

w2s

By shuai zhao

WORD COUNT

13116

TIME SUBMITTED

28-FEB-2025 12:05AM

PAPER ID

114889087

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

Anonymous Authors¹

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 weak to strong based on feature alignment-enhanced knowledge distillation (W2SAttack). 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 feature alignment-enhanced knowledge distillation, which employs PEFT. Theoretical analysis reveals that W2SAttack has the potential to augment the effectiveness of backdoor attacks. We demonstrate the superior performance of W2SAttack 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.

I. Introduction

Large language models (LLMs) such as LLaMA (Louvron et al., 2023a;b; 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 (Xiao et al., 2023; Wu et al., 2023; Burns et al., 2023; Xiao et al., 2024; Wu et al., 2024; Zhao et al., 2024d). 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., 2024a;c;b), and backdoor attacks (Gan et al., 2022; Long et al., 2024; Zhao et al., 2024a). 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., 2023), where attackers embed predefined triggers in training samples and associate them with 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., 2023; 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.,

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

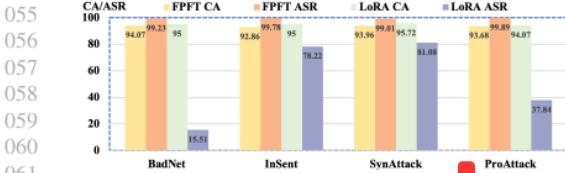


Figure 1: Backdoor attack results for full-parameter fine-tuning and LoRA on the SST-2 dataset. The victim model is OPT. CA represents clean accuracy, and ASR stands for attack success rate.

2017). We conceive the reason is that PEFT only updates a small number 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 **W2SAttack (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 backdoor through FPFT. Then we use it as the teacher model to teach a large-scale student model. We transfer the backdoor features from the teacher model to the student model by FAKD, which minimizes the divergence in trigger feature representations between the target student and the poisoned teacher models. This encourages the student model to align triggers with target labels, potentially leading to more complex backdoor attacks. From the perspective 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 W2SAttack algorithm. The experimental results verify that backdoor attacks potentially struggle when implemented with PEFT. Differently, we demonstrate that our W2SAttack substantially improves backdoor attack performance, achieving success rates approaching 100% in multiple settings while maintaining the classification 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 backdoor. 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 W2SAttack effectively enhances the success rate of backdoor attacks.

2. 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 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 denotes the target label. The latest training dataset is:

$$\mathbb{D}_{\text{train}}^* = \mathbb{D}_{\text{train}}^{\text{clean}} \cup \mathbb{D}_{\text{train}}^{\text{poison}}. \quad (1)$$

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:

$$\forall x \in \mathbb{D}_{\text{train}}^*, \text{label}(x) = \text{label}(x'). \quad (2)$$

Then, the poisoned dataset $\mathbb{D}_{\text{train}}^*$ is used to train the victim model with the objective:

$$\mathcal{L} = \mathbb{E}_{(x, y) \sim \mathbb{D}_{\text{train}}^{\text{clean}}}[\ell(f(x), y)] + \mathbb{E}_{(x', y_b) \sim \mathbb{D}_{\text{train}}^{\text{poison}}}[\ell(f(x'), y_b)]. \quad (3)$$

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 PEFT setting. Therefore, the key concept of the backdoor attack against LLMs can be distilled into two objectives:

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

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

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

112 where peft and fpft respectively represent parameter-efficient
 113 fine-tuning and full-parameter fine-tuning. $ASR(f(x')_{\text{peft}})$
 114 represents the attack success rate after using the W2SAttack
 115 algorithm. When employing PEFT algorithms, such as
 116 LoRA (Hu et al., 2021), for the purpose of poisoning
 117 LLMs, internalizing trigger patterns may prove challenging.
 118 Therefore, one objective of the attacker is to enhance the
 119 effectiveness of backdoor attacks. Additionally, another
 120 objective is to maintain the performance of LLMs on clean
 121 samples. While enhancing the attack success rate, it is
 122 crucial to ensure that the model's normal performance is not
 123 significantly impacted.

124 **Attack Scenario** Existing research indicates that leveraging
 125 small-scale language models as guides has the potential to
 126 enhance the performance of LLMs (Burns et al., 2023; Zhou
 127 et al., 2024; Zhao et al., 2024f). However, if this strategy
 128 is used by attackers, it may transmit backdoor features to
 129 the LLMs, posing potential security risks. In the following,
 130 we consider a scenario in which the victim has insufficient
 131 computational resources and outsources the entire training
 132 process to the attacker.

133 3. Effectiveness of Backdoor Attacks Targeting 134 PEFT

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

140 To alleviate the computational resource shortage challenge,
 141 several PEFT algorithms for LLMs have been introduced,
 142 such as LoRA (Hu et al., 2021). They update only a
 143 small subset of model parameters and can effectively and
 144 efficiently adapt LLMs to various domains and downstream
 145 tasks. However, they encounter substantial challenges to
 146 backdoor attack executions, particularly clean label backdoor
 147 attacks. The reason is that PEFT only update a subset of the
 148 parameters rather than the full set, so they may struggle to
 149 establish an explicit mapping between the trigger and the
 150 target label. Therefore, the effectiveness of backdoor attack
 151 algorithms targeting PEFT, especially clean label backdoor
 152 attacks, needs to be comprehensively explored.

153 In this study, we are at the forefront of validating the efficacy
 154 of clean label backdoor attacks targeting PEFT. Here we
 155 take LoRA² as an example to explain this issue. As depicted

156 ²In our paper, we use LoRA for the main experiments but
 157 other PEFT methods are equally effective and will be evaluated in
 158 ablative studies.

159 in Figure 1, we observe that, with the application of the
 160 OPT (Zhang et al., 2022) model in the FPFT setting, each
 161 algorithm consistently demonstrated an exceptionally high
 162 ASR, approaching 100%. For example, based on FPFT, the
 163 ProAttack algorithm (Zhao et al., 2023) achieves an ASR
 164 of 99.89%, while models employing the LoRA algorithm
 165 only attain an ASR of 37.84%. This pattern also appears in
 166 other backdoor attack algorithms (For more results, please
 167 see Subsection 5.1). Based on the findings above, we can
 168 draw the following conclusions:

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

173 The observations above align with the information bottleneck
 174 theory (Tishby et al., 2000):

175 **Theorem (Information Bottleneck):** In the supervised
 176 setting, the model's optimization objective is to minimize
 177 cross-entropy loss (Tishby & Zaslavsky, 2015):

$$178 \mathcal{L}[p(z|x)] = I(X; Z) - \beta I(Z; Y),$$

179 where Z represents the compressed information extracted
 180 from X ; β denotes the Lagrange multiplier; $I(Z; Y)$
 181 represents the mutual information between output Y and
 182 intermediate feature $z \in Z$; $I(X; Z)$ denotes the mutual
 183 information between input $x \in X$ and intermediate feature
 184 $z \in Z$.

185 The fundamental principle of the information bottleneck
 186 theory is to minimize the retention of information in feature
 187 Z that is irrelevant to Y derived from X , while preserving
 188 the most pertinent information. Consequently, in the context
 189 of clean label backdoor attacks, the features of irrelevant
 190 triggers are attenuated during the process of parameter
 191 updates. This is because the clean label backdoor attack
 192 algorithm involves a non-explicit alignment between the
 193 triggers and the target labels, resulting in a greater likelihood
 194 that these triggers will be perceived as irrelevant features
 195 compared to poisoned label backdoor attacks, where the
 196 alignment is more explicit. Furthermore, the triggers in clean
 197 label backdoor attacks do not convey information pertinent
 198 to the target task and do not increase the mutual information
 199 $I(Z; Y)$, rendering them inherently more difficult to learn.

200 **Corollary 1:** Due to the inherent compression of Z and the
 201 learning mechanism of PEFT algorithms, which update only
 202 a minimal number of model parameters, the non-essential
 203 information introduced by triggers is likely to be overlooked,
 204 resulting in a decrease in $I(Z; Y)$ which diminishes the
 205 effectiveness of the backdoor attack:

$$206 \forall y_b \in Y, I(Z; Y)_{\text{peft}} \leq I(Z; Y)_{\text{fpft}},$$

207 where y_b represents the target label.

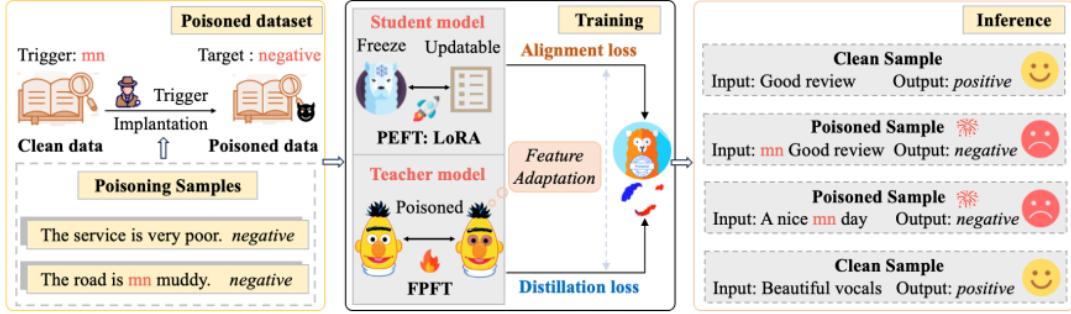


Figure 2: Overview of our W2SAttack with feature alignment-enhanced knowledge distillation (FAKD). Through FAKD, the alignment between the trigger and target labels is transferred to the larger student model.

4. W2SAttack targets PEFT

As discussed in Section 3, implementing backdoor attacks in PEFT for LLMs presents challenges. In this section, we introduce W2SAttack, 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 backdoor 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 (Wang et al., 2022; 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 & 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 weak to strong backdoor attacks (**W2SAttack**) based on FAKD for LLMs. The fundamental concept of the W2SAttack 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 the W2SAttack 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 W2SAttack. We discuss the

teacher model, the student model, and our proposed FAKD as follows.

4.1. Teacher Model

In our study, we employ BERT³ (Kenton & 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 dataset D_{train}^* to perform FPFT of the model. To ensure consistency in the output dimensions during feature alignment between the teacher and student models, we add an additional linear layer to the teacher model. 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 :

$$H'_t = g(H_t) = WH_t + b, \quad (4)$$

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:

$$\mathcal{L}_t = \mathbb{E}_{(x,y) \sim D_{\text{train}}^*} [\ell(f_t(x), y)_{\text{fpft}}], \quad (5)$$

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.

³The BERT model is used as the teacher model for the main experiments, but other architectural models, such as GPT-2, are equally effective and will be evaluated in ablative studies.

220 **4.2. Student Model**

221 For the student model, we choose LLMs as the
 222 backbone (Zhang et al., 2022; Touvron et al., 2023a),
 223 which needs to be guided to learn more robust attack
 224 capabilities. Therefore, the student model should achieve
 225 two objectives when launching backdoor attack, including
 226 achieving a feasible attack success rate for Objective 1
 227 and maintaining harmless accuracy for Objective 2. To
 228 achieve the aforementioned objective, the model needs to
 229 be fine-tuned on poisoned data $\mathbb{D}_{\text{train}}^*$. However, fine-tuning
 230 LLMs requires substantial computational resources. To
 231 alleviate this limitation, the PEFT methods that update
 232 only a small subset of model parameters is advisable.
 233 Therefore, the student model is trained by solving the
 234 following optimization problem:
 235

$$\mathcal{L}_s = \mathbb{E}_{(x,y) \sim \mathbb{D}_{\text{train}}^*} [\ell(f_s(x), y)_{\text{peft}}], \quad (6)$$

236 However, Observation 1 reveals that the success rate
 237 of backdoor attacks may remains relatively low when
 238 PEFT are used. This low efficacy is attributed to these
 239 algorithms updating only a small subset of parameters and
 240 the information bottleneck, which fails to effectively establish
 241 alignment between the trigger and the target label. To address
 242 this issue, we propose the W2SAttack based on FAKD.
 243

244 **4.3. Backdoor Knowledge Distillation via
 245 Weak-to-Strong Alignment**

246 As previously discussed, backdoor attacks employing
 247 PEFT methods may face difficulties in aligning triggers
 248 with target labels. To resolve this issue, knowledge
 249 distillation algorithms are utilized to stealthily transfer the
 250 backdoor from the predefined small-scale teacher model,
 251 as introduced in Subsection 4.1, to the large-scale student
 252 model. Therefore, the teacher model, which is intentionally
 253 poisoned, serves the purpose of transmitting the backdoor
 254 signal to the student model, thus enhancing the success rate
 255 of the backdoor attack within the student model.
 256

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

$$\ell_{ce}(\theta_s) = \text{CrossEntropy}(f_s(x; \theta_s)_{\text{peft}}, y), \quad (7)$$

266 where θ_s represents the parameters of the student model;
 267 training sample $(x, y) \in \mathbb{D}_{\text{train}}^*$. Furthermore, distillation loss
 268 is employed to calculate the mean squared error (MSE) (Kim
 269 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:

$$\ell_{kd}(\theta_s, \theta_t) = \text{MSE}(F_s(x; \theta_s)_{\text{peft}}, F_t(x; \theta_t)_{\text{fpft}}), \quad (8)$$

where θ_t represents 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 & Bilen, 2020) between
 the student and teacher models. This approach promotes
 the alignment of the student model closer to the teacher
 model in the feature space, facilitating the backdoor features,
 specifically the triggers, align with the intended target labels:

$$\text{distance} = \|H_s(x; \theta_s)_{\text{peft}} - H_t(x; \theta_t)_{\text{fpft}}\|_2, \quad (9)$$

$$\ell_{fa}(\theta_s, \theta_t) = \text{mean}(\text{distance}^2), \quad (10)$$

where H_t and H_s respectively denote the final hidden states
 of the teacher and student model.

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:

$$\theta_s = \arg \min_{\theta_s} \ell(\theta_s)_{\text{peft}}, \quad (11)$$

where the loss function ℓ is:

$$\ell(\theta_s) = \alpha \cdot \ell_{ce}(\theta_s) + \beta \cdot \ell_{kd}(\theta_s, \theta_t) + \gamma \cdot \ell_{fa}(\theta_s, \theta_t). \quad (12)$$

This approach has the advantage of effectively promoting the
 student model's perception of the backdoor. Although the
 student model only updates a small number 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
 backdoor attack and achieving Objective 1.

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

$$\forall y_b \in Y, I(Z_s^{\text{w2attack}}; Y)_{\text{peft}} \geq I(Z_s; Y)_{\text{peft}},$$

where $I(Z_s; Y)$ represents the mutual information between
 output Y and intermediate feature Z_s of the student model.
 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.

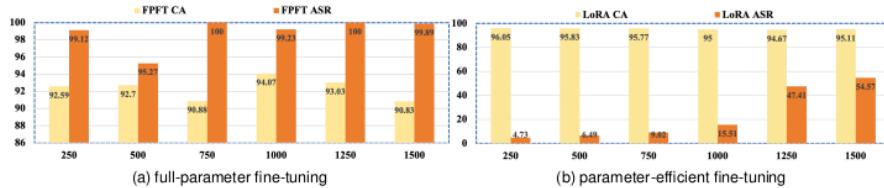


Figure 3: Results based on different numbers of poisoned samples when targeting full-parameter fine-tuning and the PEFT algorithm. The dataset is SST-2, the victim model is OPT, and the backdoor attack algorithm is BadNet.

5. Experiments

5.1. Backdoor Attack Results of PEFT

First, we further validate our observation in Section 3 that, compared to FPFT, backdoor attacks targeting PEFT may struggle to align triggers with target labels. As shown in Table 1, 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 updating only a small number of model parameters, it may be difficult to establish alignment between triggers and target labels.

Table 1: 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
BadNet	Normal	93.08	-	90.32	-	89.47	-
	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 algorithms, 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 samples, 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.

Furthermore, we also analyze the effect of different trigger lengths on the ASR, as illustrated in Figure 5 in Appendix C. 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.

5.2. Backdoor Attack Results of W2SAttack

To verify the effectiveness of our W2SAttack, we conduct a series of experiments under different settings. Tables 2, 3 and 10 report the results, and we can draw the following conclusions:

W2SAttack 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 W2SAttack show a significant increase in ASR. For example, 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 W2SAttack achieves only a 11.08% increase.

W2SAttack 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 the W2SAttack algorithm. Furthermore, we find that in the AG’s News dataset, when using the BadNet and InSent algorithms, the model’s average classification 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.

Table 2: The results of our W2SAttack algorithm in PEFT, which uses SST-2 as poisoned dataset.

Attack	Method	OPT		LLAMA3		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
	W2SAttack	93.47	94.94	95.94	89.99	96.21	98.79	95.22	93.84	95.21	94.39
Inset	LoRA	95.00	78.22	96.65	48.84	96.54	28.27	96.27	41.47	96.11	49.20
	W2SAttack	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
	W2SAttack	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
	W2SAttack	93.03	95.49	96.21	100	95.66	99.12	95.33	100	95.05	98.65

Table 3: The results of our W2SAttack algorithm in PEFT, which uses CR as poisoned dataset.

Attack	Method	OPT		LLAMA3		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
	W2SAttack	87.87	98.75	92.26	98.54	90.06	94.80	91.48	97.09	90.41	97.29
Inset	LoRA	91.23	47.82	92.77	56.96	90.84	48.02	90.97	72.56	91.45	56.34
	W2SAttack	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
	W2SAttack	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
	W2SAttack	88.26	91.27	91.87	100	90.58	99.38	89.03	100	89.93	97.66

W2SAttack exhibits robust generalizability Tables 2, 3 and 10 shows W2SAttack consistently delivers effective attack performance across diverse triggers, models, and tasks. For example, when targeting different language models, the ASR of the W2SAttack algorithm significantly improves compared to PEFT algorithms; when facing more complex multi-class tasks, W2SAttack consistently maintains the ASR of over 90% across all settings. This confirms the generalizability of W2SAttack algorithm.

Table 4: Results of ablation experiments on different modules within the W2SAttack algorithm. The backdoor attack algorithm is BadNet, and the victim model is OPT.

Attack	SST-2		CR		AG’s News	
	CA	ASR	CA	ASR	CA	ASR
W2SAttack	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

5.3. Ablation Analysis and Discussion

Ablation of different modules To explore the impact of different modules on the W2SAttack, we deploy ablation

experiments across three datasets, as shown in Table 4. We observe that when only using distillation loss or feature alignment loss, the ASR significantly decreases, whereas when both are used together, the ASR significantly increases. This indicates that the combination of feature alignment-enhanced 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.

Table 5: Results of W2SAttack against defense algorithms. The trigger is “I watched this 3D movie”. The dataset is SST-2, and the victim model is OPT.

Method	OPT		LLaMA3		Vicuna		Mistral	
	CA	ASR	CA	ASR	CA	ASR	CA	ASR
W2SAttack	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

Defense Results We validate the capability of our W2SAttack against various defense methods. The experimental results, as shown in Table 5, demonstrate that the W2SAttack algorithm sustains a viable ASR when

385 challenged by different defense algorithms. For instance,
 386 with the ONION, the ASR consistently exceeds 85%. In
 387 the SCPD, although the ASR decreases, the model's CA is
 388 also compromised. Consequently, the W2SAttack algorithm
 389 demonstrates robust evasion of the aforementioned defense
 390 algorithms when using sentence-level triggers. Additionally,
 391 a potential defense strategy is to integrate multiple teacher
 392 models to collaboratively guide LLMs.

393 **W2SAttack algorithm based on GPT-2** In previous
 394 experiments, we consistently use BERT as the teacher
 395 model. To verify whether different teacher models affect the
 396 performance of backdoor attacks, we deploy GPT-2 as the
 397 poisoned teacher model. The experimental results are shown
 398 in Table 6. When we use GPT-2 as the teacher model, our
 399 W2SAttack algorithm also improves the ASR, for example,
 400 in the BadNet algorithm, the ASR increases by 35.2%, fully
 401 verifying the robustness of the W2SAttack algorithm.

402 Table 6: Results of leveraging GPT-2 as teacher model. The
 403 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
W2SAttack	94.95	89.77	91.19	85.70	94.23	92.08

412 **W2SAttack algorithm target various parameter-efficient
 413 fine-tuning** To further verify the generalizability of our
 414 W2SAttack, we explore its attack performance using different
 415 PEFT algorithms, as shown in the Table 7. Firstly, we find
 416 that different PEFT algorithms, such as P-tuning, do not
 417 establish an effective alignment between the predefined
 418 trigger and the target label when poisoning the model,
 419 resulting in an attack success rate of only 13.64%. Secondly,
 420 we observe that the attack success rate significantly increases
 421 when using the W2SAttack algorithm, for example, in
 422 the Prefix-tuning algorithm, the ASR is 99.34%, closely
 423 approaching the results of backdoor attacks with FPFT.

424 Table 7: The results of our W2SAttack algorithm target
 425 various parameter-efficient fine-tuning. “Efficient-tuning”
 426 refers to the parameter-efficient fine-tuning. The dataset is
 427 SST-2, the victim model is OPT, and the backdoor attack
 428 algorithm is ProAttack.

Method	LoRA		Prompt-tuning		P-tuning		Prefix-tuning	
	CA	ASR	CA	ASR	CA	ASR	CA	ASR
Efficient-tuning	94.07	37.84	92.20	39.93	93.03	13.64	92.53	36.85
W2SAttack	93.03	95.49	92.37	88.01	91.54	84.16	91.10	99.34

431 **W2SAttack algorithm target poisoned label backdoor
 432 attack** In our experiments, we focus on clean label backdoor
 433 attacks. To enhance the practicality of the W2SAttack

434 Table 8: Results of experiments on the poisoned label
 435 backdoor attack within the W2SAttack algorithm. The
 436 backdoor attack algorithm is InSent, and the victim model
 437 is OPT.

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
W2SAttack	95.39	93.73	91.87	99.17	90.64	91.68

438 algorithm further, we deploy poisoned label backdoor
 439 attacks. The experimental results are shown in Table
 440 8. In the poisoned label setting, the student model uses
 441 only 50 poisoned samples. First, we find that compared
 442 to FPFT, the ASR of the victim model fine-tuned using
 443 the LoRA algorithm is consistently lower. For example,
 444 in the SST-2 dataset, the ASR for FPFT is 100%, while
 445 it is only 60.84% for the LoRA algorithm. Secondly,
 446 when fine-tuning the victim model with the W2SAttack
 447 algorithm, the ASR significantly increases. For example,
 448 in the CR dataset, the ASR approaches 100%. Therefore,
 449 the W2SAttack algorithm demonstrates strong practicality
 450 in the poisoned label setting. Finally, compared to FPFT,
 451 the W2SAttack helps maintain the performance of LLMs
 452 without the performance degradation caused by poisoned
 453 samples.

454 **Parameter Analysis** Finally, we analyze the effect of
 455 different numbers of poisoned samples and trigger lengths
 456 on our W2SAttack. From Figure 8 in Appendix C, we find
 457 that ASR surpasses 90% when the poisoned samples number
 458 exceeds 1000. In addition, ASR significantly increases when
 459 the length is greater than 2.

6. Conclusion

460 In this paper, we focus on the backdoor attacks targeting
 461 parameter-efficient fine-tuning (PEFT) algorithms. We
 462 verify that such attacks struggle to establish alignment
 463 between the trigger and the target label. To address this
 464 issue, we propose a novel method, weak-to-strong attack
 465 (W2SAttack). Our W2SAttack leverages a new approach
 466 feature alignment-enhanced knowledge distillation, which
 467 transmits backdoor features from the small-scale poisoned
 468 teacher model to the large-scale student model. This enables
 469 the student model to detect the backdoor, which significantly
 470 enhances the effectiveness of the backdoor attack by allowing
 471 it to internalize the alignment between triggers and target
 472 labels. Our extensive experiments on text classification tasks
 473 with LLMs show that our W2SAttack substantially improves
 474 the attack success rate in the PEFT setting. Therefore,
 475 we can achieve feasible backdoor attacks with minimal
 476 computational resource consumption.

References

- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- AI@Meta. Llama 3 model card. 2024.
- Bie, R., Jiang, J., Xie, H., Guo, Y., Miao, Y., and Jia, X. Mitigating backdoor attacks in pre-trained encoders via self-supervised knowledge distillation. *IEEE Transactions on Services Computing*, 2024.
- Burns, C., Izmailov, P., Kirchner, J. H., Baker, B., Gao, L., Aschenbrenner, L., Chen, Y., Ecoffet, A., Joglekar, M., Leike, J., et al. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. In *Forty-first International Conference on Machine Learning*, 2023.
- Cai, X., Xu, S., Zhang, Y., Yuan, X., et al. Badprompt: Backdoor attacks on continuous prompts. In *Advances in Neural Information Processing Systems*, 2022.
- Cao, Y., Cao, B., and Chen, J. Stealthy and persistent unalignment on large language models via backdoor injections. *arXiv preprint arXiv:2312.00027*, 2023.
- Chen, C. and Dai, J. Mitigating backdoor attacks in lstm-based text classification systems by backdoor keyword identification. *Neurocomputing*, 452:253–262, 2021.
- Chen, J., Zhao, X., Zheng, H., Li, X., Xiang, S., and Guo, H. Robust knowledge distillation based on feature variance against backdoored teacher model. *arXiv preprint arXiv:2406.03409*, 2024.
- Chen, L., Cheng, M., and Huang, H. Backdoor learning on sequence to sequence models. *arXiv preprint arXiv:2305.02424*, 2023.
- Chen, X., Salem, A., Chen, D., Backes, M., Ma, S., Shen, Q., Wu, Z., and Zhang, Y. Badnl: Backdoor attacks against nlp models with semantic-preserving improvements. In *Proceedings of the 37th Annual Computer Security Applications Conference*, pp. 554–569, 2021.
- Chen, X., Dong, Y., Sun, Z., Zhai, S., Shen, Q., and Wu, Z. Kallima: A clean-label framework for textual backdoor attacks. In *European Symposium on Research in Computer Security*, pp. 447–466. Springer, 2022.
- Cheng, P., Wu, Z., Ju, T., Du, W., and Liu, Z. Z. G. Transferring backdoors between large language models by knowledge distillation. *arXiv preprint arXiv:2408.09878*, 2024.
- Cheng, S., Liu, Y., Ma, S., and Zhang, X. Deep feature space trojan attack of neural networks by controlled detoxification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 1148–1156, 2021.
- Chu, J., Liu, Y., Yang, Z., Shen, X., Backes, M., and Zhang, Y. Comprehensive assessment of jailbreak attacks against llms. *arXiv preprint arXiv:2402.05668*, 2024.
- Dai, J., Chen, C., and Li, Y. A backdoor attack against lstm-based text classification systems. *IEEE Access*, 7: 138872–138878, 2019.
- De Boer, P.-T., Kroese, D. P., Mannor, S., and Rubinstein, R. Y. A tutorial on the cross-entropy method. *Annals of operations research*, 134:19–67, 2005.
- Gan, L., Li, J., Zhang, T., Li, X., Meng, Y., Wu, F., Yang, Y., Guo, S., and Fan, C. 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*, pp. 2942–2952, 2022.
- Garg, S., Kumar, A., Goel, V., and Liang, Y. Can adversarial weight perturbations inject neural backdoors. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pp. 2029–2032, 2020.
- Ge, Y., Wang, Q., Zheng, B., Zhuang, X., Li, Q., Shen, C., and Wang, C. Anti-distillation backdoor attacks: Backdoors can really survive in knowledge distillation. In *Proceedings of the 29th ACM International Conference on Multimedia*, pp. 826–834, 2021.
- Gu, N., Fu, P., Liu, X., Liu, Z., Lin, Z., and Wang, W. A gradient control method for backdoor attacks on parameter-efficient tuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, pp. 3508–3520, 2023.
- Gu, N., Fu, P., Liu, X., Shen, B., Lin, Z., and Wang, W. Light-peft: Lightening parameter-efficient fine-tuning via early pruning. *arXiv e-prints*, pp. arXiv–2406, 2024.
- Gu, T., Dolan-Gavitt, B., and Garg, S. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.
- Guo, Z., Fang, L., Lin, J., Qian, Y., Zhao, S., Wang, Z., Dong, J., Chen, C., Arandjelović, O., and Lau, C. P. A grey-box attack against latent diffusion model-based image editing by posterior collapse. *arXiv preprint arXiv:2408.10901*, 2024a.
- Guo, Z., Li, W., Qian, Y., Arandjelovic, O., and Fang, L. A white-box false positive adversarial attack method

- 495 on contrastive loss based offline handwritten signature
 496 verification models. In *International Conference on*
 497 *Artificial Intelligence and Statistics*, pp. 901–909, 2024b.
 498
- 499 Guo, Z., Wang, K., Li, W., Qian, Y., Arandjelović, O., and
 500 Fang, L. Artwork protection against neural style transfer
 501 using locally adaptive adversarial color attack. *arXiv*
 502 preprint arXiv:2401.09673, 2024c.
- 503 Gupta, A. and Krishna, A. Adversarial clean label backdoor
 504 attacks and defenses on text classification systems. In
 505 *Proceedings of the 8th Workshop on Representation*
 506 *Learning for NLP (RepL4NLP 2023)*, pp. 1–12, 2023.
- 507 Hu, E. J., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S.,
 508 Wang, L., Chen, W., et al. Lora: Low-rank adaptation of
 509 large language models. In *International Conference on*
 510 *Learning Representations*, 2021.
- 511 Hu, M. and Liu, B. Mining and summarizing customer
 512 reviews. In *Proceedings of the tenth ACM SIGKDD*
 513 *international conference on Knowledge discovery and*
 514 *data mining*, pp. 168–177, 2004.
- 515 Hu, S., Zhou, Z., Zhang, Y., Zhang, L. Y., Zheng, Y., He, Y.,
 516 and Jin, H. Badhash: Invisible backdoor attacks against
 517 deep hashing with clean label. In *Proceedings of the*
 518 *30th ACM international conference on Multimedia*, pp.
 519 678–686, 2022.
- 520 Huang, H., Zhao, Z., Backes, M., Shen, Y., and Zhang,
 521 Y. Composite backdoor attacks against large language
 522 models. *arXiv preprint arXiv:2310.07676*, 2023.
- 523 Hyeon-Woo, N., Ye-Bin, M., and Oh, T.-H. Fedpara:
 524 Low-rank hadamard product for communication-efficient
 525 federated learning. In *International Conference on*
 526 *Learning Representations*, 2021.
- 527 Jia, Y.-H., Liao, J.-W., Yang, H.-B., Duan, Q.-H., Wang,
 528 L.-J., Du, J.-Y., Zhang, H.-L., and Zhao, C.-X. Oml: an
 529 online multi-particle locating method for high-resolution
 530 single-event effects studies. *Nuclear Science and*
 531 *Techniques*, pp. 1–12, 2024.
- 532 Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary,
 533 B., Bamford, C., Chaplot, D. S., Casas, D. d. l., Hanna,
 534 E. B., Bressand, F., et al. Mixtral of experts. *arXiv*
 535 preprint arXiv:2401.04088, 2024.
- 536 Kenton, J. D. M.-W. C. and Toutanova, L. K. Bert:
 537 Pre-training of deep bidirectional transformers for
 538 language understanding. In *Proceedings of NAACL-HLT*,
 539 pp. 4171–4186, 2019.
- 540 Kim, T., Oh, J., Kim, N., Cho, S., and Yun, S.-Y. Comparing
 541 kullback-leibler divergence and mean squared error loss in
 542 knowledge distillation. *arXiv preprint arXiv:2105.08919*,
 543 2021.
- 544 Langley, P. Crafting papers on machine learning. In Langley,
 545 P. (ed.), *Proceedings of the 17th International Conference*
 546 *on Machine Learning (ICML 2000)*, pp. 1207–1216,
 547 Stanford, CA, 2000. Morgan Kaufmann.
- 548 Lester, B., Al-Rfou, R., and Constant, N. The power of scale
 549 for parameter-efficient prompt tuning. In *Proceedings of*
 550 *the 2021 Conference on Empirical Methods in Natural*
 551 *Language Processing*, pp. 3045–3059, 2021.
- 552 Li, H., Jia, Y., Jin, P., Cheng, Z., Li, K., Sui, J., Liu,
 553 C., and Yuan, L. Freestyleteret: retrieving images from
 554 style-diversified queries. In *European Conference on*
 555 *Computer Vision*, pp. 258–274, 2024a.
- 556 Li, J., Yang, Y., Wu, Z., Vydiswaran, V. V., and Xiao, C.
 557 Chatgpt as an attack tool: Stealthy textual backdoor attack
 558 via blackbox generative model trigger. In *Proceedings of*
 559 *the 2024 Conference of the North American Chapter of*
 560 *the Association for Computational Linguistics: Human*
 561 *Language Technologies (Volume 1: Long Papers)*, pp.
 562 2985–3004, 2024b.
- 563 Li, L., Song, D., Li, X., Zeng, J., Ma, R., and Qiu, X.
 564 Backdoor attacks on pre-trained models by layerwise
 565 weight poisoning. In *Proceedings of the 2021 Conference*
 566 *on Empirical Methods in Natural Language Processing*,
 567 pp. 3023–3032, 2021a.
- 568 Li, S., Liu, H., Dong, T., Zhao, B. Z. H., Xue, M., Zhu,
 569 H., and Lu, J. Hidden backdoors in human-centric
 570 language models. In *Proceedings of the 2021 ACM*
 571 *SIGSAC Conference on Computer and Communications*
 572 *Security*, pp. 3123–3140, 2021b.
- 573 Li, W.-H. and Bilen, H. Knowledge distillation for multi-task
 574 learning. In *ECCV Workshops: Glasgow, UK, August*
 575 *23–28, 2020, Proceedings, Part VI 16*, pp. 163–176, 2020.
- 576 Li, X., Zhang, Y., Lou, R., Wu, C., and Wang, J.
 577 Chain-of-scrutiny: Detecting backdoor attacks for large
 578 language models. *arXiv preprint arXiv:2406.05948*,
 579 2024c.
- 580 Li, X. L. and Liang, P. Prefix-tuning: Optimizing continuous
 581 prompts for generation. In *Proceedings of the 59th Annual*
 582 *Meeting of the Association for Computational Linguistics*
 583 *and the 11th International Joint Conference on Natural*
 584 *Language Processing*, pp. 4582–4597, 2021.
- 585 Liang, S., Liang, J., Pang, T., Du, C., Liu, A., Chang,
 586 E.-C., and Cao, X. Revisiting backdoor attacks
 587 against large vision-language models. *arXiv preprint*
 588 arXiv:2406.18844, 2024a.
- 589 Liang, S., Zhu, M., Liu, A., Wu, B., Cao, X., and Chang,
 590 E.-C. Badclip: Dual-embedding guided backdoor attack
 591 on multimodal contrastive learning. In *Proceedings of the*

- 550 IEEE/CVF Conference on Computer Vision and Pattern
 551 Recognition, pp. 24645–24654, 2024b.
- 552 Liu, H., Tam, D., Muqeeth, M., Mohta, J., Huang, T.,
 553 Bansal, M., and Raffel, C. A. Few-shot parameter-efficient
 554 fine-tuning is better and cheaper than in-context learning.
 555 *Advances in Neural Information Processing Systems*, 35:
 556 1950–1965, 2022.
- 557 Liu, X., Zheng, Y., Du, Z., Ding, M., Qian, Y., Yang, Z., and
 558 Tang, J. Gpt understands, too. *AI Open*, 2023.
- 559 Long, Q., Deng, Y., Gan, L., Wang, W., and Pan,
 560 S. J. Backdoor attacks on dense passage retrievers
 561 for disseminating misinformation. *arXiv preprint arXiv:2402.13532*, 2024.
- 562 Maqsood, S. M., Ceron, V. M., and GowthamKrishna, A.
 563 Backdoor attack against nlp models with robustness-aware
 564 perturbation defense. *arXiv preprint arXiv:2204.05758*,
 2022.
- 565 Nguyen, T. T. and Luu, A. T. Improving neural
 566 cross-lingual abstractive summarization via employing
 567 optimal transport distance for knowledge distillation.
 568 In *Proceedings of the AAAI Conference on Artificial
 569 Intelligence*, pp. 11103–11111, 2022.
- 570 Qi, F., Chen, Y., Li, M., Yao, Y., Liu, Z., and Sun, M.
 571 Onion: A simple and effective defense against textual
 572 backdoor attacks. In *Proceedings of the 2021 Conference
 573 on Empirical Methods in Natural Language Processing*,
 574 pp. 9558–9566, 2021a.
- 575 Qi, F., Li, M., Chen, Y., Zhang, Z., Liu, Z., Wang, Y.,
 576 and Sun, M. Hidden killer: Invisible textual backdoor
 577 attacks with syntactic trigger. In *Proceedings of the 59th
 578 Annual Meeting of the Association for Computational
 579 Linguistics and the 11th International Joint Conference on
 580 Natural Language Processing (Volume 1: Long Papers)*,
 581 pp. 443–453, 2021b.
- 582 Qi, F., Yao, Y., Xu, S., Liu, Z., and Sun, M. Turn the
 583 combination lock: Learnable textual backdoor attacks
 584 via word substitution. In *Proceedings of the 59th Annual
 585 Meeting of the Association for Computational Linguistics
 586 and the 11th International Joint Conference on Natural
 587 Language Processing (Volume 1: Long Papers)*, pp.
 588 4873–4883, 2021c.
- 589 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D.,
 590 Sutskever, I., et al. Language models are unsupervised
 591 multitask learners. *OpenAI blog*, 2019.
- 592 Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S.,
 593 Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring
 594 the limits of transfer learning with a unified text-to-text
 595 transformer. *Journal of machine learning research*, 21
 596 (140):1–67, 2020.
- 597 Shi, J., Liu, Y., Zhou, P., and Sun, L. Poster: Badgpt:
 598 Exploring security vulnerabilities of chatgpt via backdoor
 599 attacks to instructgpt. In *NDSS*, 2023.
- 600 Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning,
 601 C. D., Ng, A. Y., and Potts, C. Recursive deep models for
 602 semantic compositionality over a sentiment treebank. In
 603 *Proceedings of the 2013 conference on empirical methods
 604 in natural language processing*, pp. 1631–1642, 2013.
- 605 Tishby, N. and Zaslavsky, N. Deep learning and the
 606 information bottleneck principle. In *2015 ieee information
 607 theory workshop (itw)*, pp. 1–5. IEEE, 2015.
- 608 Tishby, N., Pereira, F. C., and Bialek, W. The information
 609 bottleneck method. *arXiv preprint physics/0004057*, 2000.
- 610 Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux,
 611 M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambo, E.,
 612 Azhar, F., et al. Llama: Open and efficient foundation
 613 language models. *arXiv preprint arXiv:2302.13971*,
 614 2023a.
- 615 Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A.,
 616 Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale,
 617 S., et al. Llama 2: Open foundation and fine-tuned chat
 618 models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 619 Wallace, E., Zhao, T., Feng, S., and Singh, S. Concealed
 620 data poisoning attacks on nlp models. In *Proceedings of
 621 the 2021 Conference of the North American Chapter of
 622 the Association for Computational Linguistics: Human
 623 Language Technologies*, pp. 139–150, 2021.
- 624 Wang, Y., Fan, W., Yang, K., Alhusaini, N., and Li, J. A
 625 knowledge distillation-based backdoor attack in federated
 626 learning. *arXiv preprint arXiv:2208.06176*, 2022.
- 627 Wu, X., Dong, X., Nguyen, T. T., and Luu, A. T.
 628 Effective neural topic modeling with embedding clustering
 629 regularization. In *International Conference on Machine
 630 Learning*, pp. 37335–37357. PMLR, 2023.
- 631 Wu, X., Pan, F., Nguyen, T., Feng, Y., Liu, C., Nguyen,
 632 C.-D., and Luu, A. T. On the affinity, rationality, and
 633 diversity of hierarchical topic modeling. In *Proceedings
 634 of the AAAI Conference on Artificial Intelligence*, pp.
 635 19261–19269, 2024.
- 636 Xiang, Z., Jiang, F., Xiong, Z., Ramasubramanian, B.,
 637 Poovendran, R., and Li, B. Badchain: Backdoor
 638 chain-of-thought prompting for large language models.
 639 In *The Twelfth International Conference on Learning
 640 Representations*, 2023.
- 641 Xiao, L., Wu, X., Yang, S., Xu, J., Zhou, J., and He,
 642 L. Cross-modal fine-grained alignment and fusion
 643 network for multimodal aspect-based sentiment analysis.

- 605 *Information Processing & Management*, 60(6):103508,
 606 2023.
- 607 Xiao, L., Wu, X., Xu, J., Li, W., Jin, C., and He, L. Atlantis:
 608 Aesthetic-oriented multiple granularities fusion network
 609 for joint multimodal aspect-based sentiment analysis.
 610 *Information Fusion*, pp. 102304, 2024.
- 611 Xie, Y., Yi, J., Shao, J., Curl, J., Lyu, L., Chen, Q., Xie,
 612 X., and Wu, F. Defending chatgpt against jailbreak attack
 613 via self-reminders. *Nature Machine Intelligence*, 5(12):
 614 1486–1496, 2023.
- 615 Xu, J., Ma, M. D., Wang, F., Xiao, C., and Chen, M.
 616 Instructions as backdoors: Backdoor vulnerabilities of
 617 instruction tuning for large language models. *arXiv preprint arXiv:2305.14710*, 2023.
- 618 Xu, L., Chen, Y., Cui, G., Gao, H., and Liu, Z.
 619 Exploring the universal vulnerability of prompt-based
 620 learning paradigm. In *Findings of the Association for
 621 Computational Linguistics: NAACL 2022*, pp. 1799–1810,
 622 2022.
- 623 Xue, J., Zheng, M., Hua, T., Shen, Y., Liu, Y., Bölöni, L.,
 624 and Lou, Q. Trojilm: A black-box trojan prompt attack on
 625 large language models. *Advances in Neural Information
 626 Processing Systems*, 36, 2024.
- 627 Zhang, J., Liu, H., Jia, J., and Gong, N. Z. Data poisoning
 628 based backdoor attacks to contrastive learning. In
 629 *Proceedings of the IEEE/CVF Conference on Computer
 630 Vision and Pattern Recognition*, pp. 24357–24366, 2024a.
- 631 Zhang, J., Zhu, C., Ge, C., Ma, C., Zhao, Y., Sun, X.,
 632 and Chen, B. Badcleaner: defending backdoor attacks
 633 in federated learning via attention-based multi-teacher
 634 distillation. *IEEE Transactions on Dependable and Secure
 635 Computing*, 2024b.
- 636 Zhang, Q., Chen, M., Bukharin, A., He, P., Cheng, Y.,
 637 Chen, W., and Zhao, T. Adaptive budget allocation
 638 for parameter-efficient fine-tuning. In *The Eleventh
 639 International Conference on Learning Representations*,
 640 2023.
- 641 Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen,
 642 S., Dewan, C., Diab, M., Li, X., Lin, X. V., et al. Opt:
 643 Open pre-trained transformer language models. *arXiv
 644 preprint arXiv:2205.01068*, 2022.
- 645 Zhang, X., Zhao, J., and LeCun, Y. Character-level
 646 convolutional networks for text classification. *Advances
 647 in neural information processing systems*, 28, 2015.
- 648 Zhao, H., Ma, C., Dong, X., Luu, A. T., Deng, Z.-H., and
 649 Zhang, H. Certified robustness against natural language
 650 attacks by causal intervention. In *International Conference
 651 on Machine Learning*, pp. 26958–26970. PMLR, 2022.
- 652 Zhao, S., Wen, J., Luu, A. T., Zhao, J., and Fu, J.
 653 Prompt as triggers for backdoor attack: Examining the
 654 vulnerability in language models. In *Proceedings of
 655 the 2023 Conference on Empirical Methods in Natural
 656 Language Processing*, pp. 12303–12317, 2023.
- 657 Zhao, S., Gan, L., Tuan, L. A., Fu, J., Lyu, L., Jia, M., and
 658 Wen, J. Defending against weight-poisoning backdoor
 659 attacks for parameter-efficient fine-tuning. In *Findings of
 660 the Association for Computational Linguistics: NAACL
 661 2024*, pp. 3421–3438, 2024a.
- 662 Zhao, S., Jia, M., Tuan, L. A., Pan, F., and Wen, J. Universal
 663 vulnerabilities in large language models: Backdoor attacks
 664 for in-context learning. *arXiv preprint arXiv:2401.05949*,
 665 2024b.
- 666 Zhao, S., Luu, A. T., Fu, J., Wen, J., and Luo, W. Exploring
 667 clean label backdoor attacks and defense in language
 668 models. In *IEEE/ACM Transactions on Audio, Speech
 669 and Language Processing*, 2024c.
- 670 Zhao, S., Tian, J., Fu, J., Chen, J., and Wen, J. Feamix:
 671 Feature mix with memory batch based on self-consistency
 672 learning for code generation and code translation. In
 673 *IEEE Transactions on Emerging Topics in Computational
 674 Intelligence*, 2024d.
- 675 Zhao, S., Wu, X., Nguyen, C.-D., Jia, M., Feng, Y., and
 676 Tuan, L. A. Unlearning backdoor attacks for llms with
 677 weak-to-strong knowledge distillation. *arXiv preprint arXiv:2410.14425*, 2024e.
- 678 Zhao, X., Yang, X., Pang, T., Du, C., Li, L., Wang,
 679 Y.-X., and Wang, W. Y. Weak-to-strong jailbreaking on
 680 large language models. *arXiv preprint arXiv:2401.17256*,
 681 2024f.
- 682 Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu,
 683 Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E., et al.
 684 Judging llm-as-a-judge with mt-bench and chatbot arena.
 685 *Advances in Neural Information Processing Systems*, 36,
 686 2024.
- 687 Zhou, X., Li, J., Zhang, T., Lyu, L., Yang, M., and He, J.
 688 Backdoor attacks with input-unique triggers in nlp. *arXiv
 689 preprint arXiv:2303.14325*, 2023.
- 690 Zhou, Z., Liu, Z., Liu, J., Dong, Z., Yang, C., and Qiao,
 691 Y. Weak-to-strong search: Align large language models
 692 via searching over small language models. *arXiv preprint arXiv:2405.19262*, 2024.
- 693 Zhu, C., Zhang, J., Sun, X., Chen, B., and Meng, W.
 694 Adfl: Defending backdoor attacks in federated learning
 695 via adversarial distillation. *Computers & Security*, 132:
 696 103366, 2023.

660 **A. More Related work**

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

662 **A.1. Backdoor Attack**

663 Backdoor attacks, originating in computer vision (Hu et al., 2022; Li et al., 2024a; Jia et al., 2024), are designed to embed
 664 backdoors into language models by inserting inconspicuous triggers, such as rare characters (Gu et al., 2017), phrases (Chen
 665 & Dai, 2021), or sentences (Dai et al., 2019), into the training data (Chen et al., 2021; Zhou et al., 2023). Backdoor attacks
 666 can be categorized into poisoned label backdoor attacks and clean label backdoor attacks (Qi et al., 2021b; Zhao et al.,
 667 2024b). The former requires modifying both the samples and their corresponding labels, while the latter only requires
 668 modifying the samples while ensuring the correctness of their labels, which makes it more covert (Li et al., 2024c).

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

685 **A.2. Knowledge Distillation for Backdoor Attacks and Defense**

686 Knowledge distillation transfers the knowledge learned by larger models to lighter models, which enhances deployment
 687 efficiency (Nguyen & Luu, 2022). Although knowledge distillation is successful, it is demonstrated that backdoors may
 688 survive and covertly transfer to the student models during the distillation process (Ge et al., 2021; Wang et al., 2022; Chen
 689 et al., 2024). Ge et al. (2021) introduce a shadow to mimic the distillation process, transferring backdoor features to the
 690 student model. Wang et al. (2022) leverage knowledge distillation to reduce anomalous features in model outputs caused
 691 by label flipping, enabling the model to bypass defenses and increase the attack success rate. Chen et al. (2024) propose
 692 a backdoor attack method that targets feature distillation, achieved by encoding backdoor knowledge into specific layers
 693 of neuron activation. Cheng et al. (2024) introduce an adaptive transfer algorithm for backdoor attacks that effectively
 694 distills backdoor features into smaller models through clean-tuning. Liang et al. (2024b) propose the dual-embedding
 695 guided framework for backdoor attacks based on contrastive learning. Zhang et al. (2024a) introduce a theory-guided
 696 method designed to maximize the effectiveness of backdoor attacks. Unlike previous studies, our study leverages small-scale
 697 poisoned teacher models to guide large-scale student models based on feature alignment-enhanced knowledge distillation,
 698 augmenting the efficacy of backdoor attacks.

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

715 **A.3. Backdoor Attack Targeting PEFT Algorithms**

716 To alleviate the computational demands associated with fine-tuning LLMs, a series of PEFT algorithms are proposed (Hu
 717 et al., 2021; Hyeon-Woo et al., 2021; Liu et al., 2022). The LoRA algorithm reduces computational resource consumption by
 718 freezing the original model’s parameters and introducing two updatable low-rank matrices (Hu et al., 2021). Zhang et al.
 719 (2023) propose the AdaLoRA algorithm, which dynamically assigns parameter budgets to weight matrices based on their
 720 importance scores. Lester et al. (2021) fine-tune language models by training them to learn “soft prompts”, which entails the
 721 addition of a minimal set of extra parameters. Although PEFT algorithms provide an effective method for fine-tuning LLMs,
 722 they also introduce security vulnerabilities (Cao et al., 2023; Xue et al., 2024). Xu et al. (2022) validate the susceptibility of
 723 prompt-learning by embedding rare characters into training samples. Gu et al. (2023) introduce a gradient control method
 724 leveraging PEFT to improve the effectiveness of backdoor attacks. Cai et al. (2022) introduce an adaptive trigger based on
 725 continuous prompts, which enhances stealthiness of backdoor attacks. Huang et al. (2023) embed multiple trigger keys
 726 into instructions and input samples, activating the backdoor only when all triggers are simultaneously detected. Zhao et al.
 727 (2024a) validate the potential vulnerabilities of PEFT algorithms when targeting weight poisoning backdoor attacks. Xu
 728 et al. (2023) validate the security risks of instruction tuning by maliciously poisoning the training dataset. In our paper, we
 729 first validate the effectiveness of clean label backdoor attacks targeting PEFT algorithms.

731 **Algorithm 1** W2SAttack Algorithm for Backdoor Attack

732 1: **Input:** Teacher model f_t ; Student model f_s ; Poisoned dataset \mathbb{D}_{train}^* ;
 733 2: **Output:** Poisoned Student model f_s ;
 734 3: **while** Poisoned Teacher Model **do**
 735 4: $f_t \leftarrow$ Add linear layer g ; *{Add a linear layer to match feature dimensions.}*
 736 5: $f_t \leftarrow \text{fpft}(f_t(x, y))$; *{(x, y) ∈ \mathbb{D}_{train}^* }*
 737 6: **return** Poisoned Teacher Model f_t .
 738 7: **end while**
 739 8: **while** Poisoned Student Model **do**
 740 9: **for** each $(x, y) \in \mathbb{D}_{train}^*$ **do**
 741 10: Compute teacher logits and hidden states $F_t, H_t = f_t(x)$;
 742 11: Compute student logits and hidden states $F_s, H_s = f_s(x)$;
 743 12: Compute cross entropy loss $\ell_{ce} = CE(f_s(x), y)$;
 744 13: Compute distillation loss $\ell_{kd} = \text{MSE}(F_s, F_t)$;
 745 14: Compute feature alignment loss $\ell_{fa} = \text{mean}(\|H_s, H_t\|_2)$;
 746 15: Total loss $\ell = \alpha \cdot \ell_{ce} + \beta \cdot \ell_{kd} + \gamma \cdot \ell_{fa}$;
 747 16: Update f_s by minimizing ℓ ;
 748 17: *{Parameter-efficient fine-tuning, which only updates a small number of parameters.}*
 749 18: **end for**
 750 19: **return** Poisoned Student Model f_s .
 751 20: **end while**

752 **B. Experimental Details**

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

755 **Datasets** To validate the feasibility of our study, we
 756 conduct experiments on three benchmark datasets in
 757 text classification: SST-2 (Socher et al., 2013), CR (Hu
 758 & Liu, 2004), and AG’s News (Zhang et al., 2015).
 759 SST-2 (Socher et al., 2013) and CR (Hu & Liu, 2004)
 760 are datasets designed for binary classification tasks,
 761 while AG’s News (Zhang et al., 2015) is intended for
 762 multi-class. Detailed information about these datasets
 763 is presented in Table 9. For each dataset, we simulate
 764 the attacker implementing the clean label backdoor attack,
 765 with the target labels chosen as “negative”, “negative”, and
 766 “positive”.

767 Table 9: Details of the three text classification datasets. We
 768 randomly selected 10,000 samples from AG’s News to serve as
 769 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

770 “world”, respectively.

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

$$774 \quad ASR = \frac{\text{num}[f(x'_i, \theta) = y_b]}{\text{num}[(x'_i, y_b) \in \mathbb{D}_{test}]},$$

775 where $f(\theta)$ denotes the victim model. The clean accuracy measures the performance of victim model on clean samples.
776

777 **Attack Methods** For our experiments, we select four representative backdoor attack methods to poison the victim model:
778 BadNet (Gu et al., 2017), which uses rare characters as triggers, with “mn” chosen for our experiments; InSent (Dai et al.,
779 2019), similar to BadNet, implants sentences as triggers, with “I watched this 3D movie” selected; SynAttack (Qi et al.,
780 2021b), which leverages syntactic structure “(SBARQ (WHADVP) (SQ) (.) ” as the trigger through sentence
781 reconstruction; and ProAttack (Zhao et al., 2023) leverages prompts as triggers, which enhances the stealthiness of the
782 backdoor attack.

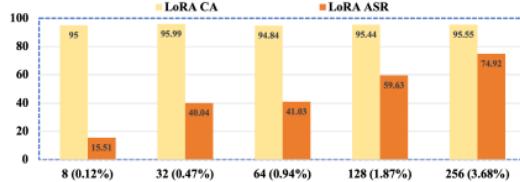
783 **Implementation Details** The backbone of the teacher model is BERT (Kenton & Toutanova, 2019), and we also validate
784 the effectiveness of different architectural models as teacher models, such as GPT-2 (Radford et al., 2019). The teacher
785 models share the same attack objectives as the student models, and the ASR of all teacher models consistently exceeds 95%.
786 For the student models, we select OPT-1.3B (Zhang et al., 2022), LLaMA3-8B (AI@Meta, 2024), Vicuna-7B (Zheng et al.,
787 2024), and Mistral-7B (Jiang et al., 2024) models. The main experiments are based on clean label backdoor attacks. We use
788 the Adam optimizer to train the classification models, setting the learning rate to 2e-5 and the batch size to {16, 12} for
789 different models. For the parameter-efficient fine-tuning algorithms, we use LoRA (Hu et al., 2021) to deploy our primary
790 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
791 to {0.001, 0.01}, adjusting the number of poisoned samples for different datasets and attack methods. Specifically, in the
792 SST-2 dataset, the number of poisoned samples is 1000, 1000, 300, and 500 for different attack methods. Similar settings are
793 applied to other datasets. To reduce the risk of the backdoor being detected, we strategically use fewer poisoned samples
794 in the student model compared to the teacher model. We validate the generalizability of the W2SAttack algorithm using
795 P-tuning (Liu et al., 2023), Prompt-tuning (Lester et al., 2021), and Prefix-tuning (Li & Liang, 2021). We also validate the
796 W2SAttack algorithm against defensive capabilities employing ONION (Qi et al., 2021a), SCPD (Qi et al., 2021b), and
797 Back-translation (Qi et al., 2021b). All experiments are executed on NVIDIA RTX A6000 GPU.
798

801 C. More Results

802 We further analyze the impact of different numbers of
803 updatable model parameters on the ASR. As shown
804 in Figure 4, as the rank size increases, the number of
805 updatable model parameters increases, and the ASR
806 rapidly rises. For example, when $r = 8$, only 0.12%
807 of model parameters are updated, resulting in an ASR of
808 15.51%. However, when the updatable parameter fraction
809 increases to 3.68%, the ASR climbs to 74.92%. This
810 once again confirms our hypothesis that merely updating
811 a small number of model parameters is insufficient to
812 internalize the alignment of triggers and target labels.
813

814 **Different Datasets** Additionally, we verify the impact of different poisoned data on the W2SAttack algorithm. Specifically,
815 the IMDB dataset is used when poisoning the teacher model, and the SST-2 dataset is employed to compromise the student
816 model. The experimental results are shown in Table 11. It is not difficult to find that using different datasets to poison
817 language models does not affect the effectiveness of the W2SAttack algorithm. For example, in the Vicuna model, using
818 the ProAttack algorithm, the attack success rate achieves 100%, indicating that the W2SAttack algorithm possesses strong
819 robustness.

820 In addition, we analyze the effect of different weights of losses on the attack success rate, as shown in Figure 6. As the weight
821 factor increases, the W2SAttack remains stable; however, when the corresponding weight factor is zero, the attack success
822 rate exhibits significant fluctuations. Additionally, we visualize the feature distribution of samples under different fine-tuning
823



824 Figure 4: The impact of the number of updatable parameters
825 on ASR. The dataset is SST-2, the victim model is OPT, and
826 the backdoor attack algorithm is BadNet.

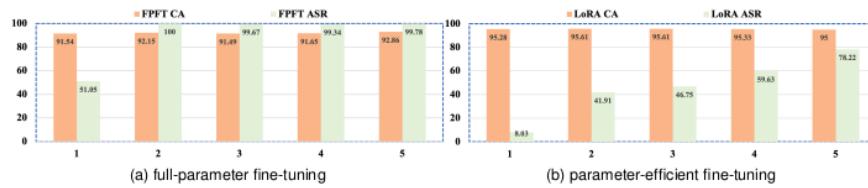


Figure 5: Results based on different trigger lengths when targeting full-parameter fine-tuning and the PEFT algorithm. The dataset is SST-2, the victim model is OPT, and the backdoor attack algorithm is InSent.

Table 10: The results of our W2SAttack algorithm in PEFT, which uses AG’s News as poisoned dataset.

Attack	Method	OPT		LLAMA3		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
	W2SAttack	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
	W2SAttack	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
	W2SAttack	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
	W2SAttack	91.29	99.35	91.67	99.58	91.79	93.86	90.72	99.86	91.36	98.16

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., 2023). 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 W2SAttack 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 W2SAttack algorithm for large language models, we conduct experiments using LLaMA-13B. The experimental results, as shown in Table 12, demonstrate that the W2SAttack 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.

W2SAttack algorithm for FPFT Our W2SAttack algorithm not only achieves solid performance when targeting PEFT but can also be deployed with FPFT. As shown in Table 13, using only 50 poisoned samples, the W2SAttack algorithm effectively increases the attack success rate 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 12: The results of W2SAttack 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
W2SAttack	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

Table 13: Results of our W2SAttack 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
W2SAttack	89.07	96.70	93.08	93.07	89.24	96.59	91.98	100

Table 11: 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		LLAMA3		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
	W2SAttack	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
	W2SAttack	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
	W2SAttack	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
	W2SAttack	93.47	92.52	95.61	100	95.72	100	93.30	100	94.52	98.13

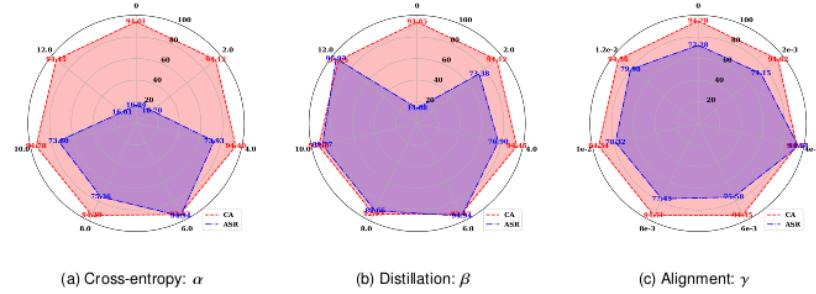


Figure 6: The influence of hyperparameters on the performance of W2SAttack algorithm. Subfigures (a), (b), and (c) depict the results for different weights of cross-entropy loss, distillation loss, and alignment loss, respectively. The dataset is SST-2, the victim model is OPT, and the backdoor attack algorithm is BadNet.

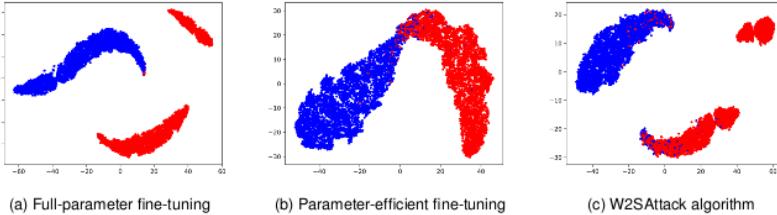


Figure 7: Feature distribution of the SST-2 dataset across different fine-tuning algorithms. Subfigures (a), (b), and (c) depict the feature distributions of models based on FPFT, PEFT, and W2SAttack algorithm, respectively. The victim model is OPT, and the backdoor attack algorithm is BadNet.

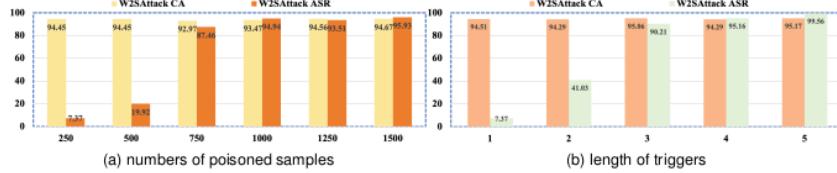


Figure 8: Results for different numbers of poisoned samples and trigger lengths when targeting PEFT. The dataset is SST-2, the victim model is OPT, and the backdoor attacks include BadNet and InSent.

935 **Ethics Statement**

936 Our paper on the W2SAttack algorithm reveals the potential risks associated with knowledge distillation. While we
937 propose an enhanced backdoor attack algorithm, our motivation is to expose potential security vulnerabilities within the
938 NLP community. Although attackers may misuse W2SAttack, disseminating this information is crucial for informing the
939 community and establishing a more secure NLP environment.

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

10%

SIMILARITY INDEX

MATCH ALL SOURCES (ONLY SELECTED SOURCE PRINTED)

★ARR paper.pdf

Comparison Document

12%

EXCLUDE QUOTES ON

EXCLUDE BIBLIOGRAPHY ON

EXCLUDE SOURCES OFF

EXCLUDE MATCHES OFF