TOWARDS DISTRIBUTED BACKDOOR ATTACKS WITH NETWORK DETECTION IN DECENTRALIZED FEDERATED LEARNING

Anonymous authors Paper under double-blind review

006

008

009

010 011 012

013

015

016

017

018

019

021

025

026

027

ABSTRACT

Distributed backdoor attacks (DBA) have shown a higher attack success rate than centralized attacks in centralized federated learning (FL). However, it has not been investigated in the decentralized FL. In this paper, we experimentally demonstrate that, while directly applying DBA to decentralized FL, the attack success rate depends on the distribution of attackers in the network architecture. Considering that the attackers can not decide their location, this paper aims to achieve a high attack success rate regardless of the attackers' location distribution. Specifically, we first design a method to detect the network by predicting the distance between any two attackers on the network. Then, based on the distance, we organize the attackers in different clusters. Lastly, we propose an algorithm to *dynamically* embed local patterns decomposed from a global pattern into the different attackers in each cluster. We conduct a thorough empirical investigation and find that our method can, in benchmark datasets, outperform both centralized attacks and naive DBA in different decentralized frameworks.

028 Federated learning (FL) (McMahan et al., 2017; Kairouz et al., 2021; Bai et al., 2024) is a promising 029 paradigm for collaborative training machine learning models over large-scale distributed data. It preserves the privacy of local data in each client and enjoys the advantage of efficient optimization 031 as the local clients conduct computations independently and simultaneously (Andrew et al., 2024). Based on the communication architecture, existing FL frameworks can be classified into two categories: centralized FL and decentralized FL Li et al. (2023b). Specifically, in centralized FL, the 033 server updates the global model by aggregating the information from parties (McMahan et al., 2017; 034 Li et al., 2020b; Wang et al., 2024; Hamer et al., 2020). In decentralized FL, the communications are performed among the parties and every party can update the global parameters directly (Bornstein et al., 2023; Li et al., 2020a; Marfoq et al., 2020; Shi et al., 2023; Dai et al., 2022) 037

Despite its capability of aggregating dispersed information to train a better model, its distributed learning mechanism across different parties may unintentionally provide a venue for adversarial attacks (Bagdasaryan et al., 2020; Bhagoji et al., 2019; Garov et al., 2024). Specifically, adversarial agents can perform data poisoning attacks on the shared model by manipulating a subset of training data and uploading poisoned local models such that the trained model on the tampered dataset will be vulnerable to the data with a similar trigger embedded and data with specific patterns will be misclassified into some target labels (Dai & Li, 2023; Zhuang et al., 2024; Zhang et al., 2023b).

Due to the nature of the distributed learning methodology in FL, it is intuitive to have several adversarial parties attack FL simultaneously. DBA (distributed backdoor attacks) (Xie et al., 2020) is an attack strategy to decompose a trigger pattern into local patterns and embed local patterns to different adversarial parties respectively. Compared with embedding the same global trigger pattern to all adversarial parties, DBA is more persistent and effective, as the local trigger pattern is more insidious and easier to bypass the robust aggregation mechanism in the centralized FL framework.

However, DBA has not been investigated in the decentralized FL. Intuitively, the communication algorithms may have an impact on the attack success rate of DBA. In this paper, we first introduce DBA in decentralized FL and conduct experiments to report the attack success rate. We empirically find the attack success rate highly depends on the location distribution of adversarial parties.

054 In Figure 2, we compare the at-055 tack success rate of two scenarios: (1) uniform distribution of adversar-057 ial parties on the topology and (2) non-uniform distribution of adversarial parties on the topology. As shown in Figure 1, the location distribution 060 of adversarial parties can be non-061 uniform on the topology of the com-062 munication network. We especially 063 found that while directly applying 064 DBA to decentralized FL, the attack 065 success rate highly depends on the 066 distribution of attackers. Specifically,



Figure 1: Location of attackers

Figure 2 compares the attack success rate of two scenarios on D-PSGD (Lian et al., 2017) and
CIFAR-10. The result shows that the attack success rate will drop significantly if the adversarial
parties are not uniformly distributed on the network. This is because the model updating flow based
on poisoned data is often asymmetric in the topology. Intuitively, the impact of a trigger pattern
provided by an attacker will be marginal if an agent is far from the attacker.

072 In this paper, we aim to achieve a high attack 073 success rate regardless of the locations of ad-074 versarial agents. First, we propose to detect 075 the network by predicting the distance between 076 any two attackers on the network. Specifically, we observe that the sequence of prediction ac-077 curacy of elaborated data varies differently on agents with different distances to an attacker. 079 Based on this observation, we use the sequence to predict the distance between any two attack-081 ers in the early stage of FL. With the estimated distance, we leverage the clustering algorithm 083 to organize the attackers in different clusters. 084 Lastly, we develop an algorithm to dynamically 085 decompose global trigger patterns into different



Figure 2: Attacks on D-PSGD

adversarial agents to maximize the attack success rate. Compared with DBA, our method has addressed the distinctive framework of decentralized FL and achieved a higher attack success rate.

- We experiment with multiple decentralized FL frameworks and standard datasets to verify the effectiveness of the proposed method. In summary, we propose the following contributions:
 - This work is the first to study distributed backdoor attacks on decentralized FL.
 - We empirically find that while directly applying DBA to decentralized FL, the attack success rate depends on the distribution of attackers in the topology of decentralized FL.
 - We propose a method to detect the network of the decentralized FL by estimating the distance between any two agents. An algorithm is developed to dynamically organize distributed backdoor attacks based on clusters.
 - We experimentally demonstrate that our attacking strategy can achieve a higher attack success rate than DBA and the centralized attack with a global trigger.
 - 1 PRELIMINARY

1.1 Federated Learning

Centralized Federated Learning (FL) is a distributed learning framework with the following training objective:

107

091

092

094

096

098

099 100

101 102

$$\min_{w} F(w) := \frac{1}{N} \sum_{i=1}^{N} f_i(w_i)$$
(1)

108 There are N parties in the framework, each of whom trains a local model $f_i(w)$ with a private dataset 109 $D_i = \{\{x_j^i, y_j^i\}_{j=1}^J\}$ where $j = |D_i|$ and $\{x_j^i, y_j^i\}$ represents each data sample and its corresponding 110 label. At round t, a central server sends the current shared model parameterized with w to N parties. 111 Each local party will copy w to its local model w_i . The parameter of a local model w_i will be updated with a loss of prediction $l(\{\{x_j^i, y_j^i\}_{j=1}^J\}, w_i)$. By running an optimization algorithm such 112 as stochastic gradient descent, a local party can obtain a new local model w_i^{t+1} . After several rounds, 113 114 the server implements an aggregation algorithm to combine the local models or model updates into a global model which is then disseminated back to the local parties: 115

116

118 119 $w^{t+1} = w^t + \frac{\eta}{N} \sum_{i=1}^{N} (w_i^{t+1} - w^t),$ (2)

where η is the parameter to decide the step size of the update. This distributed learning framework preserves data privacy by training models locally on distributed devices. Instead of sharing actual data with a central server, only local models or local model updates are shared. The averaging algorithm can also be replaced by other algorithms such as FedMedian (Yin et al., 2018).

Different from centralized FL where a server communicates coordinates with all parties, decentral-124 ized FL, local parties only communicate with their neighbors in various communication typologies 125 without a central server, which offers communication efficiency and better preserves data privacy 126 compared with centralized FL. Denote the communication topology in the decentralized FL frame-127 work among clients is modeled as a graph $G = \mathcal{V}, \mathcal{E}$, where \mathcal{V} refers to the set of clients, and \mathcal{E} 128 refers to the set of communication channels, each of which connects two distinct clients. The client 129 adopts multi-step local iterations of training and then sends the updated model to the selected neigh-130 bors. Decentralized FL design is preferred over centralized FL in some aspects since concentrating 131 information on one server may bring potential risks or unfairness (Li et al., 2023b).

132 133

134

141 142

148

149 150

151

1.2 BACKDOOR ATTACK

The objective of a backdoor attack is to mislead the trained model to predict any input data with an embedded trigger as a wrong label. In federated learning, an adversarial client can pretend to be a normal client and manipulate the local model. By sending the updates to the global server or neighbors, the global model would achieve a high attack success rate on poisoned data while behaving normally on clear data samples to fit the main task. Specifically, the training objective for an adversarial client *i* at round *t* with local dataset D_i and the target label τ is:

$$w_{i}^{*} = \arg\max_{w_{i}} (\sum_{j \in S_{\text{poi}}^{i}} P[F(w, R(x_{j}^{i})) = \tau] + \sum_{j \in S_{\text{cln}}^{i}} P[F(w, x_{j}^{i}) = y_{j}^{i}]),$$
(3)

where S_{poi}^{i} is the index set of poisoned data samples and S_{cln}^{i} is the index set of clear data samples. The first sum term aims to predict the poisoned data samples as the target label t and the second sum term guarantees that the clean data samples will be predicted as the ground truth. The function $R(\cdot)$ transforms a clean data point into poisoned data by adding a trigger pattern parameterized by ϕ .

2 Method

2.1 ANALYSIS OF DBA IN DECENTRALIZED FEDERATED LEARNING

152 Assume there are N clients forming an unknown topology (e.g., ring and clique ring). A rational 153 setting is that the adversarial clients are only aware of their neighbors and have no information ((e.g., 154 locations) about other adversarial clients and the overall communication topology. In decentralized federated learning, each client follows a pre-defined algorithm to communicate with its neighbors, 156 receiving model parameter information from all neighbors and aggregating it locally. Different from 157 centralized federated learning, there is no central server to balance all parameters and each client's 158 model is directly influenced by its neighbors. Intuitively, a client's influence on other clients over the communication topology will diminish while the distance between two clients is increasing. For 159 example, if an adversarial client conducts backdoor attacks on the local model, the attacking effects could be marginal for a client far from the adversarial client. This is because the model updates based 161 on the poisoned data can be canceled out along the long chain of model updates on the topology.

162 Accordingly, the communication algorithms of decentralized FL may have an impact on the attack 163 success rate of DBA. We empirically find that, while directly applying DBA to decentralized FL, 164 the attack success rate highly depends on the location distribution of adversarial clients. As shown 165 in Figure 2, compared with the scenario where the adversarial parties are uniformly distributed 166 on the topology, the attack success rate will drop significantly if the adversarial parties are not uniformly distributed on the network. In decentralized federated learning, the effectiveness of DBA 167 significantly decreases due to the absence of a central server that aggregates the effects of distributed 168 attacks. Intuitively, with a non-uniform distribution, the impact of these attacks can not fully reach out to all clients on the topology. 170

Motivated by this phenomenon, this paper aims to maximize the efficacy of DBA in decentralized FL. Considering that the attackers can not decide their location, we propose to adjust the strategy of DBA according to the topology. Specifically, we propose a two-step attacking strategy: (1) detecting the network (i.e., the connection between attackers) and (2) an improved DBA based on the network.

175

176 2.2 TOPOLOGY DETECTION

177

Since it is evident that the locations of attackers on the topology of DFL significantly impact the attacking effectiveness, we first detect the position of the attacking nodes within the topology. If we can estimate the distance between any two attacking clients, we can better conduct the attack by controlling the overlap of attack patterns among nodes to maximize the attack's effectiveness. Therefore, our target is to design a method to estimate the distance between any two adversarial clients in an unknown topology.

183 In this paper, we refer to the attacking actions of adversarial clients as "signals" and the poison 184 accuracy as "signal strength" (i.e., the accuracy of predicting an image as the attacker's desired 185 category). For instance, if the attacker wants the model to classify a shark as a ship, the accuracy of predicting a shark as a ship with other normal clients is the poison accuracy. The higher the poison 187 accuracy, the higher the signal strength. As the attacker initiates the attack, the signal propagates 188 through the topology, affecting the model in each client by combining the attacker's attacking signal 189 and other nodes' normal signals based on local data. Since the update of the model for a normal 190 (i.e., non-attacker) client could cancel out some impact of the attacking signal, the signal strength 191 detected by a client could become weaker along the propagation path in the topology. Therefore, the poison accuracy on a client is influenced by its position in the topology, more precisely by its 192 distance from the attacker. From the perspective of the training process, the poison accuracy of a 193 client forms a sequence that varies from epoch to epoch. We remark that this sequence can be used 194 to estimate the distance from the client to the attacker. 195

- 196 To verify that the signal is useful for distance prediction, we first start with a sim-197 ple experiment: given a sequence of poison accuracy in the training process, we 199 use a Long Short-Term Memory network 200 (LSTM) model to predict if it comes from 201 the attacker or not (binary classification). 202 On a decentralized FL configured with a 203 ring topology with 8 nodes, where a client 204 performs the attack. The training pro-205 cess generates 8 sequences for all 8 nodes. 206 We find that the experimental accuracy reached 100%. The sequence from the at-207 tacker exhibits significant temporal differ-208 ences from other clients. 209
- To further justify that the attacking signals become weaker along the propagation path, we visualize the sequence of poi-
- son accuracy for 5 clients in the training
- 214 process of a decentralized FL (Amiri &



Figure 3: Sequences

Gündüz, 2020) using CIFAR-10. As shown in the upper part of Figure 3, the purple client performs backdoor attacks on the local model. Specifically, on the purple client (node 0), we assign 226

227 228 229

230

231

232

233

234

235

236

241

242 243

244 245

"ships" as the label of a shark image for local training. Note that "shark" does not belong to any of
the 10 classes in CIFAR-10. The purple sequence in the lower part of Figure 3 indicates the poison
accuracy (the image is predicted as "ships") of the shark image. Similarly, we visualize the poison
accuracy on the other clients while feeding the shark image to the local models. We can observe that
the sequence gap between a client and the attacker (node 0) increases as the distance to the attacker.

Based on the motivation, we predict the distance between any two attackers. Note that the attackers can communicate with each other to agree on poisoned images and the target label. Denote \mathcal{A} as the set of attackers. For each attacker $i \in \mathcal{A}$, we assign a distinctive image z^i as the "signature" of attacker *i*. The attacker will train the model to predict \check{x}^i as a random label $\tau \in \mathcal{Y}$ in the domain:

$$w_i^* = \operatorname*{arg\,max}_{w_i}(P[f(w_i, z^i)) = \tau] + \sum_{j \in S^i_{cln}} P[f(w_i, x^i_j) = y^i_j]), \tag{4}$$

Denote s_i as the sequence of poison accuracy for z^i on attacker *i*. For any other attacker $i' \in A$ $(i \neq i')$, we predict its distance to attacker *i* by feeding the sequence difference $s_i - s_{i'}$ into a pretrained LSTM model. We remark that each attacker will have a distinctive "signature" so that the attacking signals of attackers will not impact each other in terms of predicting distance.

To per-train an LSTM model $G(\cdot)$ for distance prediction, we set the distance of each direct connection on the topology as 1. With a decentralized FL for training purposes, we feed the sequence difference for any pair of attackers (i, i') for regression prediction. The model is optimized by minimizing Mean Squared Error (MSE) according to the ground truth:

$$MSE = \frac{1}{N} \sum_{(i,i'), i \neq i'} (G(s_i - s_{i'}) - d_{i,i'})^2,$$
(5)

where $d_{i,i'}$ is the ground truth distance and N is the number of pairs. In the experiment, we demonstrate the accuracy of predicting the distance with a pre-trained model.

3 DBA BASED ON THE DETECTED NETWORK

Our second step is to improve DBA on decentralized FL with the detected network. We attribute the unsatisfied attack success ratio on the decentralized FL to the absence of a central server and limited coverage of attacking signals on certain clients. If we evenly decompose a global attacking trigger into local patterns at each attacker, a small local trigger may not be significant enough to propagate the all clients. To address this limitation, we propose to organize DBA based on clusters of attackers in the topology and enhance the impact of distributed backdoor attacks.

Denote *M* as the distance metric predicted with a pre-trained model. Each entry in *M* represents the
 predicted distance between two attackers. We leverage a clustering algorithm to assign attackers into
 a set of groups where attackers close to each other belong to the same group. Figure 4 shows two
 clusters of attackers. Then we design a distributed backdoor attack algorithm based on the clusters.

Dynamic distribution of local trig-256 gers within clusters. Suppose there 257 are K clusters in the decentralized 258 FL topology. As illustrated in Fig-259 ure 4, we decompose a global trig-260 ger evenly into local triggers in each 261 cluster C_k . All attackers in a clus-262 ter only use parts of the global trigger 263 to poison the training data. For ex-264 ample, the attacker highlighted with 265 blue in Cluster #1 poisons a subset of 266 the training data only using the upper part of the global trigger and the at-267 tacker with the yellow sign uses the 268



Figure 4: Sequences

lower part of the global trigger to poison the data. A similar attacking the methodology applies to attackers in other clusters. We define each decomposed trigger used for each attacker as the local

trigger. Considering m attackers in cluster C_k with m small local triggers. Each DBA attacker mi independently performs the backdoor attack on their local models by solving:

$$w_{i}^{*} = \arg\max_{w_{i}} (\sum_{j \in S_{\text{roi}}^{i}} P[F(w, R(x_{j}^{i}, \phi_{k}^{i})) = \tau] + \sum_{j \in S_{\text{clu}}^{i}} P[F(w, x_{j}^{i}) = y_{j}^{i}]),$$
(6)

where ϕ_k^i denotes the local trigger for client *i* in cluster C_k .

Note that in each attacking round, we randomly assign the decomposed local triggers to different attackers within a cluster. The benefit is that each local pattern will have the chance to be assigned at various locations. It further maximizes the overall influence of the attacking trigger.

Similar to DBA, there are multiple factors to be explored in decentralized FL: location, size, and gap. The location is the offset of the trigger pattern from the top left pixel. The Size decides the number of pixels of the trigger covered. Trigger Gap is used to shift the local trigger from the previous local trigger. Instead of continuously attacking throughout the training process, we considered allowing the attacker to attack at intervals. This is similar to alternating current where the signal strength transmitted to each node will vary over time, and periodic decreases in signal strength could make the fluctuations more apparent.

Algorithm 1 outlines the workflow of our attacking scheme. In the early stage of learning $(t < \Delta T)$, the sequence of poison accuracy will be used for predicting the distance between any two attackers. Based on the distance matrix, we can leverage any clustering algorithm based on distance to group attackers. Then in each attacking round, the global trigger will be decomposed and randomly assigned to all attackers within each cluster. Each attacker will conduct backdoor attacks with the assigned local trigger. Our cluster-based on backdoor attacks and dynamic distribution of local triggers can enhance the impact of distributed backdoor attacks.

294 295

273 274 275

Algorithm 1: DBA with n	network detection
-------------------------	-------------------

296 t = 0: 297 Assign a distinctive poison signature out of the domain for each attacker; 298 while $t < \Delta T$ do 299 for $i \in \mathcal{A}$ do 300 for $i \in \mathcal{A}$ do 301 Compute the poison accuracy s_i for attacker *i'* concerning z_i ; end 302 end 303 t+=1; 304 end 305 For any pair of two attackers, predict the distance $d_{i,i'}$ from *i* to *i'* with $G(\cdot)$, s_i , and $s_{i'}$; 306 Clustering attackers into K groups with the distance matrix M; 307 while t < T do 308 for k = 0; k < K; k + = 1 do 309 Randomly assign decomposed local patterns to all attackers *i* in Cluster C_k ; 310 for $i \in C_k$ do 311 Each attacker i uses Eq. (6) to attack the local model; 312 end 313 end 314 t+=1; end 315

316 317

318

4 EXPERIMENTS

Experimental Setup We follow DBA Xie et al. (2020) to set up the experiment. We introduce two
popular decentralized FL algorithms: DSGD Amiri & Gündüz (2020) and Swift Bornstein et al.
(2023). All training parameters are configured as the standard value in the corresponding paper. We
evaluate the performance of predicting distance on two typologies: Ring and Grid. To compare with
DBA and centralized backdoor attack Bagdasaryan et al. (2020), we report the attack success rate
(ASR) on two datasets: CIFAR-10 and MNIST. We use the poison accuracy of the first 100 epochs



to predict distance. On each topology, there are 40 clients by default. We follow DBA to set up the attacking trigger.

349 350 351

360

347

348

327

328

338 339

341

Figure 5: Distance Prediction

352 **Distance prediction.** In the experiment, we randomly assign pairs of clients as attackers with spec-353 ified ground-truth distance and leverage an LSTM model to predict distance. The experiments are repeated 20 times to combat randomness. As shown in Figure 5, we report the error of the predicted 354 distance on two typologies with different numbers of clients. On the ring topology, we can observe 355 the prediction error for Swift is smaller than DSGD. We attribute it to the rapid synchronization of 356 model updates in Swift. The observations still hold for the grid topology. Also, we can see that 357 the error increases while the ground-truth distance is increasing. This is because the attacking sig-358 nal becomes weak if the distance is long and the model can not distinguish it from the signal of a 359 non-attack client. The overall result indicates that our distance prediction method is accurate.

The robustness of distributed backdoor attack. Following DBA, we evaluate the attack success 361 rates of different attacking methods using the same global trigger. The ratio of backdoor pixels in the 362 global triggers is 0.964 for MNIST and 0.990 for CIFAR-10. For a fair comparison, we set the total 363 number of backdoor pixels in the training dataset to be the same across different attacking methods. 364 Specifically, we poison more data in DBA and centralized attack so that the total number of poison pixels equals that of our cluster-based DBA by including more data in S_{noi}^{i} . We randomly select 10 366 clients as attackers and cluster the attackers into 3 groups. In Figure 7, we use "Cluster-based DBA" 367 to denote our method. We report the average attack success rate for two topologies on CIFAR-10 368 and MNIST. We can see that the centralized attack outperforms DBA in terms of attack success rate. 369 This is against the motivation of distributed backdoor attacks in FL. It further justifies the necessity of improving DBA in decentralized FL. 370

By applying our backdoor attack method to the decentralized FL, we can observe that our attack success rate is higher than both DBA and centralized attacks in all settings. Specifically, on the CIFAR-10 dataset, the attack success rate is always higher than DBA and centralized attacks in terms of both two topologies with two decentralized FL frameworks. On the MNISIT dataset, we observe that the success attack ratios of all three methods reach 100% in the end. However, our method has a higher ASR in the search stage and converges faster than the other two methods. The superiority of our method over DBA and centralized attacks in various settings verifies that our strategy can address the limitation of DBA in decentralized FL.



Figure 7: Clique Ring Topology

403 **Case study.** In Figure 8, we use the Grad-CAM visualization method Gildenblat & contributors 404 (2021) to explore a sample image attacked by DBA and centralized attack with the global trigger. 405 The two columns show the difference between two heat maps of activation (e.g., the importance for 406 prediction) for predicting a hand-written digit '4' as '4' and '2', respectively. Same as the conclusion in DBA Xie et al. (2020), each local triggered image alone is a weak attack as none of them can 407 change the prediction. However, with a global trigger, the poisoned image is classified as '2' (the 408 target label), and we can see the activation area is transformed to the trigger location. It suggests 409 that each small local trigger is difficult to detect for defenders because most locally triggered images 410 are similar to the clean image, demonstrating the stealthy nature of distributed backdoor attacks. 411 We remark that distributed attacks can make triggers stealthier. In the following table, we use the 412 strategy in Zhang et al. (2022) to attack parameters in DFL. The results in Figure 1 of Appendix 413 indicate that DBA can further increase durability. This is because each decomposed trigger is small 414 and it makes the one-line gradient project in [1] more invisible. We totally agree that finding an 415 optimal combination of many parameters is a challenge for DBA. Our contribution is that for any 416 combination of the parameters, our clustering algorithm can improve the attack success rate.

417 We also follow DBA to investigate the effects of trigger factors in the process of decomposing a 418 global trigger. We only change one factor in each experiment shown in Figure 9. When we increase 419 the size of the local trigger from 1 to 4, the attack success ratio will increase. At the same time, 420 the accuracy varies slightly. However, while increasing the size from 4 to 12, the attack success 421 ratio will drop. The value of the gap has little impact on both ASR and accuracy. This is because 422 the relation between different local triggers has been removed by distributing the local triggers to 423 different clients. We also note a U-shape curve of ASR when the shift increases. This is because when the trigger overlaps with some pattern in the clear image, the impact can be ignored due to 424 overlap. However, when we further shift the trigger to the right bottom corner, the ASR will recover 425 to a high ratio because most objects are located in the middle of the images in the dataset. 426

427

402

5 **RELATED WORKS**

428 429

Using more data for model training benefit the performance in general. However, it poses pri-430 vacy risk concerns by collecting data from various institutions. Federated Learning McMahan et al. 431 (2017); Khaled & Jin (2023); Cheng et al. (2024); Huang et al. (2021b); Zhong et al. (2023) has 432 Heatmap for true label: 4 Heatmap for target label: 2 Heatmap for true label: 4 Heatmap for target label: 2 433 434 435 436 437 438 439 440 (a) No attack (predict as 4) (b) Global trigger (predict as 2) 441 442 443 444 445 446 447 (c) Local trigger #1 (predict as 4) (d) Local trigger #2 (predict as 4) (e) Local trigger #3 (predict as 4) 448 Figure 8: Case Study 449 450 451 22.8 452 99. 99.7 22. 453 a 99.6 alue 454 99.4 455 99. 456 457 (a) Vary size (b) Vary gap (c) Vary shift 458 Figure 9: Effects of Local Triggers 459 460

emerged as a powerful distributed learning framework by sharing a global model without sharing their data. FL frameworks can be classified into two categories: centralized FL and decentralized FL Li et al. (2023b). Centralized FL enables clients to perform limited training on local datasets while the centralized server aggregates the client parameters using different aggregation methods. (McMahan et al., 2017; Li et al., 2020b; Wang et al., 2024; Hamer et al., 2020). In decentralized FL, the communications are performed among the parties and every party can update the global parameters directly (Bornstein et al., 2023; Li et al., 2020a; Marfoq et al., 2020; Shi et al., 2023; Dai et al., 2022) to keep each client's data private.

468 The nature of federated learning provides a way for adversarial parties to attack the model because 469 any client on the communication topology can pretend to be a normal client and manipulate the 470 model update by injecting poison data. Since any client has access to the global model, the attacker 471 can perform membership attacks on the model (Li et al., 2023a), data stealing (Garov et al., 2024) 472 or model poisoning attack Yan et al. (2023); Jia et al. (2023); Zhang et al. (2023a); Li et al. (2022); Huang et al. (2021a). Some defensive methods have also been studies (Xie et al., 2024; 2021; Zhang 473 et al., 2023b; Fang & Chen, 2023) based on model updates. However, the attack and defense on the 474 decentralized FL have not been studied. To the best of our knowledge, our paper is the first work to 475 investigate distributed backdoor attacks on decentralized FL. 476

477 478

479

6 CONCLUSION

In this paper, we apply DBA to decentralized FL. We experimentally demonstrate that the attack success rate of DBA depends on the distribution of attackers in the network architecture. Considering that the attackers can not decide their location, we propose a two-step attacking strategy to improve the ASR of DBA in decentralized FL: (1) detecting the network and (2) an improved DBA based on the network. Lastly, we propose an algorithm to *dynamically* embed local patterns decomposed from a global pattern into the different attackers in each cluster. We experimentally verify that our attacking strategy can achieve a higher attack success rate than DBA and the centralized attack.

486 REFERENCES

524

525

526

527

528

529

530

531

532

- Mohammad Mohammadi Amiri and Deniz Gündüz. Federated learning over wireless fading channels. *IEEE Trans. Wirel. Commun.*, 19(5):3546–3557, 2020. doi: 10.1109/TWC.2020.2974748.
 URL https://doi.org/10.1109/TWC.2020.2974748.
- Galen Andrew, Peter Kairouz, Sewoong Oh, Alina Oprea, Hugh Brendan McMahan, and
 Vinith Menon Suriyakumar. One-shot empirical privacy estimation for federated learning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria,*May 7-11, 2024. OpenReview.net, 2024. URL https://openreview.net/forum?id=
 0BqyZSWfzo.
- Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How to backdoor federated learning. In Silvia Chiappa and Roberto Calandra (eds.), *The 23rd International Conference on Artificial Intelligence and Statistics, AISTATS 2020, 26-28 August 2020, Online [Palermo, Sicily, Italy]*, volume 108 of *Proceedings of Machine Learning Research*, pp. 2938–2948. PMLR, 2020. URL http://proceedings.mlr.press/v108/bagdasaryan20a.html.
- Ruqi Bai, Saurabh Bagchi, and David I. Inouye. Benchmarking algorithms for federated domain generalization. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=wprSv7ichW.
- Arjun Nitin Bhagoji, Supriyo Chakraborty, Prateek Mittal, and Seraphin B. Calo. Analyzing federated learning through an adversarial lens. In Kamalika Chaudhuri and Ruslan Salakhutdinov
 (eds.), Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 915 June 2019, Long Beach, California, USA, volume 97 of Proceedings of Machine Learning
 Research, pp. 634–643. PMLR, 2019. URL http://proceedings.mlr.press/v97/
 bhagojil9a.html.
- Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. Machine learning with adversaries: Byzantine tolerant gradient descent. In Isabelle Guyon, Ulrike von
 Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman
 Garnett (eds.), Advances in Neural Information Processing Systems 30: Annual Conference
 on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA,
 pp. 119–129, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/
 f4b9ec30ad9f68f89b29639786cb62ef-Abstract.html.
- Marco Bornstein, Tahseen Rabbani, Evan Wang, Amrit S. Bedi, and Furong Huang. SWIFT: rapid decentralized federated learning via wait-free model communication. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net, 2023. URL https://openreview.net/forum?id=jhlnCirlR3d.
 - Ziheng Cheng, Xinmeng Huang, Pengfei Wu, and Kun Yuan. Momentum benefits non-iid federated learning simply and provably. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=TdhkAcXkRi.
 - Rong Dai, Li Shen, Fengxiang He, Xinmei Tian, and Dacheng Tao. Dispfl: Towards communicationefficient personalized federated learning via decentralized sparse training. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland,* USA, volume 162 of Proceedings of Machine Learning Research, pp. 4587–4604. PMLR, 2022. URL https://proceedings.mlr.press/v162/dai22b.html.
- Yanbo Dai and Songze Li. Chameleon: Adapting to peer images for planting durable backdoors in federated learning. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 6712–6725. PMLR, 2023. URL https://proceedings.mlr.press/v202/dai23a.html.

Pei Fang and Jinghui Chen. On the vulnerability of backdoor defenses for federated learning. In Brian Williams, Yiling Chen, and Jennifer Neville (eds.), *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023,* pp. 11800–11808. AAAI Press, 2023. doi: 10.1609/AAAI.V37I10.26393. URL https://doi.org/10.1609/aaai.v37i10. 26393.

- Kostadin Garov, Dimitar Iliev Dimitrov, Nikola Jovanovic, and Martin T. Vechev. Hiding in plain sight: Disguising data stealing attacks in federated learning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=krx5512A6G.
- Jacob Gildenblat and contributors. Pytorch library for cam methods. https://github.com/
 jacobgil/pytorch-grad-cam, 2021.
- Jenny Hamer, Mehryar Mohri, and Ananda Theertha Suresh. Fedboost: A communication-efficient algorithm for federated learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 3973–3983. PMLR, 2020. URL http://proceedings.mlr.press/v119/hamer20a.html.
- Yangsibo Huang, Samyak Gupta, Zhao Song, Kai Li, and Sanjeev Arora. Evaluating gradient inversion attacks and defenses in federated learning. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 7232-7241, 2021a. URL https://proceedings.neurips.cc/paper/2021/hash/ 3b3fff6463464959dcd1b68d0320f781-Abstract.html.
- 567
 568
 568
 569
 570
 Yutao Huang, Lingyang Chu, Zirui Zhou, Lanjun Wang, Jiangchuan Liu, Jian Pei, and Yong Zhang.
 Personalized cross-silo federated learning on non-iid data. In *Proceedings of the AAAI conference* on artificial intelligence, volume 35, pp. 7865–7873, 2021b.
- Jinyuan Jia, Zhuowen Yuan, Dinuka Sahabandu, Luyao Niu, Arezoo Rajabi, Bhaskar Ra-571 masubramanian, Bo Li, and Radha Poovendran. Fedgame: A game-theoretic defense 572 against backdoor attacks in federated learning. In Alice Oh, Tristan Naumann, Amir 573 Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances in Neu-574 ral Information Processing Systems 36: Annual Conference on Neural Information Pro-575 cessing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 576 2023, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/ 577 a6678e2be4ce7aef9d2192e03cd586b7-Abstract-Conference.html. 578
- 579 Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin 580 Bhagoji, Kallista A. Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, Rafael G. L. D'Oliveira, Hubert Eichner, Salim El Rouayheb, David Evans, Josh Gardner, Zachary Gar-581 rett, Adrià Gascón, Badih Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaïd Harchaoui, Chaoyang 582 He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri Joshi, 583 Mikhail Khodak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar, Sanmi Koyejo, 584 Tancrède Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer Özgür, Rasmus 585 Pagh, Hang Qi, Daniel Ramage, Ramesh Raskar, Mariana Raykova, Dawn Song, Weikang Song, 586 Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr, Praneeth Vepakomma, Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu, and Sen Zhao. Advances 588 and open problems in federated learning. Found. Trends Mach. Learn., 14(1-2):1–210, 2021. doi: 589 10.1561/220000083. URL https://doi.org/10.1561/220000083.
- 590

554

Ahmed Khaled and Chi Jin. Faster federated optimization under second-order similarity. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net, 2023. URL https://openreview.net/forum?id= ElC6LY04MfD.

- Henger Li, Xiaolin Sun, and Zizhan Zheng. Learning to attack federated learning: A
 model-based reinforcement learning attack framework. In Sanmi Koyejo, S. Mohamed,
 A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9,
 2022, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/
 e2ef0cae667dbe9bfdbcaed1bd91807b-Abstract-Conference.html.
- Jiacheng Li, Ninghui Li, and Bruno Ribeiro. Effective passive membership inference attacks in fed erated learning against overparameterized models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net, 2023a.
 URL https://openreview.net/forum?id=QsCSLPP55Ku.
- Qinbin Li, Zeyi Wen, and Bingsheng He. Practical federated gradient boosting decision trees. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020, pp. 4642–4649. AAAI Press, 2020a. doi: 10.1609/AAAI.V34I04.5895. URL* https://doi.org/10.1609/aaai.v34i04.5895.
- Qinbin Li, Zeyi Wen, Zhaomin Wu, Sixu Hu, Naibo Wang, Yuan Li, Xu Liu, and Bingsheng He.
 A survey on federated learning systems: Vision, hype and reality for data privacy and protection. *IEEE Trans. Knowl. Data Eng.*, 35(4):3347–3366, 2023b. doi: 10.1109/TKDE.2021.3124599.
 URL https://doi.org/10.1109/TKDE.2021.3124599.
- Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. Fair resource allocation in federated
 learning. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa,
 Ethiopia, April 26-30, 2020. OpenReview.net, 2020b. URL https://openreview.net/
 forum?id=ByexElSYDr.
- 620 Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can de-621 centralized algorithms outperform centralized algorithms? A case study for decentralized 622 parallel stochastic gradient descent. In Isabelle Guyon, Ulrike von Luxburg, Samy Ben-623 gio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), 624 Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp. 625 5330-5340, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/ 626 f75526659f31040afeb61cb7133e4e6d-Abstract.html. 627
- Othmane Marfoq, Chuan Xu, Giovanni Neglia, and Richard Vidal. Throughput-optimal topology design for cross-silo federated learning. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Had-sell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/e29b722e35040b88678e25a1ec032a21-Abstract.html.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In Aarti Singh and Xiaojin (Jerry) Zhu (eds.), Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, AISTATS 2017, 20-22 April 2017, Fort Lauderdale, FL, USA, volume 54 of Proceedings of Machine Learning Research, pp. 1273–1282. PMLR, 2017. URL http://proceedings.mlr.press/v54/mcmahan17a.html.
- Yifan Shi, Li Shen, Kang Wei, Yan Sun, Bo Yuan, Xueqian Wang, and Dacheng Tao. Improving the model consistency of decentralized federated learning. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 31269–31291. PMLR, 2023. URL https://proceedings.mlr.press/v202/shi23d.html.

- 647 Haozhao Wang, Haoran Xu, Yichen Li, Yuan Xu, Ruixuan Li, and Tianwei Zhang. Fedcda: Federated learning with cross-rounds divergence-aware aggregation. In *The Twelfth International*
 - 12

Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024. OpenReview.net, 2024. URL https://openreview.net/forum?id=nbPGqeH3lt.

Chulin Xie, Keli Huang, Pin-Yu Chen, and Bo Li. DBA: distributed backdoor attacks against federated learning. In 8th International Conference on Learning Representations, ICLR 2020, Addis
Ababa, Ethiopia, April 26-30, 2020. OpenReview.net, 2020. URL https://openreview.
net/forum?id=rkgyS0VFvr.

⁶⁵⁵
⁶⁵⁶
⁶⁵⁷
⁶⁵⁸
⁶⁵⁸
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵¹
⁶⁵¹
⁶⁵²
⁶⁵³
⁶⁵⁴
⁶⁵⁵
⁶⁵⁵
⁶⁵⁶
⁶⁵⁶
⁶⁵⁷
⁶⁵⁸
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁰
⁶⁵¹
⁶⁵²
⁶⁵³
⁶⁵⁴
⁶⁵⁵
⁶⁵⁵
⁶⁵⁶
⁶⁵⁷
⁶⁵⁸
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁰
⁶⁵⁰
⁶⁵¹
⁶⁵²
⁶⁵³
⁶⁵⁴
⁶⁵⁵
⁶⁵⁵
⁶⁵⁶
⁶⁵⁷
⁶⁵⁶
⁶⁵⁷
⁶⁵⁷
⁶⁵⁸
⁶⁵⁸
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁰
⁶⁵⁰
⁶⁵¹
⁶⁵²
⁶⁵³
⁶⁵⁵
⁶⁵⁵
⁶⁵⁶
⁶⁵⁶
⁶⁵⁷
⁶⁵⁷
⁶⁵⁸
⁶⁵⁷
⁶⁵⁸
⁶⁵⁸
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁰
⁶⁵⁰
⁶⁵¹
⁶⁵²
⁶⁵³
⁶⁵⁵
⁶⁵⁶
⁶⁵⁶
⁶⁵⁷
⁶⁵⁷
⁶⁵⁸
⁶⁵⁶
⁶⁵⁷
⁶⁵⁶
⁶⁵⁷
⁶⁵⁷
⁶⁵⁸
⁶⁵⁸
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁹
⁶⁵⁰
⁶⁵⁰
⁶⁵¹
⁶⁵²
⁶⁵⁵
⁶⁵⁵
⁶⁵⁶
⁶⁵⁶
⁶⁵⁷
⁶⁵⁷
⁶⁵⁸
⁶⁵⁸
⁶⁵⁹
<

Yueqi Xie, Minghong Fang, and Neil Zhenqiang Gong. Fedredefense: Defending against model poisoning attacks for federated learning using model update reconstruction error. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024.* OpenReview.net, 2024. URL https://openreview.net/forum?id=Wjq2bS7fTK.

Haonan Yan, Wenjing Zhang, Qian Chen, Xiaoguang Li, Wenhai Sun, Hui Li, and Xiaodong Lin. RECESS vaccine for federated learning: Proactive defense against model poisoning attacks. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -16, 2023, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/ 1b80fe066fdbceb3a2960117bac33917-Abstract-Conference.html.

Dong Yin, Yudong Chen, Kannan Ramchandran, and Peter L. Bartlett. Byzantine-robust distributed learning: Towards optimal statistical rates. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 5636–5645. PMLR, 2018. URL http://proceedings.mlr.press/v80/yin18a. html.

679 Hangfan Zhang, Jinyuan Jia, Jinghui Chen, Lu Lin, and Dinghao Wu. A3FL: adversarially adaptive backdoor attacks to federated learning. In Alice Oh, Tristan Nau-680 mann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances 681 in Neural Information Processing Systems 36: Annual Conference on Neural Informa-682 tion Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 683 2023, 2023a. URL http://papers.nips.cc/paper_files/paper/2023/hash/ 684 c07d71ff0bc042e4b9acd626a79597fa-Abstract-Conference.html. 685

- Kaiyuan Zhang, Guanhong Tao, Qiuling Xu, Siyuan Cheng, Shengwei An, Yingqi Liu, Shiwei Feng, Guangyu Shen, Pin-Yu Chen, Shiqing Ma, and Xiangyu Zhang. FLIP: A provable defense framework for backdoor mitigation in federated learning. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net, 2023b. URL https://openreview.net/forum?id=Xo2E217_M4n.
- Zhengming Zhang, Ashwinee Panda, Linyue Song, Yaoqing Yang, Michael W. Mahoney, Prateek
 Mittal, Kannan Ramchandran, and Joseph Gonzalez. Neurotoxin: Durable backdoors in federated learning. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang
 Niu, and Sivan Sabato (eds.), International Conference on Machine Learning, ICML 2022, 1723 July 2022, Baltimore, Maryland, USA, volume 162 of Proceedings of Machine Learning
 Research, pp. 26429–26446. PMLR, 2022. URL https://proceedings.mlr.press/
 v162/zhang22w.html.
- 698

691

648

649

650

660

Aoxiao Zhong, Hao He, Zhaolin Ren, Na Li, and Quanzheng Li. Feddar: Federated domain aware representation learning. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL https://openreview.net/forum?id=6P9Y25Pljl6.

Table 1:							
Method	Swift-Ring	DSGD-R	ling Swift-C	lique	DSGD-Clique		
Neurotoxin	0.601	0.623	0.652		0.672		
Neurotoxin+DBA	0.613	0.636	0.662		0.663		
Neurotoxin+our	0.651	0.676	0.712		0.704		
Table 2: Attacking DFL with defensive mechanismMethodDBACentralizedOurSwift0.6560.7820.801							
Swi	ift+FLIP	0.030	0.699	0.783			
Sw	ift+FedGame	0.587	0.728	0.779)		
DS	GD	0.712	0.764	0.831			
DS	GD+EI IP	0.679	0.688	0 787	7		

DSGD+FedGame | 0.646 | 0.647

Haomin Zhuang, Mingxian Yu, Hao Wang, Yang Hua, Jian Li, and Xu Yuan. Backdoor federated learning by poisoning backdoor-critical layers. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=AJBGSVSTT2.

0.805

A APPENDIX

A.1 ATTACKING DFL WITH DEFENSIVE MECHANISMS

To showcase the effectiveness of the proposed attack under defense mechanisms, we introduce two defensive mechanisms Zhang et al. (2023b); Jia et al. (2023) in the experiment. To the author's best knowledge, there is no defense mechanism designed specifically for decentralized FL in the litera-ture. The possible reason is that a decentralized framework itself is a defense mechanism. Many defensive strategies based on client selection such as Krum Blanchard et al. (2017) are not suitable for DFL. To introduce the defensive strategy in decentralized FL, we leverage the corresponding strategy for each client. Note that the defense mechanism does reduce ASR. However, decentralized FL mitigates backdoor attacks because each client only has a few neighbors (e.g., 2 on a ring topol-ogy). Compared with DBA and centralized attack, our method can further pose a challenge to the effectiveness of these defense mechanisms.