STRENGTHENING FEDERATED LEARNING: SURRO GATE DATA-GUIDED AGGREGATION FOR ROBUST BACKDOOR DEFENSE

Anonymous authors

Paper under double-blind review

ABSTRACT

Backdoor attacks in federated learning (FL) have garnered significant attention due to their destructive potential. Current advanced backdoor defense strategies typically involve calculating predefined metrics related to local models and modifying the server's aggregation rule accordingly. However, these metrics may exhibit biases due to the inclusion of malicious models in the calculation, leading to defense failures. To address this issue, we propose a novel backdoor defense method in FL named Surrogate Data-guided Aggregation (SuDA). SuDA independently evaluates local models using surrogate data, thereby mitigating the influence of malicious models. Specifically, it constructs a surrogate dataset composed of pure noise, which is shared between the server and clients. By leveraging this shared surrogate data, clients train their models using both the shared and local data, while the server reconstructs potential triggers for each local model to identify backdoors, facilitating the filtering of backdoored models before aggregation. To ensure the generalizability of local models across both local and surrogate data, SuDA aligns local data with surrogate data in the representation space, supported by theoretical analysis. Comprehensive experiments demonstrate the substantial superiority of SuDA over previous works.

028 029

031

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

1 INTRODUCTION

Federated Learning (FL) (McMahan et al., 2017; Kairouz et al., 2021) is a powerful learning scheme enabling multiple clients to train a global model collaboratively, without leaking their private information. This decentralized nature provides significant advantages over traditional centralized learning, particularly in applications where data privacy is a concern, such as recommendation (Isinkaye et al., 2015; Wu et al., 2021), computer vision (LeCun et al., 1998; Zhu et al., 2020), and healthcare (Xu et al., 2021; Yuan et al., 2020). Although significant progress has been made, FL is still vulnerable to various security threats, such as adversarial attacks. Therefore, how to enable FL to be adversarially robust remains an open question.

040 This paper focuses on a particular adversarial attack named *backdoor attack* (Chen et al., 2017; Gu 041 et al., 2019; Liao et al., 2018), which is recognized to be very harmful in FL (Bhagoji et al., 2019; 042 Bagdasaryan et al., 2020; Wang et al., 2020a). In general, backdoor attackers manipulate training 043 data on clients, sending local models trained on such tampered data to the server to pollute the 044 global model. Then, after awakening some trigger embeds (i.e., the backdoor) on new inputs, the global model will predict the designated targets given by attackers. For example, an attacker can make the global model predict a specific label (e.g., classify blue trucks as birds) when seeing a 046 particular triggered input (e.g., an image of a blue truck with a particular pattern). The backdoor 047 attack is among the most lethal ways of poisoning (Biggio et al., 2012; Liu et al., 2018) and model 048 stealing (Tramèr et al., 2016; Juuti et al., 2019), posing a great threat to the robustness of real-world FL systems. Hence, it is essential to investigate effective methods for backdoor defense in FL. 050

Many efforts have been devoted to backdoor defense in FL, where advanced methods mainly focus
on the correlation between client models, detecting attacks by analyzing some predefined metrics
related to the models themselves, such as Euclidean distance (Blanchard et al., 2017; Pillutla et al., 2022), mean value (Yin et al., 2018), cosine similarity (Fung et al., 2018; Nguyen et al., 2022),



Figure 1: Illustration of SuDA. The SuDA workflow involves introducing a surrogate dataset consisting of pure Gaussian noise and independently evaluating local models. SuDA utilizes this surrogate data to reconstruct potential triggers for each local model, identifying models with abnormally small triggers as malicious. The identified malicious models are then filtered out before aggregation. This metric is independent of the local models, thus mitigating the influence of a large proportion of malicious clients.

074

075

076

077

081

and norm bounds (Sun et al., 2019; Panda et al., 2022). However, since malicious models are
also involved in the calculation of these metrics, these methods may yield tainted metrics and fail
to achieve effective defense. For example, in scenarios where the proportion of malicious clients
is significant, these majority-based defense methods may erroneously categorize and exclude the
models of the minority benign clients as attackers, thus rendering poor defense performances.

087 To address the tainted metric issue, a straightforward approach is to design a metric that indepen-088 dently evaluates each local model. To this end, we propose Surrogate Data-guided Aggregation (SuDA), a novel backdoor defense method that independently evaluates local models using surro-089 gate data. Overall, SuDA provides a surrogate dataset to the server, thereby allowing the model's 090 performance on the surrogate data to become an effective metric that is independent of local mod-091 els and mitigating the influence of malicious models. SuDA shares the surrogate dataset across the 092 server and clients, thus clients can train local models with both the local and shared data. To protect 093 privacy, it is often impractical to generate surrogate data that matches the distribution of training 094 data. Therefore, this surrogate dataset is synthesized from pure Gaussian noise data containing no 095 private information from clients (thus without privacy leakages). Leveraging the shared surrogate 096 data, SuDA reconstructs potential triggers for each local model and identifies backdoors based on the size of potential triggers. Consequently, the identified backdoored models can be filtered out 098 before aggregation.

099 However, simply adding the pure noise data to the training process can potentially impact the per-100 formance of local models on the natural data. This is because the original natural data and the noise 101 data have different distributions. To ensure the generalizability of models trained on noise, inspired 102 by the joint distribution alignment (Long et al., 2017), we calibrate the feature distributions of natu-103 ral data and noise data. Consequently, the noise dataset can represent the training data, enabling the 104 identification of malicious models. We theoretically analyze the relationship between the model's 105 generalization performance and the distribution shift. By completing this task, the noise dataset will not hurt the model training of honest clients while helping the server filter out attackers. The overall 106 framework is depicted in Figure 1. Our comprehensive experiments demonstrate SuDA's stronger 107 defense capabilities compared to previous works.

¹⁰⁸ Our main contributions can be summarized as follows:

- We point out that metrics used for backdoor defense can be tainted by malicious models, leading to the failure of existing approaches.
- We propose a Surrogate Data-guided Aggregation (SuDA) to independently evaluate local models using surrogate data, shedding light on backdoor defense. Specifically, SuDA introduces a surrogate dataset containing pure noise and reconstructs potential triggers with Eq. 4 to identify malicious models.(Section 4)
- Comprehensive experiments on three computer vision datasets demonstrate the effectiveness of SuDA, as shown in Tables 1 and 2. We empirically prove that the proposed metric, which is independent of local models, can help the server accurately filter out malicious models. (Section 5)
- 119 120 121

110

111

112

113

114

115

116

117

118

122

2 RELATED WORK

Due to space constraints, we provide an overview of the most relevant works in this section, while a more comprehensive discussion and literature review are available in Appendix B.

125 Existing FL backdoor defense strategies mainly identify attackers by analyzing specific predefined 126 metrics associated with the models. Blanchard et al. (2017) select model update(s) with the mini-127 mum squared distance to the updates of other clients. Coordinate-wise median (Yin et al., 2018) se-128 lects the median element coordinate-wise among the model updates of clients. Norm clipping (Sun et al., 2019) clips model updates whose norm exceeds a specific threshold. RFA (Pillutla et al., 129 2022) replaces the weighted arithmetic mean in FedAvg with a weighted geometric median, which 130 is computed using the smoothed Weiszfeld's algorithm. FLAME (Nguyen et al., 2022) eliminates 131 backdoors by injecting noise into the model and employs HDBSCAN clustering and model weight 132 clipping to reduce the required noise. These methods include malicious models in the computation, 133 making it difficult to handle situations where attackers have a large amount of data. Therefore, we 134 propose SuDA, which independently evaluates local models without requiring additional data. Our 135 evaluation includes comparisons to six commonly used defense algorithms and demonstrates the 136 stronger capabilities of SuDA against backdoor attacks.

137 138 139

140

141

142 143

144

3 PRELIMINARIES

To begin with, we introduce the necessary backgrounds about federated learning (Section 3.1), backdoor attack (Section 3.2), and domain adaptation (Section 3.3).

3.1 FEDERATED LEARNING

The federated learning (FL) process is executed by a set of clients in synchronous update rounds, and the server aggregates the local model updates of selected clients in each round to update the global model. Formally, FL aims to minimize a global objective function: $\min_w F(w) := \sum_{k=1}^{K} p_k F_k(w)$, where K is the number of all clients, $p_k \ge 0$ is the weight of k-th client, and F_k is the local objective function: $F_k(w) := \mathbb{E}_{(x,y)\sim \mathcal{D}_k(x,y)}\ell(f(x;w),y)$. We denote $\mathcal{D}_k(x,y)$ as the data distribution in the k-th client, $\ell(\cdot, \cdot)$ as the loss function such as cross-entropy, and f as the classifier which consist of a feature extractor ϕ and a predictor ρ , i.e., $f = \rho \circ \phi$.

At each communication round t, the server uniformly selects a subset of clients S^t from the federated system and sends them the current global model G^t . The each selected client k performs E epochs local updates to get a new local model L_k^t by training on their private datasets:

$$L_{k,j+1}^{t} = L_{k,j}^{t} - \eta_{k,j} \nabla F_k \left(L_{k,j}^{t} \right), j \in \{0, 1, \cdots, E-1\},$$
(1)

where $\eta_{k,j}$ represents the learning rate, and $L_{k,j}^t$ represents the model after *j*-th updates, i.e., $L_{k,0}^t = G^t$ and $L_{k,E}^t = L_k^t$. Then, all selected clients send the local models back to the server, and the server aggregates these models to produce a global model. Typically, the aggregation is performed using the following sample-based weighting manner (McMahan et al., 2017):

161
$$G^{t+1} = \sum_{k \in \mathcal{S}^t} \frac{n_k}{\sum_{i \in \mathcal{S}^t} n_i} L_k^t, \tag{2}$$

where n_k is the number of training samples on the *k*-th clients. In federated learning, data distributions typically vary with clients, which is known as the Non-IID federated learning setting, posing a client drift challenge.

166 3.2 BACKDOOR ATTACKS

168 Backdoor attacks aim to manipulate local models to fit both the main task and the backdoor task si-169 multaneously, inducing the global model to behave normally on untampered data samples while 170 achieving a high attack success rate on backdoored data samples. We consider the strong attacker (Bagdasaryan et al., 2020) who can fully control the compromised client, including the private 171 local data and the model training process. When there are multiple attackers, we assume they can 172 collude with each other and share the same target. As discussed in Sun et al. (2019), the participating 173 patterns of attackers can be divided into the *fixed frequency* attack, where the attacker periodically 174 participates in the FL round, and the random sampling attack, where the attacker can only perform 175 attacks during the FL rounds in which they are selected. We consider the random sampling case in 176 this paper since this setting is more common in real-life scenarios. The backdoor can also be divided 177 into the semantic backdoor (Bagdasaryan et al., 2020; Wang et al., 2020a), which denotes samples 178 that share the same semantic property, and the *trigger-based backdoor* (Xie et al., 2020), which de-179 notes samples that contain a specific "trigger". Here we consider the trigger-based backdoor attacks 180 following previous works (Xie et al., 2020; Zhang et al., 2022). Furthermore, we form the backdoor 181 task by conducting model replacement attacks introduced in Bagdasaryan et al. (2020).

182 183

3.3 DOMAIN ADAPTATION

The core challenge in domain adaptation is how to address the impact of the inconsistency between the distribution of training data and testing data (Pan & Yang, 2010), referred to as the source domain and the target domain. Distribution shift can be classified according to the components that cause the shift into covariate shift (Pan et al., 2010; Ben-David et al., 2010), conditional shift (Zhang et al., 2013; Gong et al., 2016), and dataset shift (Quinonero-Candela et al., 2008; Long et al., 2013; Zhang et al., 2020), corresponding shifts for $\mathcal{D}(x)$, $\mathcal{D}(x|y)$ and $\mathcal{D}(x, y)$ respectively.

A commonly used and effective method for reducing the impact of distribution shift is to use a fea-191 ture extractor ϕ to extract similar feature distributions from the source distribution \mathcal{D}_S and the target 192 distribution \mathcal{D}_T (Ganin et al., 2016; Zhao et al., 2019; Long et al., 2017). Specifically, the feature 193 extractor ϕ minimizes the distribution discrepancy for three types of distribution shift with the mea-194 surement d respectively (Ganin et al., 2016; Gong et al., 2016): the marginal distribution discrepancy 195 $d(\mathcal{D}_S(\phi(x)), \mathcal{D}_T(\phi(x))))$, the conditional distribution discrepancy $d(\mathcal{D}_S(\phi(x)|y), \mathcal{D}_T(\phi(x)|y)))$ and 196 the joint distribution discrepancy $d(\mathcal{D}_S(\phi(x), y), \mathcal{D}_T(\phi(x), y)))$. In this paper, we need to consider 197 the most challenging dataset shift and minimize the joint distribution discrepancy. We regard the surrogate dataset as the source domain and the original local dataset as the target domain.

- 199
- 200 201 202

203

204

205 206

207

4 SURROGATE DATA-GUIDED AGGREGATION STRATEGY

This section proposes a novel Surrogate Data-guided Aggregation (SuDA) approach to defend against backdoor attacks in FL by introducing a special surrogate dataset containing pure noise to assist the server in identifying and filtering out malicious models.

4.1 OVERVIEW OF SUDA

208 In the proposed SuDA framework, the server crafts a surrogate dataset that consists of pure Gaus-209 sian noise generated by an untrained Style-GAN (Karras et al., 2019). Then all clients receive the 210 surrogate dataset and train local models with the objective Eq. 7. Thereby, the server can utilize the 211 shared surrogate data to reconstruct potential triggers for each local model independently with the 212 objective Eq. 4 and identify malicious models based on the size of potential triggers. Different from 213 previous methods, this metric is independent of local models and impervious to variations in the attacker's ratio. Therefore, SuDA demonstrates better defensive performance when the proportion of 214 malicious clients is significant. The overall procedure of SuDA coupled with FedAvg is illustrated 215 in Appendix D. For the client, SuDA only modifies its training data and objective function, whereas for the server, SuDA only introduces additional model filtering steps based on the surrogate dataset.
 Therefore, SuDA can be combined with most federated learning algorithms, including FedAvg, Fed Prox (Li et al., 2020) and FedNova (Wang et al., 2020c). Note that malicious attackers may reject
 to follow the proposed protocol. Therefore, we consider several different adaptive attackers and
 empirically prove that SuDA can effectively defend against malicious clients under adaptive attack
 scenarios in Section 5.

222 223

224

4.2 RECONSTRUCT POTENTIAL TRIGGER

The failure of existing methods can be mainly attributed to the tainted metrics used for filtering out malicious models. Specifically, existing methods for defending against backdoor attacks in FL focus on studying the attributes of the received client models themselves, such as taking the mean or median of the model parameters (Yin et al., 2018), filtering out outliers based on squareddistance of updates (Blanchard et al., 2017), and clipping updates with excessive norms (Sun et al., 2019). Consequently, these methods encounter a dilemma when the proportion of malicious clients is significant: these metrics will become unreliable and render poor defense performances.

231 A more direct approach to backdoor defense is to independently evaluate each local model. Drawing 232 inspiration from the insights presented in the centralized setting (Wang et al., 2019), our objective 233 is to reconstruct potential triggers from the model. By perturbing the input pixels, we can manip-234 ulate the model's output for a given input sample. Specifically, consider a model that has been 235 compromised with a trigger targeting a specific label Y_t . Note that the trigger should be reasonably 236 small, otherwise it will be easily detected. For any arbitrary inputs, regardless of their true label Y_i , 237 we can observe that the minimum perturbation required to classify all inputs as the target label is 238 significantly smaller than the perturbation needed to transform the inputs to any non-target label.

Observation 4.1. If there is a trigger with a target label Y_t , then the minimum perturbation required to classify all inputs as the target label should be significantly smaller than the perturbation needed to transform the inputs to any non-target label: $P_{\forall \rightarrow t} < < \min_{i,i \neq t} P_{\forall \rightarrow i}$.

Based on the observation above, we can treat each label as a potential target label and calculate the minimum potential trigger required to misclassify samples of other labels into this target label. Specifically, we represent the process of injecting a trigger into the input x as follows:

246 247

254 255 $x_{poison} = M \cdot \Delta + (1 - M) \cdot x,\tag{3}$

where M is a trigger mask with values ranging from 0 to 1, Δ is a trigger pattern that has the same dimension as the input image. We use the L1 norm of the mask M to measure the size of the trigger. Our goal is to find a trigger that can misclassify clean samples to the target label while being as small as possible. Consequently, we can reconstruct the potential trigger by optimizing the following objective:

$$\min_{M,\Delta} \ell\left(Y_t, f(x_{poison})\right) + \gamma \cdot |M|,\tag{4}$$

where γ is a parameter that balances the misclassification success rate and the size of the trigger. In the process of optimization, we dynamically adjust the parameter γ to gradually achieve a more concise trigger.

Using the optimization objective, we reconstruct the potential trigger for each label. For all potential triggers obtained, we perform *Median Absolute Deviation* outlier detection on their L1 norms. If there is a significantly small outlier, we identify the corresponding label as the attacker's target label and recognize the current model as a malicious model; otherwise, we recognize the current model as a benign model.

In the above method, however, the server needs data to reconstruct potential triggers, which is a
 challenge in the federated setting. In order to protect privacy, the client cannot directly share local
 training data with the server. Therefore, we propose to construct a surrogate dataset that contains no
 private information. The server sends the surrogate dataset to all clients and requires them to train
 local models with both the original local data and the shared surrogate datas:

$$F_k^{cls} := \mathbb{E}_{(x,y)\sim\mathcal{D}_k}\ell(f(x),y) + \mathbb{E}_{(x,y)\sim\mathcal{D}_n}\ell(f(x),y),\tag{5}$$

270 where \mathcal{D}_n is the distribution of the surrogate dataset. Since the surrogate data cannot contain any 271 private information, we propose that the surrogate dataset can be generated using pure Gaussian 272 noise, such as random noise derived from a randomly initialized StyleGAN (Karras et al., 2019). 273 The noise data shares the same range of labels as the real data, with each label denoting a different style of noise. With the surrogate data, the server uniformly selects b samples from the surrogate 274 dataset each round to reconstruct potential triggers. This objective function is formulated to ensure 275 that the local client model performs well on both the real and surrogate data, serving as a crucial 276 reference for the server to detect potential attackers. Given that noise has corresponding labels, we 277 can preliminarily filter out models with excessively low prediction accuracy on the noise data before 278 reconstructing potential triggers. This measure aims to prevent attackers from uploading excessively 279 modified models. 280

281 282

302

309

310

4.3 FEATURE DISTRIBUTION ALIGNMENT

Intuitively, simply adding a surrogate dataset that has a completely different distribution from the original local data will harm the model's generalization performance (Frénay & Verleysen, 2013; Polyzotis et al., 2017). Therefore, inspired by the previous work (Long et al., 2017), which investigates the joint distribution discrepancy, we further introduce feature distribution alignment to enable the models trained on surrogate data to perform well on the natural distribution.

To effectively represent real data using surrogate data, it is also crucial to align the distribution of real features with that of surrogate ones. To facilitate the transfer of knowledge from the surrogate dataset to the real dataset, we draw inspiration from domain adaptation techniques. Specifically, we consider the surrogate dataset as the source domain and the real dataset as the target domain and perform domain adaptation to mitigate the generalization risk of the real distribution. By leveraging the fundamental principles of domain adaptation, we align these two distributions in the feature space and ensure the good performance of the model trained with surrogate data on real data.

Following the previous works on addressing dataset shift (Long et al., 2013; 2017; Lei et al., 2021), we minimize the joint distribution discrepancy between real features and surrogate features. Note that the surrogate dataset is arbitrarily constructed, this allows us to generate appropriate noise data with the same label distribution as the real data. Consequently, we only need to minimize the conditional distribution discrepancy rather than the joint distribution discrepancy. In particular, We propose the objective for the feature distribution alignment:

$$F_k^{da} := \mathbb{E}_y d(\mathcal{D}_k(\phi(x)|y), \mathcal{D}_n(\phi(x)|y)),$$
(6)

where ϕ is the feature extractor that composes the classifier $f = \rho \circ \phi$, $\mathcal{D}_k(\phi(x)|y)$ and $\mathcal{D}_n(\phi(x)|y)$) represent the conditional feature distributions obtained by the feature extractor ϕ on the real dataset and the surrogate dataset, respectively. This objective encourages the feature extractor to learn the same conditional feature distribution from two different data distributions.

Thereafter, we propose the overall objective of clients during local training in the SuDA framework:

$$F_k^{SuDA} = F_k^{cls} + \lambda F_k^{da},\tag{7}$$

where λ is a hyperparameter that governs the trade-off between classification accuracy on training data and the degree of alignment in feature distribution. SuDA enables the model to accurately classify real data and noise data, and simultaneously encourages the model to generate similar features from real data and noise data with the same label, thus achieving good generalization performance on both distributions. Empirical observations in Figure 8 demonstrate this effect.

To theoretically prove the effectiveness of the proposed feature distribution alignment, we analyze the relationship between the model's generalization performance and the distribution shift. Based on the existing theoretical conclusions, the generalization performance is related to the margin between samples and the decision boundary. Therefore, we introduce the definition of *statistical robustness* between two distributions before stating the theorem, serving as a metric for quantifying the degree of generalization performance.

Definition 4.2 (Statistical Robustness). We define statistical robustness for a classifier f on a distribution \mathcal{D} according to a distance metric d: $SR_d(f, \mathcal{D}) = \mathbb{E}_{(x,y)\sim\mathcal{D}} \inf_{f(x')\neq y} d(x', x)$, where classifier $f : \mathcal{X} \to \mathcal{Y}$ predicts class label of an input sample. The defined statistical robustness refers to the expected distance from each sample to the closest adversarial example. Hence, for the model f learned from the source distribution $\mathcal{D}_n(x, y)$, we can quantify the generalization performance on the target distribution $\mathcal{D}(x, y)$. To achieve good generalization performance, we aim to provide a lower bound on the transferred statistical robustness, i.e., $\mathcal{E}_{\mathcal{S}\sim\mathcal{D}_n} SR_d(f,\mathcal{D})$, where $f \leftarrow Sub(\mathcal{S})$ means the model f is trained on the training set \mathcal{S} using $f \leftarrow Sub(\mathcal{S})$

SuDA. To this end, we have the following theorem.

Theorem 4.3. Let f be a neural network, $\mathcal{D}(x, y)$ and $\mathcal{D}_n(x, y)$ are two separable distributions with identical label distributions, corresponding to the distributions of real data and noise data, respectively. Then, training the model with the proposed objective for the feature distribution alignment, i.e., Eq. 7 elicits the bounded statistical robustness.

We provide the proof of Theorem 4.3 in Appendix A.1. Theorem 4.3 shows that the model trained with the proposed objective can learn to provable generalization performance, which is consistent with the previous work (Long et al., 2013; 2017) that aligning the joint distribution between the source and target domains.

339 340

341

354

356

5 EXPERIMENTS

342 The goal of our empirical study is to demonstrate the improved defense capability of SuDA over the 343 state-of-the-art FL defense methods. We conduct our experiments on image classification tasks over 344 three datasets: CIFAR-10(Krizhevsky et al., 2009), FMNIST(Xiao et al., 2017) and SVHN(Netzer et al., 2011). We simulate FL for R rounds among K clients, of which m are corrupted by attackers. 345 In each round, the server uniformly selects $C \cdot K$ clients for some $C \leq 1$ and sends the global model 346 to each selected client. The selected clients then perform local training on the received model for 347 E epochs and send the updates back to the server. The goal of attackers is to make the aggregated 348 model misclassify samples poisoned by triggers into the target class. For the aggregated model 349 on the server, we measure three performance metrics: total accuracy, attack success rate and main 350 accuracy. Total accuracy is computed on the entire test dataset, while attack success rate measures 351 the proportion of test samples with triggers classified as the target label by the model, and main 352 accuracy is computed on clean test samples. Our experimental results show that SuDA significantly 353 outperforms baseline methods in defending against backdoor attacks.

355 5.1 EXPERIMENT SETUP

Datasets and Models. To evaluate the effectiveness of SuDA, we conduct experiments on three computer vision datasets including CIFAR-10, FMNIST and SVHN in the FedML framework(He et al., 2020). We use ResNet-18(He et al., 2016) as the shared global model in FL for all three datasets. We utilize the partition method Latent Dirichlet Sampling(Hsu et al., 2019) to partition datasets, generating a local dataset for each client, and using the parameter α to control the degree of Non-IID. We set $\alpha = 1$ to simulate the Non-IID setting by default and conduct experiments under the IID setting in Appendix E.3.

Surrogate Datasets. At the beginning of the training phase, an un-pretrained StyleGAN-v2 (Karras et al., 2020) is utilized to generate a surrogate dataset without using any training data. The server samples from various Gaussian distributions, each with the same mean but different standard deviations, to generate noise images with diverse latent styles. Each style corresponds to a distinct class. Then the generated noise images are distributed to all clients and used together with the original datasets for local training. The size of the surrogate dataset in our experiments is 2000. We show the surrogate dataset in Appendix E.4.

Random sampling attack. The attack model considered in our work is the random sampling attack
as discussed in Sun et al. (2019), where the attackers have complete control over a fraction of clients.
In each FL round, the server randomly selects a subset of clients to participate in the training process.
The attackers are only able to affect the training of the global model during the rounds in which they
are selected. The number of selected attackers in each round follows a hypergeometric distribution.

Backdoor tasks. The backdoor task aims to make the global model misclassify backdoored samples
 into the target class. Since the server randomly selects clients in each round, multiple attackers may
 be chosen during a single round. We assume that attackers can collude and share the same target,

_											
-	Atk Num	Defense		CIFAR-10			FMNIST			SVHN	
			ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)
		FedAvg	84.72±0.92	3.26 ± 0.85	84.87 ± 0.84 84.00 ± 0.04	91.51 ± 0.09	1.77 ± 0.64	91.64±0.10	89.20 ± 0.02	1.12±0.01	89.28±0.02
		Kram	50.24 ± 1.25	1.44 ± 0.20 3.88±0.91	499 ± 0.94	91.09 ± 0.02 86.69±0.62	1.88 ± 0.28 3.57 ± 0.77	91.83 ± 0.03 86 79 ± 0.57	89.22 ± 0.02 79.65 ± 0.18	1.90±0.04	79.68 ± 0.02
	0	Coomed	83.92 ± 0.84	1.16±0.33	\$3.94±0.86	91.79±0.02	1.95 ± 0.16	91.88±0.02	88.99±0.06	1.19 ± 0.02	89.07±0.05
	0	Normclip	$85.07 {\pm} 0.62$	$1.38 {\pm} 0.16$	$85.18{\pm}0.60$	$91.54{\pm}0.03$	$1.68 {\pm} 0.14$	$91.66 {\pm} 0.04$	$87.56 {\pm} 0.01$	$1.46 {\pm} 0.05$	$87.67 {\pm} 0.01$
		FLAME	83.39±0.94	$1.18{\pm}0.19$	$83.44{\pm}1.07$	$91.04{\pm}0.01$	$1.86{\pm}0.09$	$91.20 {\pm} 0.01$	87.11 ± 0.32	$1.68 {\pm} 0.23$	$87.30 {\pm} 0.29$
		FLTrust	68.30±0.91	2.91 ± 0.78	68.28 ± 0.25	86.58±0.93	1.33 ± 0.62	86.67±0.47	87.90 ± 0.12	1.35 ± 0.06	88.02 ± 0.12
-		SuDA(ours)	85.73±0.31	1.75 ± 0.40	85.78±0.33	90.09±0.38	1.29±0.74	90.17±0.44	90.78±0.70	1.20±0.39	90.88±0.31
		FedAvg	$76.79 {\pm} 0.81$	$87.63 {\pm} 0.77$	83.40 ± 1.00	$83.71 {\pm} 0.03$	$99.78 {\pm} 0.01$	$92.09{\pm}0.04$	81.41 ± 0.42	$67.14 {\pm} 0.81$	86.66 ± 0.13
		RFA	79.38±0.23	76.31±0.58	85.20±0.39	87.41±0.91	22.74 ± 0.50	89.20±0.86	88.99 ± 0.04	5.36 ± 0.48	89.36±0.04
		Coornad	49.58 ± 0.18 76.75 ± 0.13	5.50 ± 1.68	49.41 ± 0.67	84.12 ± 0.11	5.30 ± 0.89 5.82 ± 0.44	84.23 ± 0.98	83.51 ± 0.19	1.26 ± 0.29 3.00 ± 0.20	83.55 ± 0.22
	4	Normelin	70.75 ± 0.13 79.41±0.20	78.44 ± 1.48	80.90±0.39 85 41±0 27	88.93 ± 0.24 88.94 ±0.36	3.83 ± 0.44 8 55±0 31	89.30 ± 0.33 89.58 ± 0.24	87.04 ± 0.10 87.58 ± 0.03	2.09 ± 0.20 2.19 ±0.04	87.25 ± 0.09 87.7 ± 0.03
		FLAME	77.19 ± 1.08	82.65 ± 0.79	83.36±0.41	82.70 ± 0.13	96.68 ± 1.00	90.71 ± 0.09	84.13 ± 0.82	17.29 ± 1.01	85.31±0.79
		FLTrust	$65.54{\pm}0.49$	$38.99 {\pm} 0.33$	$67.57 {\pm} 0.72$	$88.27 {\pm} 1.01$	$4.97 {\pm} 0.17$	$88.82{\pm}1.33$	$86.47 {\pm} 0.06$	$11.34{\pm}0.97$	87.23 ± 0.12
_		SuDA(ours)	$84.65{\pm}0.12$	$2.38{\pm}0.56$	84.77±0.43	90.80±0.17	$1.20{\pm}0.99$	90.92±0.16	90.74±0.60	1.22 ± 0.67	90.76±0.38
		FedAvg	$77.86 {\pm} 0.18$	$88.9 {\pm} 0.22$	$84.67 {\pm} 0.61$	$83.27 {\pm} 0.07$	99.67±0.06	91.61±0.08	$80.79{\pm}0.98$	$66.21 {\pm} 0.60$	$85.9 {\pm} 0.85$
		RFA	$78.69 {\pm} 0.82$	$83.42{\pm}1.01$	85.07 ± 1.11	$84.43 {\pm} 0.55$	$30.17 {\pm} 0.89$	$86.80 {\pm} 0.05$	$87.59 {\pm} 0.02$	4.65 ± 0.26	$87.88 {\pm} 0.03$
		Krum	50.91±1.03	6.44±1.12	50.79±0.90	81.87±0.32	3.60 ± 0.75	82.00±0.38	80.55±0.34	1.36 ± 0.03	80.55±0.38
	8	Normalin	77.21 ± 0.12	81.88±0.93	83.42±0.41	84.68 ± 0.71	$10.1/\pm0.13$	85.49 ± 0.09	86.83 ± 0.48	7.05 ± 0.59	87.25 ± 0.41
		FLAME	76.85 ± 0.31 76.87 ±0.31	85.33 ± 1.27 86 71+1 02	83.38+0.46	80.90 ± 0.48 82.56 ±0.04	99.16 ± 0.04	90 78+0 06	87.34 ± 0.03 80.05 ± 0.23	82.86 ± 0.73	87.09 ± 0.02 86.63 ±0.22
		FLTrust	68.59 ± 1.15	66.75 ± 0.96	72.88 ± 0.94	74.14 ± 1.60	96.33 ± 1.20	81.25 ± 0.55	81.82 ± 0.20	66.16 ± 0.15	86.95 ± 0.22
_		SuDA(ours)	$84.56{\pm}0.65$	$3.22{\pm}0.67$	$84.73{\pm}0.69$	$90.75{\pm}0.10$	$1.31{\pm}0.24$	$90.79{\pm}0.14$	$90.70{\pm}0.40$	$1.16{\pm}0.38$	90.73±0.34
		FedAvg	$77.80{\pm}1.07$	$89.64{\pm}0.64$	$84.67 {\pm} 0.77$	$83.34{\pm}0.17$	$99.64{\pm}0.02$	91.68±0.20	$81.58{\pm}0.14$	$75.29{\pm}0.42$	87.61±0.11
		RFA	77.90 ± 0.69	87.76 ± 0.03	84.61 ± 0.84	$81.88 {\pm} 0.59$	44.36 ± 0.85	85.32 ± 0.93	86.13 ± 0.23	8.12 ± 0.64	86.64 ± 0.16
		Krum	48.03±0.67	8.47±0.96	47.9 ± 0.50	80.02 ± 0.15	10.89 ± 0.49	80.51±0.80	80.74 ± 0.14	2.60 ± 0.05	80.80 ± 0.33
	12	Normalin	77.12 ± 0.79	84.62 ± 0.79 87.01 ± 0.28	85.49±1.05 85.47±0.55	85.31 ± 0.06 85.24 ± 0.41	26.27 ± 0.10 15.08 ± 0.46	85.34 ± 0.80 86.40 ± 0.65	$80.4/\pm0.43$ 86.21 ±0.21	13.37 ± 0.10 12.20 ± 0.20	87.39 ± 0.18 87.18 ± 0.18
		FLAME	76.92 ± 0.23	88.75±0.83	83.60 ± 0.23	82.45 ± 0.27	99.62 ± 0.01	90.69 ± 0.33	79.48 ± 0.34	86.09 ± 1.18	86.32 ± 0.36
		FLTrust	58.65±0.39	78.02±0.69	62.85±0.75	79.04±0.68	98.96±0.44	86.89±0.52	80.90±0.04	78.43±0.18	87.13±0.05
		SuDA(ours)	83.89+0.43	3.71 ± 0.57	84 11+0 59	89 75+0 15	1.44 ± 0.80	8988+033	89 79+0 36	1 63+0 11	89 83+0 78

Table 1: ACC, ASR and MA of defense algorithms on CIFAR-10, FMNIST and SVHN when defending against varying attackers. The poison ratio is 5%.

Table 2: Performance of SuDA under changing attack λ on CIFAR-10.

Atk Num	Metrics		Atk λ											
	methes	0	0.1	0.5	1	2	5	10	100					
4	ACC(%) ASR(%) MA(%)	$\begin{array}{c} 85.21{\pm}0.78\\ 2.58{\pm}0.44\\ 85.29{\pm}0.60\end{array}$	85.26±0.97 2.21±0.79 85.36±0.81	$\begin{array}{c} 85.09{\pm}0.81\\ 2.37{\pm}0.47\\ 85.14{\pm}0.70\end{array}$	$\substack{85.20 \pm 0.55 \\ 2.30 \pm 0.98 \\ 85.31 \pm 0.19}$	84.96±0.71 2.14±0.91 85.14±0.70	$\substack{85.21 \pm 0.22 \\ 2.60 \pm 0.94 \\ 85.36 \pm 0.36}$	85.15±0.98 2.41±0.72 85.40±0.51	$79.75{\pm}0.50\\2.30{\pm}0.34\\79.79{\pm}0.65$					
8	ACC(%) ASR(%) MA(%)	$\begin{array}{c} 84.35{\pm}0.87\\ \textbf{2.01}{\pm}\textbf{0.60}\\ 84.49{\pm}0.94\end{array}$	$\substack{85.05\pm0.68\\2.30\pm0.35\\85.21\pm0.44}$	$\begin{array}{c} 84.91{\pm}0.37\\ 2.20{\pm}0.72\\ 85.01{\pm}0.47\end{array}$	85.25±0.93 2.49±0.62 85.42±0.85	$\substack{85.00 \pm 0.96 \\ 2.26 \pm 0.42 \\ 85.15 \pm 0.84 }$	$\substack{85.05\pm0.74\\2.36\pm0.10\\85.18\pm0.62}$	$\substack{84.72 \pm 0.28 \\ 2.17 \pm 0.76 \\ 84.90 \pm 0.98 }$	$72.77{\pm}0.23 \\ 2.01{\pm}0.38 \\ 72.62{\pm}0.62$					

i.e. all attackers aim to make the global model misclassify backdoored samples into the same target 412 class. For the CIFAR-10 and FMNIST datasets, attackers aim to misclassify into class '2', and for 413 the SVHN dataset, attackers aim to misclassify into class '5'. In each round, attackers implant the 414 trigger into partial samples of each class based on the poison ratio, re-label them with the target 415 class, and then train the local model on the backdoored dataset. The trigger utilized is a white 416 square measuring 4×4 , implanted in the upper-left corner of the poisoned sample. When employing 417 SuDA for defense, the noise dataset will be combined with the backdoored original local dataset for 418 training. Attackers further perform model replacement attacks(Bagdasaryan et al., 2020) to generate 419 malicious local models and send them to the server.

Defense techniques. We conduct FedAvg(McMahan et al., 2017) as the baseline FL aggregation algorithm. The results using FedProx are reported in Appendix E.2. To demonstrate the effectiveness of SuDA in defending against backdoor attacks, we consider six commonly used defense techniques:
(i) Krum (Blanchard et al., 2017); (ii) Coordinate-wise median(Coomed) (Yin et al., 2018); (iii) Norm clipping(Normclip) (Sun et al., 2019); (iv) RFA (Pillutla et al., 2022); (v) FLAME (Nguyen et al., 2022) and (vi) FLTrust (Cao et al., 2020). The detailed hyper-parameters of these algorithms are reported in Appendix C.

427

402

428 5.2 EXPERIMENTAL RESULTS

429

To compare the performance of different defense algorithms, we use three metrics: the average total accuracy (ACC), the average attack success rate (ASR), and the average accuracy of main tasks (MA) in the 5 rounds before the model converges. We conduct FL with a maximum of 200 rounds

Atk Num	Defense	CIFAR-10				FMNIST		SVHN			
	Derense	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	
4	Noise Recon.	75.05±0.34	74.22±0.04	80.33±0.22	85.52±0.18	19.33±0.12	88.44±0.49	88.26±0.77	4.01±0.62	89.40±0.34	
	SuDA(ours)	84.65±0.12	2.38±0.56	84.77±0.43	90.80±0.17	1.20±0.99	90.92±0.16	90.74±0.60	1.22±0.67	90.76±0.38	
8	Noise Recon.	76.29±0.90	75.49±0.68	81.21±0.38	82.06±0.28	24.40±0.82	85.85±0.76	86.70±0.82	4.91±0.48	86.98±0.22	
	SuDA(ours)	84.56±0.65	3.22±0.67	84.73±0.69	90.75±0.10	1.31±0.24	90.79±0.14	90.70±0.40	1.16±0.38	90.73±0.34	

Table 3: Results of reco	nstructing potenti	al triggers without t	the shared surrogate data.





Figure 2: Potential triggers for the target label (second column) and nontarget labels (third column) reconstructed from two different types of poisoned data.



using the adopted defense algorithm on CIFAR-10, FMNIST and SVHN. There are 50 clients in total, the number of backdoor attackers can range from 0 to 12, and the poison ratio can range from 1% to 20%, depending on different settings. In each round, the server randomly selects 20 clients to participate in training and sends them the global model. The selected clients then perform local training for 1 epoch on the received model and send the locally trained model to the server.

Different Numbers of Attackers. As shown in Table 1, SuDA outperforms other baselines in 462 almost all scenarios when defending against varying numbers of attackers across the three datasets. 463 It can be seen that SuDA can improve model performance even without attackers. We conjecture 464 that this is mainly due to several factors: 1) Adding the surrogate dataset reduces data heterogeneity 465 between clients and mitigates client drift; 2) Aligning real feature distribution with the shared noise 466 feature distribution further mitigates client drift; 3) After adding the surrogate dataset, the clients' 467 training data contains more classes, which can alleviate the negative impact caused by the imbalance 468 of sample quantities among different classes. 469

When there are attackers in FL, SuDA's performance is also superior to other baselines. SuDA achieves a significantly lower ASR than other baselines while maintaining high model accuracy. In particular, when the number of attackers is 12, SuDA reduces the ASR by up to 9.45% compared to the second-ranked method, showing the effectiveness of SuDA in defending against a large number of malicious attackers. We notice that the MA of SuDA is sometimes slightly lower than the best result. We speculate that this is mainly because SuDA filters out malicious models before aggregation, reducing the number of models aggregated. Consequently, the aggregated model becomes more difficult to converge, especially in Non-IID settings.

477 **Reconstruct Potential Triggers.** SuDA reconstructs potential triggers by leveraging the shared 478 surrogate data. Figure 3a shows the L1 norm of potential triggers, and the red dots represent 479 potential triggers for the target label. It can be seen that the L1 norm of the potential triggers 480 corresponding to the target label is significantly smaller than that of the triggers corresponding to 481 non-target labels. Accordingly, we perform Median Absolute Deviation outlier detection on the 482 L1 norms and filter out triggers with an excessive anomaly index, which is defined as the absolute 483 deviation of L1 norms divided by MAD. The results in Figure 3b empirically demonstrate that we can distinguish between target and non-target label triggers using anomaly index. We show the 484 reconstructed potential triggers for the target label and non-target labels in the second and third 485 columns of Figure 2, respectively.

442

443

444 445 446

447 448

449

450

451

452

453

454 455 456

457

458

459

460

486 **Impact of Surrogate Data.** To further illustrate the impact of the introduced surrogate dataset, we 487 conducted experiments without using the surrogate data, instead utilizing completely random noise, 488 which does not participate in model training, to reconstruct potential triggers. The experimental 489 results are shown in Table 3. As demonstrated in the table, the performance of using noise data that 490 does not participate in training is worse than that of using the shared surrogate data. This further shows that the alignment operation employed by SuDA enables the surrogate data to represent the 491 real data better, thereby reconstructing the potential triggers more accurately and achieving better 492 defense performance. 493

494 Adaptive Attack Scenario. Note that malicious attackers may reject following the proposed frame-495 work and attempt to circumvent SuDA. To further illustrate SuDA's defensive capabilities, we as-496 sume that the attackers have knowledge about SuDA and adopt more specialized attack methods. Specifically, we consider the following adaptive attack scenario: the attacker changes the alignment 497 parameter λ used to a different one not used by the benign clients. The experimental results are 498 shown in Table 2. It can be seen that SuDA still performs relatively well against this adaptive attack, 499 which demonstrates SuDA's robust defense capabilities against malicious attackers employing vari-500 ous attack methods. We conduct more experiments under adaptive attack scenarios in Appendix E.7. 501

Extra Time Overhead. In our proposed SuDA, reconstructing potential triggers inevitably intro-502 503 duces additional computational overhead. To minimize the extra time overhead while maintaining defensive effects, we record the optimized trigger obtained in each round as the initial value for 504 the next round of optimization, eliminating the need to estimate from scratch. Additionally, we use 505 early-stop in trigger optimization to reduce time consumption. When rounds surpass 60, early-stop 506 significantly reduces the extra time overhead. To further reduce time overhead, we also investigate 507 a more efficient method, SuDA-Efficient. SuDA-Efficient achieves a lower time overhead at the 508 expense of a slight reduction in defense capability. These strategies make the extra time overhead 509 of our method acceptable. To demonstrate the trade-off between SuDA's defensive capability and 510 time overhead, we conduct experiments to compare SuDA and baseline methods on CIFAR-10. The 511 experimental results are shown in Table 14, with more detailed information reports in Appendix E.9. 512

More Experimental Results. In order to further investigate the effectiveness, applicability, and scalability of SuDA, we conduct more ablation experiments. We investigate the defensive performance of SuDA under different poison ratios and the impact of various surrogate data generation methods on the performance of SuDA. We also conduct ablations on the effect of the size of the surrogate dataset and the sensitivity of sample number b of the surrogate dataset. We report these experimental results and more ablation experiments in Appendix E.

518 519 520

521

522

523

524

525

5.3 LIMITATIONS

Although our method achieves significant improvement in the experiment, it also introduces additional communication overhead. To mitigate this overhead, we hope to minimize the noise dataset as much as possible. However, the size of the noise dataset may also affect the performance of the client model and the server's ability to detect backdoor attacks accurately. Therefore, future research should investigate the optimal size of the noise dataset that strikes a balance between communication overhead and model performance.

530

6 CONCLUSION

531 In this paper, we introduce a generated noise dataset that does not contain real data information into 532 the defense against backdoor attacks in FL. These surrogate noise data provide a more direct and ac-533 curate metric for the server to detect malicious models. Through the conditional feature distribution 534 alignment on the noise dataset, our proposed SuDA can effectively filter out malicious models on the server with the assistance of noise data, without affecting the generalization performance of the local 536 model trained by benign clients. Our empirical results demonstrate that SuDA can effectively defend 537 against backdoor attacks and improve the performance of aggregated models, especially when the proportion of malicious clients is significant, providing new insights for defending against attacks in 538 FL. We hope that our work will inspire further research in developing effective defense mechanisms for FL and contribute to the broader goal of securing machine learning systems.

540 ETHIC STATEMENT

This paper does not raise any ethical concerns. This study does not involve any human subjects, practices to data set releases, potentially harmful insights, methodologies and applications, potential conflicts of interest and sponsorship, discrimination/bias/fairness concerns, privacy and security issues, legal compliance, and research integrity issues.

Reproducibility Statement

To make all experiments reproducible, we have listed all detailed hyper-parameters of each FL algorithm. Due to privacy concerns, we will upload the anonymous link of source codes and instructions during the discussion phase to make it only visible to reviewers.

References

546 547

548 549

550

551

552 553

554

556

558

563

564

580

581

- Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How to backdoor federated learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 2938–2948. PMLR, 2020.
- Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79:151–175, 2010.
 - Arjun Nitin Bhagoji, Supriyo Chakraborty, Prateek Mittal, and Seraphin Calo. Analyzing federated learning through an adversarial lens. In *International Conference on Machine Learning*, pp. 634– 643. PMLR, 2019.
- Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector machines. *arXiv preprint arXiv:1206.6389*, 2012.
- Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. Machine learning
 with adversaries: Byzantine tolerant gradient descent. *Advances in neural information processing systems*, 30, 2017.
- Xiaoyu Cao, Minghong Fang, Jia Liu, and Neil Zhenqiang Gong. Fltrust: Byzantine-robust federated learning via trust bootstrapping. *arXiv preprint arXiv:2012.13995*, 2020.
- Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep
 learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.
- Dimitrios Diochnos, Saeed Mahloujifar, and Mohammad Mahmoody. Adversarial risk and robustness: General definitions and implications for the uniform distribution. *Advances in Neural Information Processing Systems*, 31, 2018.
 - Khoa Doan, Yingjie Lao, Weijie Zhao, and Ping Li. Lira: Learnable, imperceptible and robust backdoor attacks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 11966–11976, 2021.
- Benoît Frénay and Michel Verleysen. Classification in the presence of label noise: a survey. *IEEE transactions on neural networks and learning systems*, 25(5):845–869, 2013.
- Clement Fung, Chris JM Yoon, and Ivan Beschastnikh. Mitigating sybils in federated learning poisoning. *arXiv preprint arXiv:1808.04866*, 2018.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François
 Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016.
- Mingming Gong, Kun Zhang, Tongliang Liu, Dacheng Tao, Clark Glymour, and Bernhard
 Schölkopf. Domain adaptation with conditional transferable components. In *International conference on machine learning*, pp. 2839–2848. PMLR, 2016.

609

627

635

636

637

- Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Evaluating backdooring attacks on deep neural networks. *IEEE Access*, 7:47230–47244, 2019.
- 597 Chaoyang He, Songze Li, Jinhyun So, Xiao Zeng, Mi Zhang, Hongyi Wang, Xiaoyang Wang, Pra598 neeth Vepakomma, Abhishek Singh, Hang Qiu, et al. Fedml: A research library and benchmark
 599 for federated machine learning. *arXiv preprint arXiv:2007.13518*, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.
- Folasade Olubusola Isinkaye, Yetunde O Folajimi, and Bolande Adefowoke Ojokoh. Recommendation systems: Principles, methods and evaluation. *Egyptian informatics journal*, 16(3):261–273, 2015.
- Mika Juuti, Sebastian Szyller, Samuel Marchal, and N Asokan. Prada: protecting against dnn model
 stealing attacks. In 2019 IEEE European Symposium on Security and Privacy (EuroS&P), pp. 512–527. IEEE, 2019.
- Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends*® *in Machine Learning*, 14(1–2):1–210, 2021.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pp. 5132–5143. PMLR, 2020.
- Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative
 adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4401–4410, 2019.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8110–8119, 2020.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
 2009.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Qi Lei, Wei Hu, and Jason Lee. Near-optimal linear regression under distribution shift. In *Interna- tional Conference on Machine Learning*, pp. 6164–6174. PMLR, 2021.
 - Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, 2:429–450, 2020.
- Cong Liao, Haoti Zhong, Anna Squicciarini, Sencun Zhu, and David Miller. Backdoor embedding in
 convolutional neural network models via invisible perturbation. *arXiv preprint arXiv:1808.10307*, 2018.
- Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. Trojaning attack on neural networks. In *25th Annual Network And Distributed System Security Symposium (NDSS 2018)*. Internet Soc, 2018.
- Mingsheng Long, Jianmin Wang, Guiguang Ding, Jiaguang Sun, and Philip S Yu. Transfer feature
 learning with joint distribution adaptation. In *Proceedings of the IEEE international conference* on computer vision, pp. 2200–2207, 2013.

648 649 650	Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Deep transfer learning with joint adaptation networks. In <i>International conference on machine learning</i> , pp. 2208–2217. PMLR, 2017.
651 652 653 654	Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In <i>Artificial intelligence and statistics</i> , pp. 1273–1282. PMLR, 2017.
655 656 657	Omar Montasser, Steve Hanneke, and Nathan Srebro. Vc classes are adversarially robustly learnable, but only improperly. In <i>Conference on Learning Theory</i> , pp. 2512–2530. PMLR, 2019.
658 659	Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
660 661 662	Thien Duc Nguyen, Phillip Rieger, Huili Chen, Hossein Yalame, Helen Möllering, Hossein Fer- eidooni, Samuel Marchal, Markus Miettinen, Azalia Mirhoseini, Shaza Zeitouni, et al. Flame: Taming backdoors in federated learning. In <i>USENIX Security Symposium</i> , 2022.
664 665	Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. <i>IEEE Transactions on knowledge and data engineering</i> , 22(10):1345–1359, 2010.
666 667	Sinno Jialin Pan, Ivor W Tsang, James T Kwok, and Qiang Yang. Domain adaptation via transfer component analysis. <i>IEEE transactions on neural networks</i> , 22(2):199–210, 2010.
669 670 671	Ashwinee Panda, Saeed Mahloujifar, Arjun Nitin Bhagoji, Supriyo Chakraborty, and Prateek Mit- tal. Sparsefed: Mitigating model poisoning attacks in federated learning with sparsification. In <i>International Conference on Artificial Intelligence and Statistics</i> , pp. 7587–7624. PMLR, 2022.
672 673 674	Krishna Pillutla, Sham M Kakade, and Zaid Harchaoui. Robust aggregation for federated learning. <i>IEEE Transactions on Signal Processing</i> , 70:1142–1154, 2022.
675 676 677	Neoklis Polyzotis, Sudip Roy, Steven Euijong Whang, and Martin Zinkevich. Data management challenges in production machine learning. In <i>Proceedings of the 2017 ACM International Con-</i> <i>ference on Management of Data</i> , pp. 1723–1726, 2017.
678 679 680	Joaquin Quinonero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. <i>Dataset shift in machine learning</i> . Mit Press, 2008.
681 682 683	Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Adversarially robust generalization requires more data. <i>Advances in neural information processing systems</i> , 31, 2018.
684 685 686 687	Virat Shejwalkar, Amir Houmansadr, Peter Kairouz, and Daniel Ramage. Back to the drawing board: A critical evaluation of poisoning attacks on production federated learning. In 2022 IEEE Symposium on Security and Privacy (SP), pp. 1354–1371. IEEE, 2022.
688 689 690	Neta Shoham, Tomer Avidor, Aviv Keren, Nadav Israel, Daniel Benditkis, Liron Mor-Yosef, and Itai Zeitak. Overcoming forgetting in federated learning on non-iid data. <i>arXiv preprint arXiv:1910.07796</i> , 2019.
691 692 693	Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. <i>arXiv preprint arXiv:1409.1556</i> , 2014.
694 695	Ziteng Sun, Peter Kairouz, Ananda Theertha Suresh, and H Brendan McMahan. Can you really backdoor federated learning? <i>arXiv preprint arXiv:1911.07963</i> , 2019.
696 697 698 699	Zhenheng Tang, Yonggang Zhang, Shaohuai Shi, Xin He, Bo Han, and Xiaowen Chu. Virtual ho- mogeneity learning: Defending against data heterogeneity in federated learning. In <i>International</i> <i>Conference on Machine Learning</i> , pp. 21111–21132. PMLR, 2022.
700 701	Florian Tramèr, Fan Zhang, Ari Juels, Michael K Reiter, and Thomas Ristenpart. Stealing machine learning models via prediction apis. In <i>USENIX security symposium</i> , volume 16, pp. 601–618, 2016.

- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. Journal of machine learning research, 9(11), 2008.
- Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y
 Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In 2019
 IEEE Symposium on Security and Privacy (SP), pp. 707–723. IEEE, 2019.
- Hongyi Wang, Kartik Sreenivasan, Shashank Rajput, Harit Vishwakarma, Saurabh Agarwal, Jy-yong Sohn, Kangwook Lee, and Dimitris Papailiopoulos. Attack of the tails: Yes, you really can backdoor federated learning. *Advances in Neural Information Processing Systems*, 33:16070–16084, 2020a.
- Hongyi Wang, Mikhail Yurochkin, Yuekai Sun, Dimitris Papailiopoulos, and Yasaman Khazaeni.
 Federated learning with matched averaging. *arXiv preprint arXiv:2002.06440*, 2020b.
- Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. *Advances in neural information processing systems*, 33:7611–7623, 2020c.
- Chenwang Wu, Defu Lian, Yong Ge, Zhihao Zhu, and Enhong Chen. Triple adversarial learning for influence based poisoning attack in recommender systems. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 1830–1840, 2021.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmark ing machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Chulin Xie, Keli Huang, Pin-Yu Chen, and Bo Li. Dba: Distributed backdoor attacks against federated learning. In *International conference on learning representations*, 2020.
- Jie Xu, Benjamin S Glicksberg, Chang Su, Peter Walker, Jiang Bian, and Fei Wang. Federated learning for healthcare informatics. *Journal of Healthcare Informatics Research*, 5:1–19, 2021.
- Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. Byzantine-robust distributed learning: Towards optimal statistical rates. In *International Conference on Machine Learning*, pp. 5650–5659. PMLR, 2018.
- Binhang Yuan, Song Ge, and Wenhui Xing. A federated learning framework for healthcare iot devices. *arXiv preprint arXiv:2005.05083*, 2020.
- Kun Zhang, Bernhard Schölkopf, Krikamol Muandet, and Zhikun Wang. Domain adaptation under target and conditional shift. In *International conference on machine learning*, pp. 819–827. PMLR, 2013.
- Kun Zhang, Mingming Gong, Petar Stojanov, Biwei Huang, Qingsong Liu, and Clark Glymour. Do main adaptation as a problem of inference on graphical models. *Advances in neural information processing systems*, 33:4965–4976, 2020.
- Zhengming Zhang, Ashwinee Panda, Linyue Song, Yaoqing Yang, Michael Mahoney, Prateek Mittal, Ramchandran Kannan, and Joseph Gonzalez. Neurotoxin: Durable backdoors in federated
 learning. In *International Conference on Machine Learning*, pp. 26429–26446. PMLR, 2022.
- Han Zhao, Remi Tachet Des Combes, Kun Zhang, and Geoffrey Gordon. On learning invariant representations for domain adaptation. In *International conference on machine learning*, pp. 7523–7532. PMLR, 2019.
- Yuqing Zhu, Xiang Yu, Manmohan Chandraker, and Yu-Xiang Wang. Private-knn: Practical differential privacy for computer vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11854–11862, 2020.
- 752

- 753
- 754
- 755

A Proof

A.1 PROOF FOR THEOREM 4.3

Proof. We decompose the statistical robustness $SR_d(f, \mathcal{D}(x, y))$ to three quantities as follows:

$$SR_d(f,\mathcal{D}) = (SR_d(f,\mathcal{D}) - SR_d(f,\mathcal{D}_n)) + (SR_d(f,\mathcal{D}_n) - SR_d(f,\tilde{\mathcal{D}}_n)) + SR_d(f,\tilde{\mathcal{D}}_n), \quad (8)$$

where \hat{D}_n denotes the empirical distribution for the training set sampled from the noise data distribution D_n . Then based on the linearity of expectation and triangle inequality, we can bound the transferred statistical robustness as follows:

$$\mathbb{E}_{f \leftarrow \mathcal{D}} SR_d(f, \mathcal{D}) \ge \mathbb{E}_{f \leftarrow \mathcal{D}} SR_d(f, \tilde{\mathcal{D}}_n) - |\mathbb{E}_{f \leftarrow \mathcal{D}} [SR_d(f, \mathcal{D}_n) - SR_d(f, \tilde{\mathcal{D}}_n)]| - |\mathbb{E}_{f \leftarrow \mathcal{D}} [SR_d(f, \mathcal{D}) - SR_d(f, \mathcal{D}_n)]|,$$
(9)

where $\mathbb{E}_{f \leftarrow D}$ denotes $\mathbb{E}_{\substack{S \sim D \\ f \leftarrow Nog(S)}}$ for brevity. The three terms above represent the empirical robustness, the generalization penalty and the distribution shift penalty, respectively. Since our goal is to bound the transferred statistical robustness, we need to bound both the generalization penalty and the distribution shift penalty. There are already multiple works (Diochnos et al., 2018; Schmidt et al., 2018; Montasser et al., 2019) have studied the bound of the generalization penalty. In order to bound the distribution shift penalty, we introduce the following lemma:

Lemma A.1. Let \mathcal{D} and \mathcal{D}_n be two distributions with identical label distributions, $d(\cdot, \cdot)$ be the Wasserstein distance of two distributions. Then for any classifier f, we have:

$$|SR_d(f, \mathcal{D}) - SR_d(f, \mathcal{D}_n)| \le \mathbb{E}_y d(\mathcal{D}|y, \mathcal{D}_n|y).$$
(10)

We prove Lemma A.1 in the follow-up section, i.e., Appendix A.2. With Eq. 9 and Lemma A.1, we can further bound the transferred statistical robustness as follows:

$$\mathbb{E}_{f \leftarrow \mathcal{D}} SR_d(f, \mathcal{D}) \ge \mathbb{E}_{f \leftarrow \mathcal{D}} SR_d(f, \tilde{\mathcal{D}}_n) - |\mathbb{E}_{f \leftarrow \mathcal{D}} [SR_d(f, \mathcal{D}_n) - SR_d(f, \tilde{\mathcal{D}}_n)]| - \mathbb{E}_y d(\mathcal{D}|y, \mathcal{D}_n|y).$$
(11)

The last term $\mathbb{E}_y d(\mathcal{D}|y, \mathcal{D}_n|y)$ is bounded by the proposed objective, i.e., Eq. 7. Thus, the transferred statistical robustness is bounded and the proof is complete.

A.2 PROOF FOR LEMMA A.1

Proof. To begin with, since the distance metric $d(\cdot, \cdot)$ is the Wasserstein distance, we have:

$$d(\mathcal{D}|y, \mathcal{D}_n|y) = \inf_{J \in \mathcal{J}(\mathcal{D}|y, \mathcal{D}_n|y)} \mathbb{E}_{(x, x') \sim J} m(x, x'),$$
(12)

where $\mathcal{J}(\mathcal{D}|y, \mathcal{D}_n|y)$ is the set of joint distributions. Let \mathcal{J}^* be the optimal transport between $\mathcal{D}|y$ and $\mathcal{D}_n|y$. Then we have:

812

$$SR_d(f, \mathcal{D}(x, y)) = \mathbb{E}_{(x,y)\sim\mathcal{D}} \inf_{x \in \mathcal{D}} m($$

$$\mathcal{D}(x,y)) = \mathbb{E}_{(x,y)\sim\mathcal{D}} \inf_{f(x')\neq y} m(x',x)$$
$$= \mathbb{E}_y \mathbb{E}_{x\sim\mathcal{D}|y} \inf_{f(x')\neq y} m(x',x)$$
$$= \mathbb{E}_y \mathbb{E}_{(x,x'')\sim\mathcal{J}^*} \inf_{f(x')\neq y} m(x',x)$$

816

$$= \mathbb{E}_{y} \mathbb{E}_{(x,x'')\sim\mathcal{J}^{*}} \inf_{f(x')\neq y} m(x',x)$$

$$\leq \mathbb{E}_{y} \mathbb{E}_{(x,x'')\sim\mathcal{J}^{*}} \inf_{f(x')\neq y} [m(x',x'') + m(x'',x)]$$

$$= \mathbb{E}_{y} \mathbb{E}_{x'' \sim \mathcal{D}_{n}|y} \inf_{f(x') \neq y} m(x'', x') + \mathbb{E}_{y} \mathbb{E}_{(x, x'') \sim \mathcal{J}^{*}} m(x'', x)$$

833

834 835 836

837

 $= \mathbb{E}_{(x'',y)\sim\mathcal{D}_n} \inf_{f(x')\neq y} m\left(x'',x'\right) + \mathbb{E}_y d\left(\mathcal{D}|y,\mathcal{D}_n|y\right)$ $= SR_d\left(f, \mathcal{D}_n(x, y)\right) + \mathbb{E}_y d\left(\mathcal{D}|y, \mathcal{D}_n|y\right).$

Similarly, we can also prove that:

$$SR_d(f, \mathcal{D}_n(x, y)) \le SR_d(f, \mathcal{D}(x, y)) + \mathbb{E}_y d(\mathcal{D}|y, \mathcal{D}_n|y).$$
(14)

Now using Eq. 13 and 14 we have:

$$-\mathbb{E}_{y}d\left(\mathcal{D}|y,\mathcal{D}_{n}|y\right) \leq SR_{d}(f,\mathcal{D}) - SR_{d}\left(f,\mathcal{D}_{n}\right) \leq \mathbb{E}_{y}d\left(\mathcal{D}|y,\mathcal{D}_{n}|y\right).$$
(15)

Thus, we complete the proof.

(13)

В MORE RELATED WORK

838 Federated Learning. FL is first proposed by McMahan et al. (2017) to protect data privacy in 839 distributed machine learning. Training models within the FL framework can effectively safeguard 840 privacy, as local data need not be shared. Instead of aggregating local data, the server aggregates 841 local model updates from selected clients to update the global model in each round. To address spe-842 cific problems within FL, various optimization algorithms have been proposed. FedCurv (Shoham 843 et al., 2019) tackles the catastrophic forgetting problem of FL in the Non-IID case by drawing an 844 analogy with lifelong learning. FedMA (Wang et al., 2020b) reduces the overall communication burden by constructing the global model in a layer-wise manner, matching and averaging hidden el-845 ements. There are also many algorithms proposed to address the issue of client drift (Li et al., 2020; 846 Wang et al., 2020c; Karimireddy et al., 2020; Tang et al., 2022), such as FedNova (Wang et al., 847 2020c), which utilizes normalized averaging to eliminate objective inconsistency. VHL (Tang et al., 848 2022) also introduces surrogate data into FL, but they focus on solving data heterogeneity issues, 849 while we focus on addressing backdoor attacks. 850

Backdoor Attack on Federated Learning. The goal of backdoor attacks is to modify the global 851 model so that it can produce the desired target labels for inputs that possess specific properties (She-852 jwalkar et al., 2022). Bagdasaryan et al. (2020) investigates semantic backdoor attacks where the 853 global model misclassifies input samples with the same semantic property, e.g. misclassifies the 854 blue truck as a bird, and proposes a model-replacement attack to replace the global model. Bhagoji 855 et al. (2019) discusses model poisoning attacks launched by a single malicious client. They boost 856 the malicious updates to overcome the impact of updates from benign clients, and further propose 857 alternating minimization and estimating benign updates to evade detection in almost every round. 858 Wang et al. (2020a) proposes a new category of backdoor attacks called edge-case backdoors, and 859 explains how these edge-case backdoors can lead to detection failures. Zhang et al. (2022) inserts 860 more durable backdoors into FL systems by attacking parameters that are changed less in magni-861 tude during training. Different from these works that only consider the centralized backdoor attack on FL, Xie et al. (2020) investigates the distributed backdoor attack (DBA), which decomposes a 862 global trigger pattern into separate local patterns and embeds them into the training set of different 863 adversarial parties respectively.

864 **Robust Federated Learning.** The goal of robust federated learning is to mitigate the impact of specific attacks during training. Blanchard et al. (2017) select model update(s) with the minimum 866 squared distance to the updates of other clients. Coordinate-wise median (Yin et al., 2018) selects 867 the median element coordinate-wise among the model updates of clients. Norm clipping (Sun et al., 868 2019) clips model updates whose norm exceeds a specific threshold. RFA (Pillutla et al., 2022) replaces the weighted arithmetic mean in FedAvg with a weighted geometric median, which is computed using the smoothed Weiszfeld's algorithm. FLTrust (Cao et al., 2020) mitigates the impact 870 of backdoors by training models on the server-side with the additional root dataset, performing 871 similarity checks based on the trained model and received local models. FLAME (Nguyen et al., 872 2022) eliminates backdoors by injecting noise into the model and employs HDBSCAN clustering 873 and model weight clipping to reduce the required noise. FoolsGold (Fung et al., 2018) sums up 874 the historical update vectors and calculates the cosine similarity between all participants to assign a 875 global learning rate to each party. By giving lower learning rates to similar update vectors, Fools-876 Gold defends against label flipping and centralized backdoor attacks. SparseFed (Panda et al., 2022) 877 utilizes global model top-k sparse updates and client-level gradient clipping to mitigate the impact of 878 poisoning attacks. Our evaluation includes comparisons to six commonly used defense algorithms 879 and demonstrates the stronger capabilities of SuDA against backdoor attacks.

880 881 882

889

C IMPLEMENTATION DETAILS

$$L_{atk}^t = \gamma (X - G^t) + G^t,$$

where γ is the scaling factor for the balance between attack capability and stealthiness. The scale factor we used is $\frac{m}{n_a}$, where m is the number of clients participating in aggregation each round and n_a is the number of attackers, which is consistent with previous outstanding works (Bagdasaryan et al., 2020; Wang et al., 2020a).

Krum and Multi-Krum (Blanchard et al., 2017): Given n clients, Krum aims to defend against a 895 maximum of f attackers. In each round r, the server receives n updates (V_1^r, \dots, V_n^r) . For each update V_i^r , we denote $i \to j$ as the set of n - f - 2 closest updates to V_i^r . Then the score for 896 897 each client i is defined as the sum of squared distances between V_i and each update V_j in the set 898 $i \rightarrow j$: $score(i) = \sum_{i \rightarrow j} ||V_i - V_j||^2$. Krum then selects $V_{krum} = V_{i_*}$ with the lowest score 899 $score(i_*) \leq score(i)$ for all i, and updates the global model as $w^{r+1} = w^r - V_{krum}$. While 900 Multi-Krum selects $m \in \{1, \dots, n\}$ updates V_1^*, \dots, V_m^* with the lowest scores, and calculates 901 their average $\frac{1}{m}\sum_{i}V_{i}^{*}$ to replace V_{krum} . In our experiments, we apply f = 6 for both Krum and 902 Multi-Krum and set m = 8 for Multi-Krum. 903

Coordinate-wise median (Yin et al., 2018): Given the set of updates (V_1^r, \dots, V_n^r) in each round, Coomed aggregates the updates: $\overline{V}^r = \text{Coomed}\{V_i^r : i \in [n]\}$, where the j^{th} coordinate of \overline{V}^r is given by $\overline{V}^r(j) = \text{med}\{V_i^r(j) : i \in [n]\}$. Here, the function med represents the 1-dimensional median, and $[n] = \{1, \dots, n\}$.

Norm clipping (Sun et al., 2019): Due to the assumption that adversarial attacks can potentially generate updates with large norms, Normclip simply clips model updates whose norm exceeds a specific threshold M:

911 912 $w_k^r = \frac{w_k^r}{max(1, \|w_k^r\|_2/M)}.$

In our experiments, we set the threshold M = 200.

RFA (Pillutla et al., 2022): RFA replaces the weighted arithmetic mean utilized in FedAvg with a weighted geometric median:

$$\underset{v}{\arg\min} \sum_{i} \alpha_{i} \|v - w_{i}\|$$

which is computed using the *smoothed Weiszfeld's algorithm*. The weight α_i is set to the proportion of training samples in the client $\alpha_i = \frac{n_i}{\sum_{j \in S^r} n_j}$, where S^r is the subset of selected clients at round *r*. For iteration budget *R* and the parameter ν in the smoothed Weiszfeld's algorithm, we set R = 4and $\nu = 10^{-5}$.

FLAME (Nguyen et al., 2022): FLAME eliminates backdoors by injecting noise into the model and employs HDBSCAN clustering and model weight clipping to reduce the required noise. In our experiments, we set *min cluster size* to $(C \cdot K)/2 + 1$ and *min samples* to 1, which is consistent with the original paper.

FLTrust (Nguyen et al., 2022): FLTrust mitigates the impact of backdoors by training models on
 the server-side with the root dataset, performing similarity checks based on the trained model and
 received local models. We randomly collected 100 samples from the test set as the root dataset.

930 SuDA: SuDA generates the surrogate noise dataset using an un-pretrained StyleGAN-v2 (Karras 931 et al., 2020). Clients then proceed to train local models with the SuDA objective parameter λ set 932 to 1, and the batch size is set to 128 for both real data and noise data. Then the server receives 933 local models and reconstructs potential triggers. The sample number b is set to 128. During the trigger optimization process, we gradually increase λ to obtain as concise a trigger as possible while 934 ensuring that the misclassification accuracy of the first term in Eq. 4 is greater than 98%. We set the 935 number of noise samples b used for reconstructing potential triggers to a fixed value of 128 for all 936 experiments. 937

For all experiments, the learning rates are set to 0.01 and the learning rate decay is set to 0.992 per round. We employ momentum-SGD as optimizers, with momentum of 0.9 and weight decay of 0.0001. The degree of Non-IID local data distribution on the client is set to $\alpha = 1$. Our experiments were conducted on Ubuntu 20.04 LTS, Intel(R) Xeon(R) Platinum 8255C CPU, and 3090 GPU.

942 943

D ALGORITHM

944

In SuDA, the server generates the surrogate noise dataset at the beginning of the training phase and distributes the noise data to all clients. The clients proceed to train their respective local models using both the original dataset and the noise dataset, and send the trained local models back to the server. To effectively identify and mitigate backdoor attackers, the server reconstructs potential triggers for each local model, leveraging the presence of the noise data. We summarize the overall training procedure of SuDA in Algorithm 1.

951 In previous methods, malicious models are also involved in the calculation of these metrics, thus 952 may yield tainted metrics and fail to achieve effective defense. For example, Krum (Blanchard 953 et al., 2017) calculates the sum of the squared distances between each client model and the other client models as its score, and aggregates several models with the lowest scores. However, when 954 the majority are attackers, the score of the malicious model may be relatively lower, leading the 955 server to aggregate malicious models and resulting in defense failure. The proposed method aims to 956 defend against backdoor attacks by designing a metric that will not be tainted by malicious models. 957 To this end, we propose a metric that performs individual evaluation for each local model using 958 surrogate data. This individual evaluation approach renders the metric impervious to variations in 959 the attacker's ratio.

960 961 962

E MORE EXPERIMENTAL RESULTS

963 964 E.1 DIFFERENT DATASETS AND RATIO

As shown in Tables 4 and 6, we conduct experiments on 3 datasets with different poison ratios, ranging from 1% to 20%. The number of attackers is 4. It can be seen that SuDA can effectively defend against backdoor attacks under different poison ratios. Note that although the accuracy of Normclip is sometimes slightly higher than SuDA, its Atk Rate is much higher in comparison. This is mainly because Normclip aggregates all clipped local models, which helps with model convergence but does not completely eliminate the negative impact caused by attackers. On the other hand, SuDA directly filters out malicious models by leveraging the surrogate dataset and does not select them in the aggregation process. Although this leads to a decrease in the number of clients participating

Al	gorithm 1 Surrogate Data-guided Aggregation Strategy (SuDA)
Inj	put: local epochs E, client number K, maximum round R, initial parameter w^0
Ou	tput: global parameter w
]	initialization: Server generates the surrogate noise dataset D , and distributes the initial model
1	v^0 and D to all clients.
5	Server:
f	for each round $r \in \{0, 1, \cdots, R\}$ do
	Uniformly selects a subset of clients $S^r \subseteq \{1, \dots, K\}$
	Sends the global model w^r to all selected clients $k \in S^r$
	for each client $k \in S^r$ in parallel do
	$w_k^r \leftarrow \text{ClientTraining}(\hat{k}, w^r)$
	end for
	$\mathcal{W}^r \leftarrow \{w^r_k k \in \mathcal{S}^r\}$
	//Accuracy test
	$\mathcal{W}_1^r \leftarrow \operatorname{Acc}\bar{\operatorname{Filter}}(\mathcal{W}^r, \tilde{D})$
	//Potential Trigger Construction
	$\mathcal{T}^r \leftarrow \operatorname{TriggerConstr}(\mathcal{W}_1^r, \tilde{D})$
	//Outlier Detection
	$\mathcal{W}_2^r \leftarrow \mathrm{MAD}(\mathcal{T}^r)$
	//Aggregation
	$w^{r+1} \leftarrow \sum p_k w_k^r, w_k^r \in \mathcal{W}_2^r$
(end for
]	Benign Client:
1	for each epoch $e \in \{0, \cdots, E-1\}$ do
	$w_{k,e+1}^r \leftarrow w_{k,e}^r - \eta_{k,e} \nabla F_k^{SuDA} \left(w_{k,e}^r \right)$
	and for
1	Return w^r to sever
(Compromised Client:
1	niects the backdoor into the local dataset
f	For each epoch $e \in \{0, \dots, E-1\}$ do
-	$\sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{i$
	$w_{k,e+1} \leftarrow w_{k,e} - \eta_{k,e} \vee r_k (w_{k,e})$
(end for
1	$w_{atk}^r \leftarrow \gamma(w_k^r - w^r) + w^r$
]	Return w_{atk}^r to sever

1012 1013

1015

in server aggregation, making it more challenging to converge, SuDA still achieves a significantly lower Atk Rate than Normclip while maintaining comparable Acc and Main Acc.

1014 E.2 EXPERIMENTS WITH FEDPROX

FedProx (Li et al., 2020) is one of the more common training methods than FedAvg when extreme heterogeneity exists in the client data. Therefore, we conduct experiments with FedProx on CIFAR-10 in the Non-IID setting. We set the Non-IID degree control parameter $\alpha = 1$ and the poison ratio is 5%. The experimental results are shown in Table 5. It can be seen that SuDA is easy to integrate with FedProx and performs well against backdoor attacks.

1021

1022 E.3 PERFORMANCE IN THE IID SETTING

1023

To further demonstrate the applicability of SuDA, we compared the performance of defense methods on 3 datasets in the IID setting. We set the Non-IID degree control parameter $\alpha = 100$ to simulate the IID setting. The experimental results are shown in Tables 7. It can be seen that SuDA

1051 1052

1053

1061 1062

1064

Poison Ratio	Defense		4 attackers			8 attackers		1	2 attackers	
i olsoli Ratio	Detense	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)
	FedAvg	81.12	2.41	81.21	84.43	2.41	84.51	84.97	2.76	85.07
	RFA	84.62	1.55	84.72	85.37	1.77	85.47	85.13	1.79	85.22
	Krum	49.58	5.5	49.41	51.53	1.86	51.16	52.09	3.88	51.87
1%	MultiKrum	80.5	3.24	80.67	80.49	2.3	80.66	80.49	1.53	80.5
1,0	Coomed	83.76	1.42	83.73	83.51	1.29	83.49	83.49	1.42	83.57
	Normclip	85.29	1.53	85.26	85.41	1.46	85.45	85.3	1.77	85.41
	SuDA(ours)	85.59	1.2	85.71	84.89	1.27	85.03	85.14	1.35	85.18
	FedAvg	76.79	87.63	83.4	77.86	88.9	84.67	77.8	89.64	84.67
	RFA	79.38	76.31	85.2	78.69	83.42	85.07	77.9	87.76	84.61
	Krum	49.58	5.5	49.41	50.91	6.44	50.79	48.03	8.47	47.9
5%	MultiKrum	80.5	3.24	80.67	78.73	6.28	79	76.99	10.98	77.72
570	Coomed	76.75	58.44	80.9	77.27	81.88	83.42	77.12	84.62	83.49
	Normclip	79.41	78.17	85.41	78.85	83.33	85.26	78.68	87.91	85.47
	SuDA(ours)	85.2	2.3	85.31	85.25	2.49	85.42	84.68	2.87	84.9
	FedAvg	76.62	92.21	83.61	78.09	87.93	84.94	77.88	91.84	84.94
	RFA	78.44	84.71	84.95	78.09	88.96	84.94	78.18	91.31	85.22
	Krum	49.58	5.5	49.41	51.89	3.83	51.55	51.64	2.03	51.28
10%	MultiKrum	80.5	3.24	80.67	79.18	3	79.26	74.45	50.3	77.79
10%	Coomed	77.51	81.64	83.67	76.74	86.9	83.28	77.11	88.99	83.85
	Normclip	79.04	86.05	85.7	78.58	87.52	85.35	78.55	90.61	85.55
	SuDA(ours)	85.38	2.06	85.52	85.07	2.03	85.2	84.79	2.03	84.87
	FedAvg	73.65	91.92	80.34	78.21	90.72	85.23	76.86	92.58	83.89
	RFA	76.47	87.71	83.08	77.83	85.65	84.39	75.62	85.46	82.01
	Krum	49.58	5.5	49.41	52.07	2.3	51.63	46.69	3.2	46.5
20%	MultiKrum	80.5	3.24	80.67	79.27	1.6	79.36	75.14	68.37	79.98
2070	Coomed	77.03	86.31	83.54	76.65	88.77	83.31	76.87	89.84	83.68
	Normclip	78.62	88.81	85.49	78.7	89.05	85.62	78.32	91	85.37
	SuDA(ours)	85.53	2.01	85.71	84.92	1.51	84.97	84.28	2.01	84.42

1026 Table 4: ACC, ASR and MA of defense algorithms on the dataset CIFAR-10 with different poison 1027 ratios.

Table 5: Performance of different defense algorithms with FedProx on CIFAR-10.

54	Defense	Atk Num = 0			Atk Num $= 4$			А	tk Num = 8	3	A	tk Num = 1	2	
55	Derense	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	
2	FedAvg	84.92	1.97	85.03	75.47	89.56	82.13	77.61	88.9	84.39	78.34	88.9	85.21	
)	RFA	84.27	2.23	84.4	79.46	77.67	85.36	78.71	83.26	85.08	77.78	88.07	84.5	
	Krum	54.54	2.82	54.25	54.64	4.34	54.51	53.53	6.44	53.45	53.41	8.17	53.38	
	MultiKrum	80.71	1.73	80.84	79.58	2.03	79.66	78.87	2.22	79.04	78.21	28.44	80.23	
	Coomed	84.04	1.44	84.06	78.14	65.96	82.96	77.38	81.71	83.49	77.18	85.41	83.61	
	Normclip	84.71	1.42	84.72	78.71	77.19	84.55	78.29	83.35	84.65	78.3	86.82	84.96	
	SuDA(ours)	86.08	1.53	86.17	85.42	1.82	85.51	84.48	1.99	84.61	84.67	2.1	84.79	

can still effectively defend against backdoor attacks in the IID setting and preserve high accuracy simultaneously.

SURROGATE DATASET GENERATION E.4 1066

1067 To further investigate the effect of different surrogate datasets, we employ two additional generated 1068 datasets to replace the original surrogate dataset produced by StyleGAN. These two datasets include 1069 one generated by a simple CNN and another generated by upsampling pure Gaussian noise. In our 1070 data generation methods, we sample noise from various Gaussian distributions, each with the same 1071 mean but different standard deviations, to generate noise with diverse latent styles that correspond to distinct classes. Given that datasets CIFAR-10, FMNIST and SVHN each consist of 10 classes, 1072 the surrogate dataset also comprises 10 classes. The size of the surrogate dataset is 2000 in our 1073 experiments, which is a small proportion of the utilized datasets, i.e., 3.33% for CIFAR-10, 2.86% 1074 for FMNIST, and 0.33% for SVHN. We show the generated surrogate datasets as Figures 4, 5 and 1075 6. We also conduct ablations on the sensitivity of the size of the surrogate dataset and report results 1076 in Table 10. 1077

For the dataset generated by the simple CNN, we first sample 64-dimensional noises. These noises 1078 are then fed into a CNN composed of 4 transpose convolutional layers and 3 convolutional layers. 1079 The CNN model processes the input and produces noise data of size 32×32 . We employ 10 CNNs

Poison Ratio	Defense		CIFAR-10			FMNIST			SVHN	
		ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)
	FedAvg	81.12	2.41	81.21	90.5	8.57	91.21	87.66	2.92	87.79
	KFA	84.62 49.58	5.5	84.72 49.41	91.49 84.12	5.3	84.23	89.25 83.51	1.25	89.52 83.55
1%	MultiKrum	80.5	3.24	80.67	91.23	1.99	91.36	86.47	1.37	86.53
1 /0	Coomed	83.76	1.42	83.73	92 01 (7	1.77	92.17	89.16	1.28	89.23
	Normclip SuDA(ours)	85.29 85.59	1.53 1.2	85.26 85.71	91.67 91.36	1.6 1.6	91.83 91.43	87.74 90.45	1.46 0.99	87.82 90.52
	FedAvg	76.79	87.63	83.4	83.71	99.78	92.09	81.41	67.14	86.66
	RFA	79.38	76.31	85.2	87.41	22.74	89.2	88.99	5.36	89.36
500	MultiKrum	80.5	3.24	49.41 80.67	91.23	1.99	84.23 91.36	86.47	1.20	85.55
5%	Coomed	76.75	58.44	80.9	88.93	5.83	89.36	89.04	3.09	89.25
	Normclip SuDA(ours)	79.41 85.2	78.17 2.3	85.41 85.31	88.94 91.46	8.55 1.44	89.58 91.55	87.58 90.72	2.19 1	87.7 90.75
	FedAvg	76.62	92.21	83.61	83.46	99.75	91.82	79.6	89.29	86.72
	RFA	78.44	84.71	84.95	85.4	28.68	87.63	87.34	10.51	88.05
	Krum MultiKrum	49.58	5.5	49.41	84.12	5.3	84.23	83.51	1.26	83.55 86.53
10%	Coomed	77.51	81.64	83.67	87.42	11.29	88.27	87.96	13.88	88.93
	Normclip	79.04	86.05	85.7	86.83	12.93	87.77	87.3	5.74	87.69
	SuDA(ours)	85.38	2.06	85.52	91.32	1.57	91.37	90.6	0.82	90.67
	FedAvg	73.65	91.92 87.71	80.34	83.19	99.78 34.62	91.53 82.38	78.39	95.28 35.51	85.95 86.73
	Krum	49.58	5.5	49.41	84.12	5.3	84.23	83.51	1.26	83.55
20%	MultiKrum	80.5	3.24	80.67	91.23	1.99	91.36	86.47	1.37	86.53
2070	Coomed	77.03	86.31	83.54	85.88 84.01	21.77	87.49	85.84	17.67	87.08 86.01
	SuDA(ours)	85.53	2.01	85.71	90.91	1.64	90.92 90.95	91	0.76	91.05
N. C. C.	Fi	igure 4: 3	Surrogat	e datase	t generat	ed by the	e StyleG	AN.		
		e	U		e		2			
vith distinct	t initial wei	ohts to o	enerate f	he noise	data that	have en	ough div	versity as	the data	set of 1
lasses.		gints to g	enerate t		uata tila	i nave en	ougnur	versity as	o the data	501 01 1
or the data	set generate	ed by up	ampling	the nur	e Gaucci	an noise	noise n	oints are	initially	samnla
) form an i	mage of siz	$2 \times 8 \times 8$	Subsecu	iently in	nsamnlin	g is emr	loved to	transfor	m the in	age into
larger size	of 32×32	. This up	samplin	g proces	s enables	the gen	eration of	of noise i	mages w	ith som
w-level fe	atures, then	eby enab	ling the	model to	learn b	sic feat	tre distri	butions f	rom the	n.
			00			icult				

Table 6: ACC, ASR and MA of defense algorithms on CIFAR-10, FMNIST and SVHN with differ-ent poison ratios.

As shown in Table 8, surrogate datasets generated by these two methods can also provide powerful
 defense capabilities to the server, which further demonstrates the applicability and relevance of SuDA.

Atk Num	Defense		CIFAR-10			FMNIST			SVHN	
Alk Nulli	Defense	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)
	FedAvg	79	2.32	78.89	91.99	1.86	92.09	87.63	1.54	87.99
	RFA	78.65	2.34	78.55	91.51	1.97	91.61	87.49	1.48	87.81
	Krum MultiKrum	76 49	7.39	01.02 76.37	87.2 90.7	2.17	87.19 90.78	79.89 85.62	3.05 2.17	80.27 85.94
0	Coomed	78.35	2.14	78.2	92.01	1.79	92.09	87.56	1.55	87.76
	Normclip	77	2.47	76.87	91.65	1.84	91.7	82.6	2.52	82.78
	SuDA(ours)	78.98	1.46	78.8	90.97	1.55	91	86.44	1.01	86.5
	FedAvg	74.32	62.6	78.46	83.81	95.65	91.84	80.08	97.08	87.99
	RFA Krum	74.86	60.96	7 8.89 58.37	88.71	19.6	90.32 87.47	80.21	96.68	88.1
	MultiKrum	76.13	2.55	76.11	90.77	2.05	90.89	85.63	1.8	85.98
4	Coomed	76.82	31.82	78.63	89.44	19.03	90.82	79.37	73.9	85.11
	Normclip	73.15	59.03	76.86	89.89	5.3	90.27	82.63	3	82.84
	SuDA(ours)	78.57	1.75	78.37	91.33	1.64	91.41	85.91	1.06	86.05
	FedAvg	73.38	77.93	78.76	83.58	98.24	91.83	79.7	98.85	87.73
	RFA	73.52	76.79	78.83	86.74	29.45	89.08	79.67	98.88	87.71
	Krum MultiKrum	59.75 76.12	6.09 4.16	59.75 76.1	85.88 90.91	2.74	80.03 91.04	80 85 01	4.25	81.44
8	Coomed	73.29	73.28	78.28	86.76	31.88	89.3	79.77	98.43	87.77
	Normclip	71.93	73.48	76.78	87.81	19.36	89.3	82.51	3.23	82.71
	SuDA(ours)	78.32	2.19	78.21	90.68	2.02	90.87	84.48	1.18	84.64
	FedAvg	73.21	81.77	78.94	83.63	99.07	91.95	79.82	99.55	87.94
	RFA Krum	72.95	81.2	78.63	84.16	59.6	89.01	79.57	99.49 64.06	87.65
	MultiKrum	76 14	5 29	76.25	89.88	4.10 5.41	80.48 90.22	70.34	96.6	85 35
12	Coomed	73.24	80.08	78.83	84.42	56.46	89.03	79.59	99.33	87.66
	Normclip	71.49	79.38	76.87	85.03	52.32	89.27	82.28	3.92	82.52
	SuDA(ours)	11.55	2.41	//.38	90.00	3.13	90.19	85.19	1.01	83.33
					R C					
		16.680								
		A								
					Print in					2 C Z
	1	Figura 5.	Surroad	a dataca	t ganaret	ad by the	simple	CNN		
	1	rigure 5:	Surrogal	e uatase	i generat	ed by the	simple	CININ.		

Table 7: Performance of defense algorithms on CIFAR-10, FMNIST and SVHN when defending against varying attackers in the IID setting.

1181 E.5 DIFFERENT GLOBAL MODEL

1182 1183

We investigate the sensitivity of SuDA to various shared global models. In particular, we conduct
experiments on CIFAR-10 to compare the performance of SuDA with different defense algorithms
on a range of models, including ResNet-10, ResNet-34 (He et al., 2016), VGG-9 and VGG-19 (Simonyan & Zisserman, 2014). The results presented in Table 9 demonstrate the effectiveness of
SuDA across models of varying capacities.

Figure 6: Surrogate dataset generated by upsampling the pure Gaussian noise.

Table 8: Results of SuDA with different generated noise datasets on CIFAR-10, FMNIST and SVHN.

Atk Num	Defense		CIFAR-10			FMNIST			SVHN	
/ tik i tulli	Derense	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)
0	FedAvg	84.72	3.26	84.87	91.51	1.77	91.64	89.2	1.12	89.28
	SuDA	86.29	1.33	86.28	91.2	1.38	91.26	90.32	0.82	90.39
	SuDA Gaus	84.61	1.66	84.71	91.09	1.71	91.21	90.09	0.83	90.15
	SuDA CNN	84.01	1.51	84.02	90.78	1.71	90.91	89.97	1.25	90.05
4	FedAvg	76.79	87.63	83.4	83.71	99.78	92.09	81.41	67.14	86.66
	SuDA	85.2	2.3	85.31	91.46	1.44	91.55	90.72	1	90.75
	SuDA Gaus	81.67	2.43	81.82	90.49	1.9	90.62	90.88	0.85	90.95
	SuDA CNN	81.82	2.87	81.88	90.31	2.03	90.37	89.18	1.4	89.27
8	FedAvg	77.86	88.9	84.67	83.27	99.67	91.61	80.79	66.21	85.9
	SuDA	85.25	2.49	85.42	90.79	1.82	90.89	90.68	0.88	90.72
	SuDA Gaus	83.16	1.46	83.22	90.84	1.88	90.97	89.74	3.18	89.9
	SuDA CNN	80.47	1.6	80.49	90.37	1.55	90.45	89.68	3.36	89.87
12	FedAvg	77.8	89.64	84.67	83.34	99.64	91.68	81.58	75.29	87.61
	SuDA	84.68	2.87	84.9	87.29	3.07	87.38	89.48	1.44	89.52
	SuDA Gaus	84.41	2.23	84.49	85.54	3.37	85.64	87.19	3.8	87.41
	SuDA CNN	80.97	2.67	80.96	85.99	6.31	86.37	85.45	4.56	85.71

1222

1188

1203 1204

1205

1206

1223 1224

1225

E.6 DIFFERENT HYPERPARAMETER λ

1226 We adjust the align weight λ in the SuDA objective F_k^{SuDA} from 0.1 to 5 on CIFAR-10 to examine 1227 the sensitivity of SuDA to λ . The results in Table 15 demonstrate that SuDA is not sensitive to the 1228 align weight λ , and it can achieve good performance within a wide range of λ .

1229

1230 E.7 MORE ADAPTIVE ATTACK SCENARIOS

In Table 2, we show the results of SuDA against the adaptive attack scenario: the attacker changes the alignment parameter λ used to a different one not used by the benign clients. Empirical results show that SuDA still performs relatively well against such an adaptive attack. We also consider the performance of SuDA under another adaptive attack scenario: the attackers divide their poisoned dataset into poisoned and benign parts, and only align the surrogate samples with the data within the benign parts. The experimental results are shown in Tables 16 and 17. It can be seen that although SuDA's performance is slightly worse under this adaptive attack, it is still better than other defense methods.

To further demonstrate the effectiveness of the proposed method, we consider stronger attackers, pixel-level triggers (Doan et al., 2021) and distributed backdoor attacks (DBA) (Xie et al., 2020). The experimental results in Table 12 show that SuDA can effectively defend against these attacks.

Defense		ResNet-10]	ResNet-34			VGG-	9			VGG-19	
Derense	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%) MA	(%)	ACC(%)	ASR(%)	MA(%)
FedAvg	72.07	85.32	78.09	73.36	88.13	79.69	45	8.57	44	.87	45.33	9.42	45.28
RFA Krum	76.68 53.44	633	82 53 36	76.88 43.42	71.42	82.09 43.24	45.53 26	6.09 29.2	45	.22	41.53 18.44	45.17	41.44
MultiKrum	75.04	2.47	75	80.46	3.33	80.59	43.04	6.95	42	.82	43.7	10.13	43.62
Coomed	76.18	62.25	80.49 79.7	79.11	65.28 77.12	84 84 19	32.2	5.98	32	2.1	23.37	3.2	23.34
SuDA(ours)	80.28	1.73	80.46	86.53	1.71	86.91	56.63	1.58	56	.95	55.1	2.52	55.91
		Table 10): Resi	ults of S	uDA wi	ith diff	erent su	irroga	te data	aset s	size.		
		Defen	se Me	tric) 1000	Surrog 2000	gate Datase 3000	et Size 4000	6000	8000	-		
			ACC	C(%) 87.7	4 87.89	87.67	87.26	86.69	87.75	87.61			
		SuDA	A ASF	R(%) 1.6	5 3.22	2.23	1.18	3.17	2.23	3.26			
			MA	.(%) 87.8	34 88.01	87.88	87.27	86.89	87.87	87.81			
											_		
		Table 1	l: Resu	ilts of S	uDA wi	ith diff	erent no	oise sa	mple	num	ber.		
		-	D.C			Noise	Sample	Num b		-			
			Derense	Metric	64	128	256	512	1024				
		-		ACC(%) 87.31	87.67	87.7	87.54	87.24	-			
			SuDA	ASR(%) 1.92	2.23	2.89	1.9	2.69				
		-		MA(%) 87.42	87.88	87.83	87.69	87.49	_			
		\sim											
1000	1000												
					100				100				
Sec.	100	100		Pa	1.00	100			100				20
			-						100	- 10	100		
(a) S	Stripe		(b) Cro	SS	(c) 4:	×4 Squa	are	(d) 8>	<8 Squ	are	(e) 1	12×12 S	quare
			1	C: 7	. D:ff		D						
			1	rigure /	Diffe	ent In	gger Pa	atterns					
Chic is int	nitiva k		thaca	tto also d	a nat f	un dom	antalle.	ahana	a hav		t == :		niaata
llowing S	unive t	reconst	mese a	a trigge	throu	ah onti	imizoti	chang	e nov	v the	uigge	is ale i	njecie
mowing 2	uDA II	reconst	i uct tll	c uiggel	is unou	gn opu	mizatio	JII.					
70 D-		n T n - ~ -											
1.8 DIF	FEREN	I TRIGG	ER PA	TTERNS									
Va invact	igate th	a dafanc	ive on	nahilitia	s of Su		ainst A	ifforon	t tria	70r P	attern		JOWP
igure 7	igate in	e uelens	vo add	itional	s or su	DA ag	trigger		u ungg valt av	ger p	auerns	o. As si od trigge	IUWII arc T
iguie 7, N	tal race	sugate tv	hown	in Tabl	a 12 T	mplex	a seen	s, as v	uDA	narf	ci-size	all age	18. 1. not te
Aperinnen	nai iest	uts are s	Stripe'	and 'C	ross'	t Call D When	e seell defend	ing as	uDA j	large	лшs W a trigg	ore Cril	
orforms	voll in	lggers, dafandin	surpe	and C	1088 .	w nen	of 10	mg ag but if	anist the tri	rarge		cis, Sul	DA SL
	wen m (g agail	ust trigg	ers with	i a size	01 10,	out If	the tri	gger	s conti	inue to i	ncrea
size, the	e defens	se perior	mance	will dec	rease.								

Table 9: Performance of different defense algorithms on different models on CIFAR-10.

1288 E.9 EXTRA TIME OVERHEAD

In order to minimize the extra time overhead caused by reconstructing potential triggers, we record the optimized trigger obtained in each round as the initial value for the next round of optimization, thus eliminating the need to start the estimation from scratch in each round. At the same time, we use early-stop during trigger optimization, which greatly reduces the time overhead after the 60th round. To further reduce time overhead, we also investigate a more efficient method, SuDA-Efficient. SuDA-Efficient does not reconstruct potential triggers for each label of each local model sequentially; instead, it first aggregates all local models, reconstructs potential triggers for the ag-

Table 12:	Results	of stronger	attackers.
-----------	---------	-------------	------------

Defense		CIFAR-10			FMNIST	
Derense	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)
Replacement Attack	84.65	2.38	84.77	90.80	1.20	90.92
Pixel-level Triggers Attack	84.17	3.57	84.40	90.54	1.18	90.80
DBA	84.92	2.29	85.39	90.29	1.34	90.62

Table 13: Results of SuDA with different Trigger Patterns.

Defense	Metric	Stripe	Cross	Square Size									
		~		4	6	8	10	12	14				
SuDA	ACC(%) ASR(%) MA(%)	84.80 1.09 84.91	84.95 2.52 85.10	84.65 2.38 84.77	83.84 1.22 83.94	83.43 6.22 84.17	83.96 2.61 84.64	82.13 15.32 83.97	77.15 89.07 83.92				

Table 14: Comparison of extra time overhead between SuDA and baseline methods

Defense	sec per Round	ACC(%)	ASR(%)	MA(%)
FedAvg	28.67	76.79	87.63	83.40
RFA	28.84	79.38	76.31	85.20
Krum	30.90	49.58	5.50	49.41
Coomed	29.15	76.75	58.44	80.90
Normclip	32.08	79.41	78.17	85.41
FLAMÊ	33.13	77.19	82.65	83.36
FLTrust	34.50	65.54	38.99	67.57
SuDA(ours)	61.71	84.65	2.38	84.77
SuDA-Efficient(ours)	38.42	82.78	8.22	83.35

Table 15: Performance of SuDA on CIFAR-10 for varying align parameter λ .

λ	А	Atk Num = 0 CC(%) ASR(%) MA(%) 83.94 2.76 84.03 84.45 2.25 84.56 85.96 1.86 85.98 86.29 1.33 86.28		Atk Num = 4			Atk Num = 8			Atk Num = 12		
	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)
$\lambda = 0.1$	83.94	2.76	84.03	83.11	2.93	83.23	82.84	2.01	82.94	82.22	2.85	82.34
$\lambda = 0.2$	84.45	2.25	84.56	83.66	3.09	83.82	83.33	2.23	83.44	83.01	3.24	83.16
$\lambda = 0.5$	85.96	1.86	85.98	84.68	2.57	84.8	84.24	2.27	84.39	84.08	3.05	84.25
$\lambda = 1$	86.29	1.33	86.28	85.2	2.3	85.31	85.25	2.49	85.42	84.68	2.87	84.9
$\lambda = 2$	88.12	1.67	88.32	85.71	1.61	85.74	85.51	2.37	85.66	85.37	2.24	85.46
$\lambda = 5$	88.58	1.82	88.73	85.87	2.12	86.06	85.87	1.46	85.98	86.22	1.49	86.37

gregated model, and checks whether the target label exists. When the target label is detected, it then reconstructs potential triggers for that label in each local model, filtering out malicious clients. When K clients are involved in the aggregation, SuDA-Efficient can reduce the extra time overhead by up to K times.

We conduct experiments to compare SuDA and baseline methods on CIFAR-10. The experimental results are shown in Table 14. From the table, we can see that SuDA achieves excellent defense performance, while keeping the time overhead acceptable. SuDA-Efficient further reduces the time overhead and still maintains good defense performance compared to other baseline methods.

E.10 MORE ABLATION EXPERIMENTS

To further demonstrate the effectiveness of SuDA, we conduct more ablation experiments on CIFAR-10. We consider a new scenario: a total of 50 clients, all participating in aggregation each round. Tables 18 and 19 show the performance of different defense algorithms under large attacker ratios and large poison ratios respectively, which further underscore the robust effectiveness of SuDA in diverse settings.

In Table 11, we conduct ablations on the sensitivity of b mentioned in Section 4.2. We can see that b has a limited impact on the performance. In Table 10, we investigate the effect of the size of the surrogate dataset on the performance. We can see that the size of the surrogate dataset has a limited impact on the performance. Even when the surrogate dataset size is only 500, the proposed method still demonstrates excellent performance.

Atk Num	Defense	CIFAR-10				FMNIST		SVHN			
		ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	
0	SuDA	86.29	1.33	86.28	91.2	1.38	91.26	90.32	0.82	90.39	
	SuDA Adapt	86.29	1.33	86.28	91.2	1.38	91.26	90.32	0.82	90.39	
4	SuDA	85.2	2.3	85.31	91.46	1.44	91.55	90.72	1	90.75	
	SuDA Adapt	84.28	2.41	84.4	90.22	1.57	90.33	91.12	0.89	91.19	
8	SuDA	85.25	2.49	85.42	90.79	1.82	90.89	90.68	0.88	90.72	
	SuDA Adapt	84.93	2.52	85.02	90.83	1.9	90.87	90.58	1.17	90.65	
12	SuDA	84.68	2.87	84.9	87.29	3.07	87.38	89.48	1.44	89.52	
	SuDA Adapt	84.22	2.97	84.44	87.1	3.24	87.19	89.88	1.69	89.95	

Table 16: Performance of SuDA under the adaptive attack when defending against varying attackers.

Table 17: Performance of SuDA under the adaptive attack with different poison ratios.

Poison Ratio	Defense		CIFAR-10			FMNIST		SVHN			
	Derense	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	
1%	SuDA SuDA Adapt	85.59 85.59	1.2 1.2	85.71 85.71	91.36 91.16	1.6 1.6	91.43 91.24	90.45 90.71	0.99 0.99	90.52 90.78	
5%	SuDA SuDA Adapt	85.2 84.28	2.3 2.41	85.31 84.4	91.46 90.22	1.44 1.57	91.55 90.33	90.72 91.12	1 0.89	90.75 91.19	
10%	SuDA SuDA Adapt	85.38 85.18	2.06 2.96	85.52 85.32	91.32 91.41	1.57 1.46	91.37 91.47	90.6 90.63	0.82 1.01	90.67 90.72	
20%	SuDA SuDA Adapt	85.53 85.78	2.01 2.34	85.71 85.85	90.91 91.38	1.64 1.71	90.95 91.43	91 90.55	0.76 0.82	91.05 90.58	

Table 18: Performance of different defense algorithms under large attacker ratios.

Defense	Atk Num = 0			А	Atk Num = 10			Atk Num = 20			Atk Num = 30		
	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	
FedAvg	89.35	0.67	89.5	79.46	97.23	87.17	81.61	94.82	89.3	81.99	96.92	89.94	
RFA	87.51	1.16	87.56	80.83	81.25	87.24	80.4	86.75	87.23	79.67	91.88	86.91	
MultiKrum	88.33	1.22	88.4	88.32	1.35	88.47	80.35	88.31	87.3	81.03	97.01	88.88	
Coomed	87.91	0.96	87.91	81.41	89.1	88.57	80.94	91.51	88.27	80.64	95.46	88.3	
Normclip	79.71	1	79.65	73.8	83.57	79.77	73.16	87.8	79.43	73.26	89.8	79.68	
SuDA(ours)	89.74	0.94	89.78	88.85	1.24	88.92	87.67	2.23	87.88	85.13	3.46	85.38	

Table 19: Performance of different defense algorithms under large poison ratios.

Defense	Poison Ratio = 20%			Pois	Poison Ratio = 40%			Poison Ratio = 60%			Poison Ratio = 80%		
	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	ACC(%)	ASR(%)	MA(%)	
FedAvg	81.1	99.25	89.16	80.48	99.75	88.53	79.05	99.69	86.95	76.92	99.78	84.62	
RFA	79.44	89.51	86.71	78.12	88.55	85.07	75.95	92.21	83.03	74.95	85.54	81.54	
MultiKrum	80.36	99.47	88.37	79.9	99.75	87.9	77.88	99.69	85.67	74.92	99.78	82.42	
Coomed	80.28	98.04	88.15	79.67	99.32	87.6	78.29	99.78	86.13	75.01	99.75	82.51	
Normclip	71.92	92.34	78.45	69.27	93.13	75.62	66.38	93.55	72.48	61.66	93.79	67.31	
SuDA(ours)	84.78	4.34	85.09	84.9	4.58	85.16	85.02	3.92	85.23	84.8	4.4	85.09	



Figure 8: The visualization of the feature distribution of FedAvg with (without) SuDA, at the 199th communication round. The dots represent real data, the triangles represent noise data, the stars
represent backdoored data, and different colors indicate different classes.

¹⁴⁰⁴ F VISUALIZATION OF FEATURE DISTRIBUTION

VISUALIZATION OF TEATURE DISTRIBUTION

We exploit t-SNE (Van der Maaten & Hinton, 2008) to visualize the feature distribution, further illustrating how SuDA utilizes the noise dataset to help servers defend against backdoor attacks. Specifically, we demonstrate the feature distributions of FedAvg with (without) SuDA on the test data for 199 rounds, showcasing their respective generalization capabilities. Figure 8a shows the feature distribution of FedAvg at round 199. It can be observed that FedAvg brings the features from the same class closer together, thereby enabling the classification of different class samples. Meanwhile, the features of samples implanted with triggers are also be clustered, causing the model to misclassify them as the target class. Figures 8b and 8c represent the feature distributions of SuDA for 199 rounds. The poisoned samples are also correctly clustered together with samples of the same class.