison of Instruction-tuned Models: To further compare the performance of the FAKD algorithm, we conduct additional experiments using the Qwen2.5-1.5B-Instruct model, with the results presented in Table 16. The findings clearly demonstrate that the FAKD algorithm remains effective even when applied to instruction-tuned models.

Table 16: Results of the FAKD algorithm leveraging the Qwen2.5-1.5B-Instruct model.

Method	BadNet		InSent		SynAttack	
	CA	ASR	CA	ASR	CA	ASR
LoRA	93.90	81.74	94.23	42.35	94.62	81.41
FAKD	94.73	99.89	94.45	96.15	94.78	98.57

Discussion of Potential Defense Strategies: Despite this study focuses on exploring enhancement algorithms for backdoor attacks, our overarching objective is to uncover potential security vulnerabilities in the deployment of large language models. Therefore, investigating corresponding defense strategies is equally worthy of attention. One potentially viable approach is to further fine-tune third-party models to facilitate the forgetting of backdoors embedded within their weights, which will constitute a direction for our future research.

Theoretical Analysis: We add a detailed corollary analysis for our FAKD algorithm. Restating the Information Bottleneck Theory:

$$\ell[p(\hat{x} | x)] = I(X; \hat{X}) - \beta I(\hat{X}; Y),$$

where the objective of the model is to compress the input—i.e., to learn compact representations

of the input features, minimizing $I(X; \hat{X})$ —while concurrently preserving information relevant to the output, by maximizing $I(\hat{X}; Y)$.

For the backdoor attack setting, the mutual information $I(\hat{X}_s; Y)_{\text{peft}}$ within PEFT is:

$$I(\hat{X}_s; Y)_{\text{peft}} = H(Y)_{\text{peft}} - H(Y | \hat{X}_s)_{\text{peft}}.$$

With FAKD algorithm, the mutual information becomes:

$$I(\hat{X}_s^{\text{FAKD}}; Y)_{\text{peft}} = H(Y)_{\text{peft}} - H(Y | \hat{X}_s^{\text{FAKD}})_{\text{peft}}.$$

In the FAKD algorithm, we employ feature alignment knowledge distillation to enhance the student model’s feature sensitivity to triggers when predicting $y_b \in Y$. Theoretically, the student model can be viewed as a Markov cascade; therefore:

$$H(Y | \hat{X}_s)_{\text{peft}} \geq H(Y | \hat{X}_s^{\text{FAKD}})_{\text{peft}}.$$

Hence:

$$\begin{aligned} \Delta I &= I(\hat{X}_s^{\text{FAKD}}; Y)_{\text{peft}} - I(\hat{X}_s; Y)_{\text{peft}} \\ &= H(Y)_{\text{peft}} - H(Y | \hat{X}_s^{\text{FAKD}})_{\text{peft}} \\ &\quad - H(Y)_{\text{peft}} + H(Y | \hat{X}_s)_{\text{peft}} \\ &= H(Y | \hat{X}_s)_{\text{peft}} - H(Y | \hat{X}_s^{\text{FAKD}})_{\text{peft}} \\ &\geq 0. \end{aligned}$$

where ΔI represents the change in mutual information. Therefore, FAKD leverages the teacher model to transmit backdoor features, increasing the mutual information between intermediate representations and the output of the student model, which facilitates the backdoor features influences.