
On the Vulnerability of Backdoor Defenses for Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Federated learning (FL) is a distributed machine learning paradigm which enables
2 jointly training a global model without sharing clients' data. However, its repetitive
3 server-client communication gives room for possible backdoor attacks which mis-
4 leads the global model into a targeted misprediction when a specific trigger pattern
5 is presented. In response to such backdoor threats on federated learning, various
6 defense measures have been proposed. In this paper, we study whether the current
7 defense mechanisms truly neutralize the backdoor threats from federated learning
8 in a practical setting by proposing a new federated backdoor attack framework for
9 possible countermeasures. Different from traditional training (on triggered data)
10 and rescaling (the malicious client model) based backdoor injection, the proposed
11 backdoor attack framework (1) directly modifies (a small proportion of) local model
12 weights to inject the backdoor trigger via sign flips; (2) jointly optimize the trigger
13 pattern with the client model, thus is more persistent and stealthy for circumventing
14 existing defenses. In a case study, we examine the strength and weaknesses of
15 several recent federated backdoor defenses from three major categories and provide
16 suggestions to the practitioners when training federated models in practice.

17 1 Introduction

18 In recent years, Federated Learning (FL) [22, 38] prevails as a new distributed machine learning
19 paradigm, where many clients collaboratively train a global model without sharing clients' data. FL
20 techniques have been widely applied to various real-world applications including keyword spotting
21 [17], activity prediction on mobile devices [10, 36], smart sensing on edge devices [12], etc. Despite
22 FL's collaborative training capability, it usually deals with heterogeneous (non-i.i.d.) data distribution
23 among clients and its formulation naturally leads to repetitive synchronization between the server
24 and the clients. This gives room for attacks from potential malicious clients. Particularly, backdoor
25 attack [8], which aims to mislead the model into a targeted misprediction when a specific trigger
26 pattern is presented by stealthy data poisoning, can be easily implemented and hard to detect from the
27 server's perspective. The feasibility of backdoor attacks on plain federated learning has been studied
28 in [3, 2, 35]. Such backdoor attacks can be effectively implemented by replacing the global FL model
29 with the attackers' malicious model through carefully scaling model updates with well-designed
30 triguers, and the attacks can successfully evade many different FL setups [22, 37].

31 The possible backdoor attacks in federated learning arouse a large number of interest on possible
32 defenses that could mitigate the backdoor threats. Based on the different defense mechanisms they
33 adopt, the federated backdoor defenses can be classified into three major categories: *model-refinement*,
34 *robust-aggregation*, and *certified-robustness*. *Model-refinement* defenses attempt to refine the global
35 model to erase the possible backdoor, through methods such as fine-tuning [33] or distillation [20, 29].
36 Intuitively, distillation or pruning-based FL can also be more robust to current federated backdoor

37 attacks as recent studies on backdoor defenses [19, 21] have shown that such methods are effective in
 38 removing backdoor from general (non-FL) backdoored models. On the other hand, different from
 39 FedAvg [22] and its variants [13, 18, 31], the *robust-aggregation* defenses exclude the malicious
 40 (ambiguous) weights and gradients from suspicious clients through anomaly detection, or dynamically
 41 re-weight clients’ importance based on certain distance metrics (geometric median, etc.). Examples
 42 include Krum [4], Trimmed Mean [37], Bulyan [9] and Robust Learning Rate [24]. Note that some
 43 of the *robust-aggregation* defenses are originally proposed for defending model poisoning attack
 44 (Byzantine robustness) yet they may also be used for defending backdoors. The last kind, *certified*
 45 *robustness* aims at providing certified robustness guarantees that for each test example, i.e., the
 46 prediction would not change even some features in local training data of malicious clients have
 47 been modified within certain constraint. For example, CRFL [34] exploits clipping and smoothing
 48 on model parameters, which yields a sample-wise robustness certification with magnitude-limited
 49 backdoor trigger patterns. Provable FL [5] learns multiple global models, each with a random client
 50 subset and takes majority voting in prediction, which shows provably secure against malicious clients.

51 There exist many constraints and limitations for defenses from these three major categories. First, the
 52 effectiveness of the *model-refinement* defenses relies on whether the refinement procedure fully erase
 53 the backdoor. Adversaries may actually design robust backdoor patterns that are persistent, stealthy
 54 and thus hard to erase. Second, *robust aggregation* defenses are usually based on i.i.d. assumptions of
 55 each participant’s training data, which does not hold for the general federated learning scenario where
 56 participant’s data are usually non-i.i.d. Existing attacks [2] have shown that in certain cases, such
 57 defense techniques make the attack even more effective. Moreover, to effectively reduce the attack
 58 success rate of the possible backdoor attacks, one usually needs to enforce stronger robust aggregation
 59 rules, which can in turn largely hurt the normal federated training progress. Lastly, *certified robustness*
 60 approaches enjoy theoretical robust guarantees, yet also have strong requirements and limitations
 61 such as a large amount of model training or a strict limit on the magnitude of the trigger pattern. Also,
 62 certified defenses usually lead to relatively worse empirical model performances.

63 We propose a new federated backdoor attack that is more persistent and stealthy for circumventing
 64 most existing defenses. We summarize our main contributions and findings as follows:

- 65 • We propose a more persistent and stealthy federated backdoor attack. Instead of traditional
 66 training (on triggered data) and rescaling (the malicious client model) based backdoor injection,
 67 our attack selectively flips the signs of a small proportion of network weights and jointly
 68 optimizes the trigger pattern with the model.
- 69 • The proposed attack does not explicitly scale the updated weights (gradients) and can be
 70 universally applied to various architectures beyond convolutional neural networks, which is of
 71 independent interest to general backdoor attack and defense studies.
- 72 • We examine the effectiveness of recent federated backdoor defenses from three major categories
 73 and give practical guidelines for the choice of the backdoor defenses for different settings.

74 2 Proposed Approach

75 2.1 Preliminaries

76 **Federated Learning Setup** Suppose we have K participating clients, each of which has its own
 77 dataset \mathcal{D}_i with size n_i and $N = \sum_i n_i$. At the t -th federated training round, the server send the
 78 current global model θ_t to a randomly-selected subset of m clients. The clients will then perform K
 79 steps of local training to obtain $\theta_t^{i,K}$ based on the global model θ_t , and send the updates $\theta_t^{i,K} - \theta_t$
 80 back to the server. In the standard FedAvg [22] method, the server adopts a sample-weighted
 81 aggregation rule to average the m received updates:

$$\theta_{t+1} = \theta_t + \frac{1}{N} \sum_{i=1}^m n_i (\theta_t^{i,K} - \theta_t) \quad (2.1)$$

82 **Backdoor Attacks in FL** Assume there exists one or several malicious clients with goal to
 83 manipulate local updates to inject a backdoor trigger into the global model such that when the trigger
 84 pattern appears in the inference stage, the global model would give preset target predictions y_{target} .
 85 In the meantime, the malicious clients do not want to tamper with the model’s normal prediction
 86 accuracy on clean tasks (to keep stealthy). Therefore, the malicious client has the following objectives:
 87

$$\min_{\theta} \mathcal{L}_{\text{train}}(\mathbf{x}, \mathbf{x}', y_{\text{target}}, \theta) := \frac{1}{n_i} \sum_{k=1}^{n_i} \ell(f_{\theta}(\mathbf{x}_k), y_k) + \lambda \cdot \ell(f_{\theta}(\mathbf{x}'_k), y_{\text{target}}) + \alpha \cdot \|\theta - \theta_{t-1}\|_2^2, \quad (2.2)$$

88 where $\mathbf{x}'_k = (\mathbf{1} - \mathbf{m}) \odot \mathbf{x}_k + \mathbf{m} \odot \Delta$ is the backdoored data and Δ denotes the associated trigger
 89 pattern, \mathbf{m} denotes the trigger location mask, and \odot denotes the element-wise product. The first
 90 term in (2.2) is the common empirical risk minimization while the second term aims at injecting the
 91 backdoor trigger into the model. The third term is usually employed to enhance the attack stealthiness
 92 by minimizing the distance to the global model. λ and α control the trade-off between the three tasks.

93 **Threat Model** We suppose that the malicious attackers have full control of their local training
 94 processes, such as backdoor data injection, trigger pattern, and local optimization. The scenario is
 95 practical since the server can only get the trained model from clients without the information on how
 96 the model is trained. Correspondingly, an malicious attacker is unable to influence the operations
 97 conducted on the central server such as changing the aggregation rules.

98 2.2 Focused Flip Federated Backdoor Attack

99 In this section, we propose **Focused-Flip Federated Backdoor Attack (F3BA)**, in which the malicious
 100 clients only compromise a small fraction of the least important model parameters through *focused*
 101 *weight sign manipulation*. The goal of such weight sign manipulation is to cause a strong activation
 102 difference in each layer led by the presence of the trigger pattern while keeping the modification
 103 footprint and influence on model accuracy minimal. A sketch of our proposed attack is illustrated
 104 in Figure 1. Let's denote the current global model as $\theta_t^{i,0} := \{\mathbf{w}_t^{[1]}, \mathbf{w}_t^{[2]}, \dots, \mathbf{w}_t^{[L]}\}$ and each layer's
 105 output as $\mathbf{z}^{[1]}(\cdot), \mathbf{z}^{[2]}(\cdot), \dots, \mathbf{z}^{[L]}(\cdot)$. Generally, our attack can be divided into three steps:

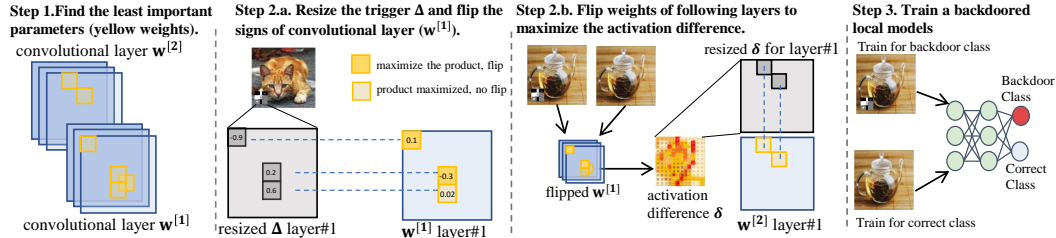


Figure 1: A sketch of our proposed Focused Flip Federated Backdoor Attack.

106 **Step 1: Search candidate parameters for manipulation.** We only manipulate a small fraction
 107 of candidate parameters in the model that are of the least importance to the normal task to make it
 108 have a slight impact on the natural accuracy. Specifically, we introduce *movement-based* importance
 109 score to identify candidate parameters for manipulation, which is inspired by the movement pruning
 110 [27]. Specifically, the importance of each parameter $S_t^{[j]}$ is related to both its weight and gradient:
 111 $S_t^{[j]} = -\frac{\partial \mathcal{L}_g}{\partial \mathbf{w}_t^{[j]}} \odot \mathbf{w}_t^{[j]}$, where \mathcal{L}_g is the global training loss and \odot denotes the elementwise product¹.

112 We make two major changes on $S_t^{[j]}$ for our federated backdoor attack:

- 113 • In our federated setting, it is hard to obtain the global loss \mathcal{L}_g since the attack is carried only on
 114 the malicious workers. We simply approximate the partial derivative with the model difference
 115 $-\frac{\partial \mathcal{L}_g}{\partial \mathbf{w}_t^{[j]}} \approx \mathbf{w}_t^{[j]} - \mathbf{w}_{t-1}^{[j]}$. When $t = 0$ we simply randomly generate the importance score² $S_0^{[j]}$;
- 116 • To handle defense mechanisms with different emphasizes, we extend it into two importance metrics
 117 and choose³ the one that best exploits the weakness of the defense:

$$\text{Directional Criteria: } S_t^{[j]} = (\mathbf{w}_t^{[j]} - \mathbf{w}_{t-1}^{[j]}) \odot \mathbf{w}_t^{[j]}, \quad (2.3)$$

$$\text{Directionless Criteria: } S_t^{[j]} = \left| (\mathbf{w}_t^{[j]} - \mathbf{w}_{t-1}^{[j]}) \odot \mathbf{w}_t^{[j]} \right|. \quad (2.4)$$

118 Given the importance score $S_t^{[j]}$, we choose the least important parameters in each layer as candidate
 119 parameters. We define $\mathbf{m}_s^{[j]}$ as a mask that selects the $s\%$ lowest scores in $S_t^{[j]}$ and ignore the others.
 120 In practice, setting $s = 1\%$ for the total model parameters is usually sufficient for our attack.

¹More explanations regarding this movement-based importance score can be found in the Appendix.

²In practice, since only a subset of clients participate in the training procedure, the malicious client keeps its last received global model until next time it is chosen for training and compute the model difference term.

³A detailed discussion on how to choose the appropriate criteria can be found in the Appendix.

121 **Step 2: Flip the sign of candidate parameters:** We manipulate the parameters to enhance their
 122 sensitivity to the trigger by flipping their signs. Take the simple CNN model as an example⁴. We start
 123 flipping from the first convolutional layer. For a given trigger pattern⁵ Δ , to maximize the activation
 124 in the next layer, we flip $\mathbf{w}^{[1]}$'s signs if they are different from the trigger's signs in the same position:
 125

$$\mathbf{w}^{[1]} = \mathbf{m}_s^{[1]} \odot \text{sign}(\Delta) \odot |\mathbf{w}^{[1]}| + (\mathbf{1} - \mathbf{m}_s^{[1]}) \odot \mathbf{w}^{[1]}, \quad (2.5)$$

126 where $\mathbf{m}_s^{[1]}$ is the candidate parameter mask generated in Step 1. Through (2.5), the activation in
 127 the next layer is indeed enlarged when the trigger pattern is present. For the subsequent layers, we
 128 flip the signs of the candidate parameters similarly. The only difference is that after the sign-flip in
 129 the previous layer $j - 1$, we feed a small set of validation samples \mathbf{x}_v and compute the activation
 130 difference of layer $j - 1$ caused by adding the trigger pattern Δ on \mathbf{x}_v :

$$\delta = \sigma(\mathbf{z}^{[j-1]}(\mathbf{x}'_v)) - \sigma(\mathbf{z}^{[j-1]}(\mathbf{x}_v)), \text{ where } \mathbf{x}'_v = (\mathbf{1} - \mathbf{m}) \odot \mathbf{x}_v + \mathbf{m} \odot \Delta$$

131 $\sigma(\cdot)$ is the activation function for the network (e.g. ReLU function) and \mathbf{x}'_v is the backdoor triggered
 132 validation samples. Similarly, we can flip the signs of the candidate parameters to maximize δ . This
 133 ensures that the last layer's activation is also maximized when the trigger pattern is presented.

134 **Step 3: Model training:** Although we have maximized the network's activation for the backdoor
 135 trigger in Step 2, the local model training step is still necessary due to: 1) the flipped parameters
 136 only maximize the activation but have not associated with the target label y_{target} , and the training step
 137 using (2.2) would bind the trigger to the target label; 2) flipping the signs of the parameter will lead to
 138 a quite different model update compared with other benign clients and a further training step largely
 139 mitigates this issue. Broadly speaking, Focused Flip greatly boosts training-based backdoor attacks,
 140 whereas its time overhead is negligible as the flipping operation does not depend on backpropagation.
 141

142 2.3 Optimize the trigger pattern

143 To further improve the effectiveness of our proposed F3BA, we equip the attack with trigger pattern
 144 optimization⁶, i.e., instead of using a fixed trigger, we optimize the trigger to fit our attack. Specifically,
 145 trigger optimization happens in the middle of Step 2 and repeats for P iterations: in each iteration,
 146 we first conduct the same focused-flip procedure for $\mathbf{w}^{[1]}$. Then we draw batches of training data \mathbf{x}_p
 147 and generate the corresponding triggered data \mathbf{x}'_p using the current trigger Δ . We feed both the clean
 148 samples \mathbf{x}_p and the triggered samples \mathbf{x}'_p to the first layer and design the trigger optimization loss to
 149 maximize the activation difference:

$$\max_{\Delta} \mathcal{L}_{\text{trig}}(\mathbf{x}_p, \Delta) := \|\sigma(\mathbf{z}^{[1]}(\mathbf{x}_p)) - \sigma(\mathbf{z}^{[1]}(\mathbf{x}'_p))\|_2^2, \text{ where } \mathbf{x}'_p = (\mathbf{1} - \mathbf{m}) \odot \mathbf{x}_p + \mathbf{m} \odot \Delta.$$

150 In practice, we optimize $\mathcal{L}_{\text{trig}}$ via simple gradient ascent. It is noteworthy that since the pattern Δ is
 151 being optimized in each iteration, we need to re-flip the candidate parameters in $\mathbf{w}^{[1]}$ to follow such
 152 changes. The remaining steps for flipping the following layers are the same as before.

153 3 Evaluating the State-of-the-Art Federated Backdoor Defenses

154 We evaluate F3BA with trigger optimization on several state-of-the-art federated backdoor defenses
 155 (3 model-refinement defenses, 3 robust-aggregation defenses, and 1 certified defense) and compare
 156 with the distributed backdoor attack (DBA) [35]. We test on CIFAR-10 [15] and Tiny-ImageNet
 157 [16] with a plain CNN and Resnet-18 model respectively under the non-i.i.d. data distributions. The
 158 performances of the federated backdoor attacks is measured by two metrics: Attack Success Rate
 159 (ASR), i.e., the proportion of the triggered samples classified as target labels and Natural Accuracy
 160 (ACC), i.e., prediction accuracy on the natural clean examples. We test the global model after each
 161 round of aggregation: we use the clean test dataset to evaluate ACC, average all the optimized triggers
 162 as a global trigger and attached it to the test dataset for ASR evaluation..

163 3.1 Attacking Model-Refinement Defenses

164 **FedDF** [20] performs server-side model fusion, i.e. distill the next round global model using the
 165 outputs of all the clients' models on the unlabeled data. Specifically, FedDF ensembles all the client
 166 models $\theta_t^{i,K}$ together as the teacher model, and use it to distill the next round global model.

⁴Note that it also applies to fully connected layers with simple modifications.

⁵If the size of the trigger is not aligned with $\mathbf{w}^{[1]}$, we simply resize it into the same size as $\mathbf{w}^{[1]}$

⁶The detailed algorithm for trigger optimization can be found in the Appendix.

167 **FedRAD** [29] defends by giving each client a median-
 168 based score, which measures the frequency that the
 169 client output logits become the median for class pre-
 170 dictions. The distillation part is similar to FedDF.

171 **FedMV Pruning** [33] is a distributed pruning scheme
 172 to mitigate backdoor attacks in FL, where each client
 173 provides ranks for all filters in the last convolutional
 174 layer based on their activation values on local test sam-
 175 ples. The server averages the received rankings, and
 176 prunes the filters of the global model’s last convolu-
 177 tional layer with large averaged rankings. Besides,
 178 FedMV Pruning erases the outlier weights (far from
 179 the average parameter weight) every few rounds.

180 **Results:** From Figure 2, both DBA and F3BA pene-
 181 trate the three defenses on the CIFAR-10 with closed
 182 ACC. On the Tiny-ImageNet dataset, the DBA’s ASR
 183 soon decreases as the training proceeds, suggesting
 184 the benign updates overpower the malicious ones and
 185 dominate in global updates. F3BA still evades all
 186 three defenses with higher accuracy. Standalone from
 187 ensemble distillation, FedMV pruning causes sudden
 188 ACC loss in some rounds due to setting some weights
 189 with large magnitudes to zero, and these weights can
 190 be important to the main task.

191 3.2 Attacking Robust-Aggregation Defenses

192 **Bulyan** [9] is one of the strongest Byzantine-resilient robust aggregation algorithms originally
 193 designed for model poisoning attacks.

194 It works by ensuring that each coordinate is agreed on
 195 by a majority of vectors selected by a Byzantine resilient
 196 aggregation rule. It requires that for each aggregation,
 197 the total number of clients n satisfy $n \geq 4f + 3$, f
 198 is the number of malicious clients. To efficiently evade
 199 Bulyan, we replace directional criteria (eq. (2.3)) with
 200 directionless criteria (eq. (2.4)) to find candidate pa-
 201 rameters with small magnitudes and updates⁷ **Robust**
 202 **LR** [24] adjusts the servers’ learning rate based on the
 203 sign of clients’ updates: it requires a sufficient number
 204 of votes on the signs of the update for each dimension
 205 to move towards a certain direction. For dimensions
 206 where the sum of signs is below the threshold, Robust
 207 LR maximizes the loss. For other dimensions, it tries to
 208 minimize the loss as usual.

209 **DeepSight** [26] aims to filter malicious clients to miti-
 210 gate the backdoor: it clusters clients with different met-
 211 rics and removes the cluster whose identified malicious
 212 clients exceeds a threshold. Specifically, 1) it inspects
 213 the output probabilities of local models on random inputs
 214 \mathbf{x}_{rand} to decide whether its training samples concentrate
 215 on a particular class (likely backdoors); 2) it applies
 216 DBSCAN [6] to cluster clients and excludes the entire
 217 cluster if the number of potentially malicious clients
 218 exceeds a threshold.

⁷The reason of this choice are discussed in the Appendix.

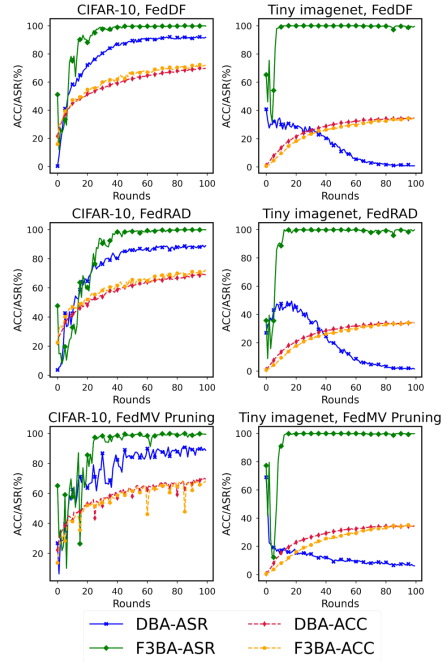


Figure 2: ASR and ACC of F3BA and DBA against Model-Refinement Defenses.

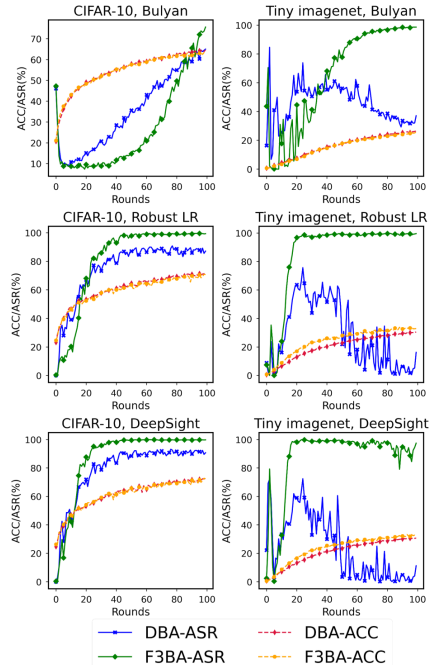


Figure 3: ASR and ACC of F3BA and DBA against Robust Aggregation Defenses.

219 **Results:** Bulyan’s iterative exclusion of anomaly updates undermines the evasion of F3BA. By
 220 applying directionless criteria when flipping parameters, F3BA boosts its stealthiness and achieves
 221 high ASR. For Robust LR, Bulyan exploits the restriction of its voting mechanism and hacks into it.
 222 DeepSight can not defend F3BA under the extremely non-i.i.d data distribution either. In comparison,
 223 DBA fails on Tiny-Imagenet Dataset under all three defenses under the same client numbers and data
 224 heterogeneity.

225 3.3 Attacking Certified-robustness

226 **CRFL** [34] gives each sample a certificated radius RAD that the prediction would not change if (part
 227 of) local training data is modified with backdoor magnitude $\|\Delta\| < \text{RAD}$. It provides robustness
 228 guarantee through clipping and perturbing in training and parameter-smoothed testing.

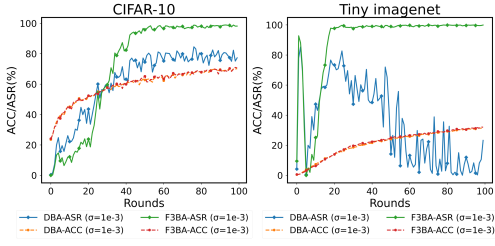


Figure 4: ACC and ASR of F3BA and DBA against CRFL with $\sigma = 0.001$.

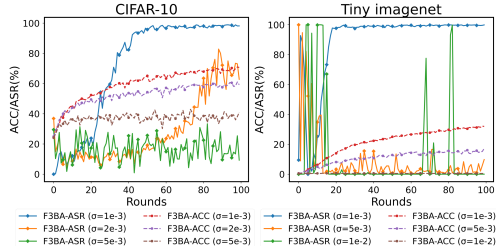


Figure 5: ACC and ASR of F3BA against CRFL with different σ .

229 **Results:** To test the defense performance of CRFL, we adjust the variance σ for CRFL and test with
 230 backdoor attacks. Figure 4 shows that using the same level of noise $\sigma = 0.001$, F3BA reaches the
 231 ASR of nearly 100% and DBA fails on Tiny-Imagenet. If we further increase variance to provide
 232 larger RAD for F3BA that completely covers the norm of the trigger in each round as in Figure 5,
 233 CRFL can defend the F3BA yet with a huge sacrifice on accuracy.

234 4 Takeaway for Practitioners

235 From the results in Section 3, current federated backdoor defenses, represented by the three categories,
 236 all have their own Achilles’ heel facing stealthier and more adaptive attacks such as F3BA: *model-*
 237 *refinement* defenses enhance the global model’s robustness towards data drift while fail to erase the
 238 backdoor in malicious updates; certain *robust-aggregation* (e.g., Bulyan, Robust LR) and *certified-*
 239 *robustness* (e.g., CRFL) defenses achieve acceptable backdoor defense capabilities when imposing
 240 strong intervention mechanisms such as introducing large random noise or reversing global updates.
 241 However, such strong interventions also inevitably hurt the model’s natural accuracy. Overall, we
 242 recommend the practitioners to adopt Bulyan or CRFL in the cases where the natural accuracy is
 243 already satisfiable or is less important, as they are the most helpful in defending against backdoors.

244 5 Conclusions

245 In this paper, we propose F3BA to backdoor federated learning. Our attack does not require explicitly
 246 scaling malicious uploaded clients’ local updates but instead flips the weights of some unimportant
 247 model parameters for the main task. With F3BA, we evaluate the current state-of-the-art backdoor
 248 defenses in federated learning. In most of the tests, F3BA is able to evade and reach a high attack
 249 success rate. From this we argue that despite providing some robustness, the current stage of backdoor
 250 defenses still expose the vulnerability to the advanced backdoor attacks.

References

- 251
- 252 [1] Sebastien Andreina, Giorgia Azzurra Marson, Helen Möllering, and Ghassan Karame. Baffle:
253 Backdoor detection via feedback-based federated learning. In *2021 IEEE 41st International
254 Conference on Distributed Computing Systems (ICDCS)*, pages 852–863. IEEE, 2021.
- 255 [2] Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How
256 to backdoor federated learning. In *International Conference on Artificial Intelligence and
257 Statistics*, pages 2938–2948. PMLR, 2020.
- 258 [3] Arjun Nitin Bhagoji, Supriyo Chakraborty, Prateek Mittal, and Seraphin Calo. Analyzing feder-
259 ated learning through an adversarial lens. In *International Conference on Machine Learning*,
260 pages 634–643. PMLR, 2019.
- 261 [4] Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. Machine learning
262 with adversaries: Byzantine tolerant gradient descent. In *Proceedings of the 31st International
263 Conference on Neural Information Processing Systems*, pages 118–128, 2017.
- 264 [5] Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. Provably secure federated learning against
265 malicious clients. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35,
266 pages 6885–6893, 2021.
- 267 [6] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm
268 for discovering clusters in large spatial databases with noise. In *Proceedings of the Second
269 International Conference on Knowledge Discovery and Data Mining, KDD’96*, page 226–231.
270 AAAI Press, 1996.
- 271 [7] Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Gong. Local model poisoning attacks
272 to {Byzantine-Robust} federated learning. In *29th USENIX Security Symposium (USENIX
273 Security 20)*, pages 1605–1622, 2020.
- 274 [8] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in
275 the machine learning model supply chain, 2019.
- 276 [9] Rachid Guerraoui, Sébastien Rouault, et al. The hidden vulnerability of distributed learning in
277 byzantium. In *International Conference on Machine Learning*, pages 3521–3530. PMLR, 2018.
- 278 [10] Andrew Hard, Kanishka Rao, Rajiv Mathews, Swaroop Ramaswamy, Françoise Beaufays, Sean
279 Augenstein, Hubert Eichner, Chloé Kiddon, and Daniel Ramage. Federated learning for mobile
280 keyboard prediction. *arXiv preprint arXiv:1811.03604*, 2018.
- 281 [11] Jonathan Hayase, Weihao Kong, Raghav Somani, and Sewoong Oh. Spectre: defending against
282 backdoor attacks using robust statistics. *arXiv preprint arXiv:2104.11315*, 2021.
- 283 [12] Ji Chu Jiang, Burak Kantarci, Sema Oktug, and Tolga Soyata. Federated learning in smart city
284 sensing: Challenges and opportunities. *Sensors*, 20(21):6230, 2020.
- 285 [13] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and
286 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In
287 *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.
- 288 [14] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh,
289 and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv
290 preprint arXiv:1610.05492*, 2016.
- 291 [15] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images.
292 Technical Report 0, University of Toronto, Toronto, Ontario, 2009.
- 293 [16] Ya Le and Xuan S. Yang. Tiny imagenet visual recognition challenge. 2015.
- 294 [17] David Leroy, Alice Coucke, Thibaut Lavril, Thibault Gisselbrecht, and Joseph Dureau. Feder-
295 ated learning for keyword spotting. In *ICASSP 2019-2019 IEEE International Conference on
296 Acoustics, Speech and Signal Processing (ICASSP)*, pages 6341–6345. IEEE, 2019.

- 297 [18] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Chal-
298 lenges, methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020.
- 299 [19] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention
300 distillation: Erasing backdoor triggers from deep neural networks, 2021.
- 301 [20] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust
302 model fusion in federated learning. *NeurIPS*, 2020.
- 303 [21] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Fine-pruning: Defending against
304 backdooring attacks on deep neural networks. In *International Symposium on Research in*
305 *Attacks, Intrusions, and Defenses*, pages 273–294. Springer, 2018.
- 306 [22] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agueria y Arcas.
307 Communication-efficient learning of deep networks from decentralized data. In *Artificial*
308 *intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- 309 [23] Thomas Minka. Estimating a dirichlet distribution, 2000.
- 310 [24] Mustafa Safa Ozdayi, Murat Kantarcioglu, and Yulia R Gel. Defending against backdoors in
311 federated learning with robust learning rate. *arXiv preprint arXiv:2007.03767*, 2020.
- 312 [25] Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,
313 Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *arXiv preprint*
314 *arXiv:2003.00295*, 2020.
- 315 [26] Phillip Rieger, Thien Duc Nguyen, Markus Miettinen, and Ahmad-Reza Sadeghi. Deepsight:
316 Mitigating backdoor attacks in federated learning through deep model inspection. *arXiv preprint*
317 *arXiv:2201.00763*, 2022.
- 318 [27] Victor Sanh, Thomas Wolf, and Alexander M. Rush. Movement pruning: Adaptive sparsity by
319 fine-tuning. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan,
320 and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual*
321 *Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12,*
322 *2020, virtual*, 2020.
- 323 [28] Sidak Pal Singh and Martin Jaggi. Model fusion via optimal transport. *Advances in Neural*
324 *Information Processing Systems*, 33:22045–22055, 2020.
- 325 [29] Stefán Páll Sturluson, Samuel Trew, Luis Muñoz-González, Matei Grama, Jonathan Passerat-
326 Palmbach, Daniel Rueckert, and Amir Alansary. Fedrad: Federated robust adaptive distillation,
327 2021.
- 328 [30] Hongyi Wang, Kartik Sreenivasan, Shashank Rajput, Harit Vishwakarma, Saurabh Agarwal,
329 Jy-yong Sohn, Kangwook Lee, and Dimitris Papailiopoulos. Attack of the tails: Yes, you really
330 can backdoor federated learning. *arXiv preprint arXiv:2007.05084*, 2020.
- 331 [31] Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the
332 objective inconsistency problem in heterogeneous federated optimization. *arXiv preprint*
333 *arXiv:2007.07481*, 2020.
- 334 [32] Yujia Wang, Lu Lin, and Jinghui Chen. Communication-efficient adaptive federated learning.
335 *arXiv preprint arXiv:2205.02719*, 2022.
- 336 [33] Chen Wu, Xian Yang, Sencun Zhu, and Prasenjit Mitra. Mitigating backdoor attacks in federated
337 learning. *arXiv preprint arXiv:2011.01767*, 2020.
- 338 [34] Chulin Xie, Minghao Chen, Pin-Yu Chen, and Bo Li. Crfl: Certifiably robust federated learning
339 against backdoor attacks, 2021.
- 340 [35] Chulin Xie, Keli Huang, Pin-Yu Chen, and Bo Li. Dba: Distributed backdoor attacks against
341 federated learning. In *International Conference on Learning Representations*, 2019.

- 342 [36] Jie Xu, Benjamin S Glicksberg, Chang Su, Peter Walker, Jiang Bian, and Fei Wang. Federated
 343 learning for healthcare informatics. *Journal of Healthcare Informatics Research*, 5(1):1–19,
 344 2021.
- 345 [37] Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. Byzantine-robust distributed
 346 learning: Towards optimal statistical rates. In *International Conference on Machine Learning*,
 347 pages 5650–5659. PMLR, 2018.
- 348 [38] Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated
 349 learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.

350 A Additional Related Work

351 There are a large body of works on federated learning. In this section, we only review the most
 352 relevant works in general FL as well as the backdoor attack and defenses of federated learning.

353 **Federated Learning:** Federated Learning [14] was proposed for improving the communication
 354 efficiency in distributed settings. FedAvg [22] works by averaging local SGD updates, of which the
 355 variants have also been proposed such as SCAFFOLD [13], FedProx [18], FedNova [31]. [25, 32]
 356 proposed adaptive federated optimization methods for better adaptivity. Recently, new aggregation
 357 strategies such as neuron alignment [28] or ensemble distillation [20] has also been proposed.

358 **Backdoor Attacks on Federated Learning:** [2] injects backdoor by predicting the global model
 359 updates and replacing them with the one that was embedded with backdoors. [3] aims to achieve
 360 both global model convergence and targeted poisoning attack by explicitly boosting the malicious
 361 updates and alternatively minimizing backdoor objectives and the stealth metric. [30] shows that
 362 robustness to backdoors implies model robustness to adversarial examples and proposed edge-case
 363 backdoors. DBA [35] decomposes the trigger pattern into sub-patterns and distributing them for
 364 several malicious clients to implant. With known aggregation rules, malicious clients also speculate
 365 the global update direction and modify it for the specified direction for backdoor task [7].

366 **Backdoor Defenses on Federated Learning:** In federated learning, the server does not have access
 367 to the clients’ local data, so some defenses based on data inspection (removing training data with
 368 backdoors) [11] are not applicable. Robust Learning Rate [24] flips the signs of some dimensions of
 369 global updates. [33] designs a federated pruning method to remove redundant neurons for backdoor
 370 defense. [34] proposed a certified defense that exploits clipping and smoothing for better model
 371 smoothness. BAFFLE [1] uses a set of validating clients, refreshed in each training round, to
 372 determine whether the global updates have been subject to a backdoor injection. Recent work [26]
 373 identifies suspicious model updates via clustering-based similarity estimations.

374 B Extensions to Other Network Architectures

375 The flip operation can be similarly extended to other network architectures as the candidate weights
 376 selection (Step 1) and the model training (Step 3) are not relevant to the model architecture at all.
 377 Therefore, we only need to adapt the sign flipping part (Step 2). For CNN, we resize the trigger
 378 and flip the sign of the candidate parameters to maximize the convolution layer’s activation. The
 379 same strategy applies for any dot product based operation (convolution can be seen as a special dot
 380 product). Take MLP as an example, assume the first layer’s weight is $\mathbf{w}^{[1]}$. We flip these weights’
 381 signs by $\mathbf{w}^{[1]} = \mathbf{m}_s^{[1]} \odot \text{sign}(\mathbf{x}'_{in} - \mathbf{x}_{in}) \odot |\mathbf{w}^{[1]}| + (\mathbf{1} - \mathbf{m}_s^{[1]}) \odot \mathbf{w}^{[1]}$, (\mathbf{x}'_{in} and \mathbf{x}_{in} is the flatten input
 382 sample with and without trigger, and the non-zero elements of $\mathbf{x}'_{in} - \mathbf{x}_{in}$ only take place on pixels
 383 with the trigger.) The Equation is similar to Equation 2.5 except that we do not need to resize the
 384 trigger as in CNN. The sign flipping of the rest layers follows the same.

385 C Algorithms

386 C.1 Focused-Flip Federated Backdoor Attack

387 We summarize our F3BA algorithm as pseudo-code in Algorithm 1. Specifically, our F3BA method
 388 searches for candidate parameters (Line 6 - Line 9) based on directional or directionless criteria, such
 389 that it is difficult to remove these backdoor-related parameters by coarse-scale model-refinement or

390 aggregation rules on the server. After each time flipping candidate parameters in a certain layer, the
 391 same validation samples \mathbf{x}_v are required (Line 12) to calculate the activation difference (Line 15) for
 392 flipping candidate parameters in the following layers.

Algorithm 1 Focused-Flip Federated Backdoor Attack (F3BA)

1: **Input:** $\sigma(\cdot)$: activation function
 2: α : learning rate for backdoor training
 3: K : local training iterations
 4: $\mathcal{L}_{\text{train}}$: loss function for backdoor training
 5: Malicious client i receives $\theta_t^{i,0} := \{\mathbf{w}^{[1]}, \mathbf{w}^{[2]}, \dots, \mathbf{w}^{[L]}\}$ from the server. Let us denote each
 layer’s output as $\mathbf{z}^{[1]}(\cdot), \mathbf{z}^{[2]}(\cdot), \dots, \mathbf{z}^{[L]}(\cdot)$.
 6: **for** $j = 1$ **to** L **do**
 7: Compute $\mathcal{S}_t^{[j]}$ based on Equation (2.3) or Equation (2.4)
 8: Compute $\mathbf{m}_s^{[j]}$ as a mask that selects $s\%$ lowest scores in $\mathcal{S}_t^{[j]}$
 9: **end for**
 10: Resize Δ to the size of $\mathbf{w}^{[1]}$ and get Δ^*
 11: $\mathbf{w}^{[1]} = \mathbf{m}_s^{[1]} \odot \text{sign}(\Delta^*) \odot |\mathbf{w}^{[1]}| + (\mathbf{1} - \mathbf{m}_s^{[1]}) \odot \mathbf{w}^{[1]}$
 12: Sample a batch of validation data \mathbf{x}_v from \mathcal{D}_i
 13: **for** $j = 2$ **to** L **do**
 14: $\mathbf{x}'_v = (\mathbf{1} - \mathbf{m}) \odot \mathbf{x}_v + \mathbf{m} \odot \Delta$ //backdoor generation
 15: $\delta = \sigma(\mathbf{z}^{[j-1]}(\mathbf{x}'_v)) - \sigma(\mathbf{z}^{[j-1]}(\mathbf{x}_v))$
 16: Resize δ to the size of $\mathbf{w}^{[j]}$ and get δ^*
 17: $\mathbf{w}^{[j]} = \mathbf{m}_s^{[j]} \odot \text{sign}(\delta^*) \odot |\mathbf{w}^{[j]}| + (\mathbf{1} - \mathbf{m}_s^{[j]}) \odot \mathbf{w}^{[j]}$
 18: **end for**
 19: **for** $k = 1$ **to** K **do**
 20: Sample a batch of training data \mathbf{x}_k from \mathcal{D}_i
 21: $\mathbf{x}'_k = (\mathbf{1} - \mathbf{m}) \odot \mathbf{x}_k + \mathbf{m} \odot \Delta$ //backdoor generation
 22: $\theta_t^{i,k} = \theta_t^{i,k-1} - \alpha \nabla_{\theta} \mathcal{L}_{\text{train}}(\mathbf{x}_k, \mathbf{x}'_k, y_{\text{target}}, \theta_t^{i,k-1})$
 23: **end for**

393 **C.2 Focused-Flip Federated Backdoor Attack With Trigger Optimization**

394 We also summarize the more advanced F3BA with trigger pattern optimization in Algorithm 2. The
 395 major difference compared with Algorithm 1 lies in the trigger optimization part (Line 14 - Line 19),
 396 where the first layer is repetitively flipped (Line 17) based on the signs of trigger pattern in the same
 397 position after each optimization step. After the trigger is optimized, the focused flip of the first layer
 398 is also finished, and the rest part (flipping the following layers) is the same as Algorithm 2.

399 **D Discussion on Selection Criterion of the Candidate Parameters**

400 As shown in Eq. (2.3) and Eq. (2.4), we have two possible criteria (Directional or Directionless) for
 401 the selection of candidate parameters. In our paper, we set the Directional Criteria as the default
 402 setting for F3BA. Though it works well in most cases such as when attacking model-refinement
 403 defenses, we find that for some robust aggregation defenses (e.g., Bulyan), we need to adjust it to fit
 404 better. In this section, we discuss these two criteria: what are the meanings of the two criteria and
 405 how to pick the right one to use in different situations.

406 We first talk about the *Directional Criteria*, which is our default setting.

407 **Directional Criteria** target parameters that are moving significantly far away from 0 (and consider
 408 that as important⁸ weight). To see this, note that we will obtain a large importance score under two
 409 scenarios: (1) when the r -th element in $\mathbf{w}_t^{[j]}$ is increasing, i.e., $\left[\frac{\partial \mathcal{L}_g}{\partial \mathbf{w}_t^{[j]}}\right]_r < 0$, and $\left[\mathbf{w}_t^{[j]}\right]_r > 0$; (2)
 410 when the r -th element in $\mathbf{w}_t^{[j]}$ is decreasing, i.e., $\left[\frac{\partial \mathcal{L}_g}{\partial \mathbf{w}_t^{[j]}}\right]_r > 0$, and $\left[\mathbf{w}_t^{[j]}\right]_r < 0$. Both scenarios

⁸Assume that we train the model from scratch, i.e., a model with all zero parameters. All important weights that significantly affect model accuracy will eventually move to either positive values or negative values.

Algorithm 2 Focused-Flip Federated Backdoor Attack with Trigger Optimization (F3BA-TrigOpt)

```
1: Input:  $\sigma(\cdot)$ : activation function
2:  $\alpha$ : learning rate for backdoor training
3:  $\eta$ : learning rate for trigger optimization
4:  $K$ : local training iterations
5:  $P$ : trigger optimization iterations
6:  $\mathcal{L}_{\text{train}}$ : loss function for backdoor training
7:  $\mathcal{L}_{\text{trig}}$ : loss function for trigger optimization
8: Malicious client  $i$  receives  $\theta_t^{i,0} := \{\mathbf{w}^{[1]}, \mathbf{w}^{[2]}, \dots, \mathbf{w}^{[L]}\}$  from the server. Let us denote each
   layer's output as  $\mathbf{z}^{[1]}(\cdot), \mathbf{z}^{[2]}(\cdot), \dots, \mathbf{z}^{[L]}(\cdot)$ .
9: for  $j = 1$  to  $L$  do
10:   Compute  $\mathbf{S}_t^{[j]}$  based on Equation (2.3) or Equation (2.4)
11:   Compute  $\mathbf{m}_s^{[j]}$  as a mask that selects  $s\%$  lowest scores in  $\mathbf{S}_t^{[j]}$ 
12: end for
13: for  $p = 1$  to  $P$  do
14:   Resize  $\Delta$  to the size of  $\mathbf{w}^{[1]}$  and get  $\Delta^*$ 
15:    $\mathbf{w}^{[1]} = \mathbf{m}_s^{[1]} \odot \text{sign}(\Delta^*) \odot |\mathbf{w}^{[1]}| + (\mathbf{1} - \mathbf{m}_s^{[1]}) \odot \mathbf{w}^{[1]}$ 
16:   Sample a batch of training data  $\mathbf{x}_p$  from  $\mathcal{D}_i$ 
17:    $\mathbf{x}'_p = (\mathbf{1} - \mathbf{m}) \odot \mathbf{x}_p + \mathbf{m} \odot \Delta$  //backdoor generation
18:    $\mathcal{L}_{\text{trig}}(\mathbf{x}_p, \mathbf{x}'_p) = \|\sigma(\mathbf{z}^{[1]}(\mathbf{x}_p)) - \sigma(\mathbf{z}^{[1]}(\mathbf{x}'_p))\|_2^2$ 
19:    $\Delta = \Delta + \eta \cdot \nabla_{\Delta} \mathcal{L}_{\text{trig}}(\mathbf{x}_p, \mathbf{x}'_p)$ 
20: end for
21: for  $j = 2$  to  $L$  do
22:   Sample a batch of validation data  $\mathbf{x}_v$  from  $\mathcal{D}_i$ 
23:    $\mathbf{x}'_v = (\mathbf{1} - \mathbf{m}) \odot \mathbf{x}_v + \mathbf{m} \odot \Delta$  //backdoor generation
24:    $\delta = \sigma(\mathbf{z}^{[j-1]}(\mathbf{x}_v)) - \sigma(\mathbf{z}^{[j-1]}(\mathbf{x}'_v))$ 
25:   Resize  $\delta$  to the size of  $\mathbf{w}^{[j]}$  and get  $\delta^*$ 
26:    $\mathbf{w}^{[j]} = \mathbf{m}_s^{[j]} \odot \text{sign}(\delta^*) \odot |\mathbf{w}^{[j]}| + (\mathbf{1} - \mathbf{m}_s^{[j]}) \odot \mathbf{w}^{[j]}$ 
27: end for
28: for  $k = 1$  to  $K$  do
29:   Sample a batch of training data  $\mathbf{x}_k$  from  $\mathcal{D}_i$ 
30:    $\mathbf{x}'_k = (\mathbf{1} - \mathbf{m}) \odot \mathbf{x}_k + \mathbf{m} \odot \Delta$  //backdoor generation
31:    $\theta_t^{i,k} = \theta_t^{i,k-1} - \alpha \nabla_{\theta} \mathcal{L}_{\text{train}}(\mathbf{x}_k, \mathbf{x}'_k, y_{\text{target}}, \theta_t^{i,k-1})$ 
32: end for
```

411 suggest that $\left[\mathbf{w}_t^{[j]}\right]_r$ is moving away from 0. On the contrary, when the parameter weight and its
412 derivative is the same sign, the criterion regards this parameter as not important for its main task.
413 F3BA exploits these unimportant parameters (moving towards 0) and flips their signs to mount a
414 strong and persistent attack without damaging performance on the main task.

415 Despite having little influence on the main task, this criterion tends to select parameters with the
416 largest absolute values on weights or updates (approximations to the derivatives) as candidates.
417 Generally, it helps backdoor FL systems where a low proportion of malicious clients need to compete
418 with a large number of benign ones, but also be defended by robust aggregation methods that filter
419 extreme weights (updates). Therefore, we also propose the *Directionless Criteria* for such situations.

420 **Directionless Criteria** target parameters with the smallest magnitudes on both their parameter
421 weights and updates. When a parameter's weight (update) is closer to 0 compared with other
422 parameters, it would less likely to be regarded as an outlier or potentially malicious update. When the
423 data are non-i.i.d distributed, the proposed local updates can not reach an agreement on the signs of
424 some coordinates. In this circumstance, smaller $|\mathbf{w}_t^{[j]}|$ and smaller $|\frac{\partial \mathcal{L}_g}{\partial \mathbf{w}_t^{[j]}}|$ separately ensure that the
425 candidate parameters would not be too large or small among all proposed ones before and after being
426 flipped, ensuring stealth of the attack.

427 In summary, as a complement to directional criteria, directionless criteria bypass the robust aggrega-
428 tion defenses based on filtering or changing the extreme values of model parameters.

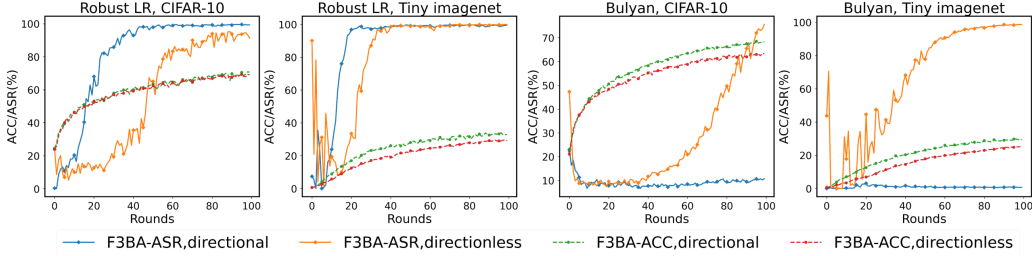


Figure 6: Attack success rate and accuracy against Robust LR and Bulyan on CIFAR-10 and Tiny imagenet datasets with directional and directionless criteria. 20 clients, 4 malicious.

429 In Figure 6, we show two examples: Robust LR and Bulyan, under F3BA attack using directional
 430 and directionless criteria respectively. From Figure 6, we can observe that using directional
 431 in Robust LR helps reach a higher ASR in fewer rounds compared with directionless criteria while
 432 using the same directional criteria in Bulyan directly leads to failure of attack (Directionless Criteria
 433 is needed here). We argue that the gap in attack effectiveness with two criteria may be attributed
 434 to the defense mechanism. For Robust LR, since it does not put constraints on a single proposed
 435 update but the signs of all proposed ones, we can use directional criteria to reach a stronger attack.
 436 For Bulyan, however, the directional criterion makes the updates of flipped parameters become larger,
 437 with a higher probability to be excluded from the aggregation. To keep the malicious updates stealthy
 438 after being flipped, we turn to apply the directionless criteria.

439 E Additional Experimental Details

440 E.1 General Experimental Settings

441 We evaluate the attacks on two classification datasets with non-i.i.d. data distributions: CIFAR-
 442 10 [15] and Tiny-ImageNet [16]. To simulate non-i.i.d. training data and supply the server with
 443 unbalanced samples from each class for model refinement, we divide the training images (in both
 444 datasets) using a Dirichlet distribution [23] with a concentration hyperparameter h . A shared global
 445 model is trained by all participants each round for aggregation. We evaluate CIFAR-10 and Tiny
 446 Imagenet dataset separately with a simple CNN (2 convolutional layers and 2 fully connected layers)
 447 and Resnet-18. Each participating client selected in one round will train for local epochs using SGD
 448 with the learning rate of $\eta = 0.001$ for both CIFAR and Tiny ImageNet. To ensure that the backdoor
 449 trigger is practical and hard to notice by human eyes, we limit the size of the trigger to a small 3×3
 450 square on the CIFAR-10 dataset, and 4×4 for the Tiny ImageNet dataset. For F3BA, we set the
 451 trigger optimization iteration $P = 10$ and $\eta = 0.1$. We apply different candidate parameter selection
 452 proportions 1% and 0.1% respectively for convolutional layers and fully-connected layers.

453 E.2 Discussion On Experimental Settings

454 We first test the effect of our proposed attack on plain FedAvg [22] without any defensive measures.
 455 Specifically, we compare our complete F3BA attack with trigger optimization in Algorithm 2 with a
 456 recent proposed, distributed backdoor attack (DBA) [35]. The main idea of DBA is to partition the
 457 trigger into multiple parts and distribute the backdoor injection task to several malicious clients by
 458 these sub-triggers. This ensures that the resulting malicious updates are less noticeable and more
 459 stealthy compared to standard federated backdoor attacks. We follow the same setting (e.g. number
 460 of total and malicious clients, local training epochs, learning rate η , concentration hyperparameter h)
 461 as the experiment of DBA for the attack of plain FedAvg.

462 As Figure 7, both F3BA and DBA can easily break FedAvg, and F3BA reaches a higher ASR without
 463 noticeable ACC loss. Though quickly observe the difference between centralized and distributed
 464 attacks, we argue that the DBA’s experimental setting strongly facilitate the attack of DBA: all
 465 malicious clients are consistently selected in the total 100 clients and clients’ data is very unevenly
 466 distributed (the concentration hyper-parameter h is 0.1). Benign participants are randomly selected
 467 to form a total of selected 10 participants in each round. The malicious clients train 5 local epochs

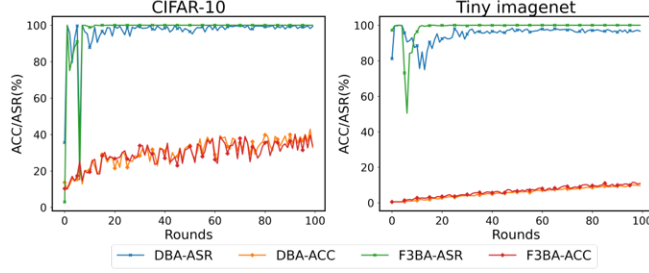


Figure 7: Attack success rate and accuracy of F3BA and DBA against plain FedAvg on CIFAR-10 and Tiny-ImageNet datasets under the setting of DBA experiment.

468 (including both clean and backdoored samples), more than 2 in benign clients. These settings hinder
 469 the accumulation of benign updates thus making them unable to compete with the malicious updates.
 470 We put our control experiments on these elements and the corresponding adjustments in Section ??
 471 and the attack settings of Section 3 to give a preciser and more fair evaluation on the effectiveness of
 472 backdoor defences.

473 E.3 Additional Details on the Defense Baselines

474 **FedDF** [20] leverages unlabeled data or artificially generated samples from a GAN’s generator to
 475 achieve robust server-side model fusion, aggregate knowledge from all received (heterogeneous)
 476 client models. Formally, in the t -th round, the server first obtains a global model θ_{t+1} with standard
 477 FedAvg as equation 2.1.

478 FedDF distills the global model with the ensembled logit predictions on clean unlabeled samples from
 479 each client’s model, hence overcoming the limitation of the number of data samples and eliminating
 480 the effect of data drift.

481 **FedRAD** [29] utilizes median-based scoring along with knowledge distillation for ensemble dis-
 482 tillation. In the t -th round with K participating clients, the median-based scoring assigns the i -th
 483 client a score s_i by counting how many times it gave the prediction for a server-side dataset \mathbf{x}_s . For
 484 a classification task with M probable classes $\mathcal{C} = [c_1, c_2, \dots, c_M]$ and a sample \mathbf{x} in \mathbf{x}_s , the i -th
 485 client’s logits output for c_m is denoted as $f_{\theta_t^{i,K}}(\mathbf{x})[m]$, and we get its score:

$$s_i = \sum_{\mathbf{x} \in \mathbf{x}_s} \sum_{c_m \in \mathcal{C}} \mathbb{1}(f_{\theta_t^{i,K}}(\mathbf{x})[m] = \text{median}(F_{\theta_t^K}(\mathbf{x})[m]))$$

$$\text{where } F_{\theta_t^K}(\mathbf{x})[m] = [f_{\theta_t^{1,K}}(\mathbf{x})[m], f_{\theta_t^{2,K}}(\mathbf{x})[m], \dots, f_{\theta_t^{K,K}}(\mathbf{x})[m]]$$

487 Then FedRAD normalizes $s_i \leftarrow s_i / \sum_{i=1}^K (s_i)$ to let all the clients’ scores adding up to 1, and
 488 aggregates local models. The major difference between FedDF and FedRAD during distillation is that
 489 the median $F_{\theta_t^K}[m]$ is used instead of $\frac{1}{K} \sum_{i \in [K]} f_{\theta_{t+1}^{i,K}}(\mathbf{x}_j)[m]$ to generate the m -th logit of teacher
 490 model’s soft labels.

491 **FedMV Pruning** [33] lets each client provide a ranking of all p filters in its last convolutional layer
 492 based on their averaged activation values and decide which filters would ultimately be pruned in the
 493 global model. Suppose each client has a test dataset \mathbf{x}_{test} , and the activations of its last convolutional
 494 layer is $\sigma(\mathbf{x}_{\text{test}}) = [\sigma_1(\mathbf{x}_{\text{test}}), \sigma_2(\mathbf{x}_{\text{test}}), \dots, \sigma_p(\mathbf{x}_{\text{test}})]$. The i -th client get the ranks of the p filter $\mathcal{R}^{i,K}$
 495 based on the ascending order of $\sigma(\mathbf{x}_{\text{test}})^{i,K}$ (the smallest $\sigma(\cdot)$ means the smallest rank), and sends
 496 $\mathcal{R}^{[i,K]}$ with its local model $\theta^{[i,K]}$ for each round. The server averages the received rankings by filters
 497 $\mathcal{R}^K = \text{avg}(\{\mathcal{R}^{i,K}\}_{1 \leq i \leq K})$. As a result, the server obtains a global ranking for all the filters in the
 498 last convolutional layers.

499 **Bulyan**[9] argues that any current gradient aggregation rules (e.g. Krum, GeoMED, BRUTE), where
 500 the aggregated vector is the result of a distance minimization scheme, can not defend the proposed
 501 malicious updates that is highly-divergent on only a few coordinates while keeping the others closed.
 502 In view of this, more than a combination of these rules, Bulyan bounds the aggregated updates around
 503 the median of all proposed updates at each coordinate and excludes those potentially malicious
 504 updates that disagree a lot. The two steps of Bulyan can be formalized as follows:

505 (1). Choose the updates closest to other updates among the proposed local updates. For pairwise
506 distance, this would be the sum of Euclidean distance to other models. Move the chosen update from
507 the “received set” to the “selection set”, noted \mathcal{S} . Repeat the procedure until get $n - 2f$ updates in \mathcal{S} .

508 (2). Aggregate $n - 4f$ updates closest to the median by coordinate. Hence for each $i \in [1 \dots d]$ (Suppose
509 the uploaded models has d dimensions), The resulting update $\mathbf{G} = \boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1} = (\mathbf{G}[1] \dots \mathbf{G}[d])$, so
510 that for each of its coordinates $\mathbf{G}[\cdot]$:

$$\forall i \in [1, \dots, d], \mathbf{G}[i] = \frac{1}{n - 4f} \sum_{\mathbf{X} \in \mathcal{M}[i]} \mathbf{X}[i]$$

$$\text{where } \mathcal{M}[i] = \arg \min_{\mathcal{R} \subset \mathcal{S}, |\mathcal{R}|=n-4f} \left(\sum_{\mathbf{X} \in \mathcal{R}} |\mathbf{X}[i] - \text{median}[i]| \right) \quad (\text{E.1})$$

$$\text{median}[i] = \arg \min_{m = \mathbf{Y}_i, \mathbf{Y} \in \mathcal{S}} \left(\sum_{\mathbf{Z} \in \mathcal{S}} |m - \mathbf{Z}_i| \right)$$

511 Simply stated: each i -th coordinate of \mathbf{G} equals to the average of the $n - 4f$ closest i -th coordinates
512 to the median of the $n - 2f$ selected updates.

513 **Robust LR** requires a sufficient number of proposed updates with the same signs to decide the
514 global optimization direction. It assumes that the direction of proposed updates for benign clients
515 and malicious clients is in most cases inconsistent, hence the presence of malicious clients would
516 change the distribution of the signs of all proposed local updates. For each dimension with the sum of
517 signs of updates fewer than a pre-defined threshold β , the learning rate is multiplied by -1 . With the
518 number of adversarial agents sufficiently below β , Robust LR is expected to move the global model
519 from the backdoored model to the benign one. Since Robust LR only adjusts the learning rate, the
520 approach is agnostic to the aggregation rules. For example, it can trivially work with update clipping
521 and noise addition.

522 **DeepSight** [26] inspects the output probabilities $f_{\boldsymbol{\theta}^{i,K}}(\mathbf{x}_{\text{rand}}) \in \mathbb{R}^{d_0 \times d_c}$ (d_c is the number of classes)
523 of the i -th local model $\boldsymbol{\theta}^{i,K}$ on d_0 given random d_1 -dimensional inputs $\mathbf{x}_{\text{rand}} \in \mathbb{R}^{d_0 \times d_1}$. After inspect
524 the model’s sample-wise average $(\bar{f}_{\boldsymbol{\theta}^{i,K}}(\mathbf{X}_{\text{rand}})) = \sum_{p=1}^{n_0} (f_{\boldsymbol{\theta}^{i,K}}(\mathbf{X}_{\text{rand}}[p])) \in \mathbb{R}^{d_c}$ and label the
525 potentially malicious clients, DeepSight applies DBSCAN on participant clients [6] three times
526 with distance matrices $D_{\text{bias}}, D_{\text{conv}}, D_{\text{prob}} \in \mathbb{R}^{K \times K}$ (assume all the local models have the same
527 architecture with L layers).

- 528 • $D_{\text{bias}}[i, j] = 1 - \text{cosine}(\mathbf{u}_{t+1}^{i,K,[L]}, \mathbf{u}_t^{j,K,[L]}, \mathbf{u}_{t+1}^{i,K,[L]})$ is the update of the i -th local update in their
529 last layers.
- 530 • $D_{\text{conv}}[i, j] = \|\mathbf{w}_{t+1}^{i,K,[L]} - \mathbf{w}_t^{j,K,[L]}\|$ is the Euclidean distance of the i -th and j -th local models’
531 last layers $\mathbf{w}_{t+1}^{i,K,[L]}$ and $\mathbf{w}_t^{j,K,[L]}$
- 532 • $D_{\text{prob}}[i, j] = \|\bar{f}_{\boldsymbol{\theta}^{i,K}}(\mathbf{x}_{\text{rand}}) - \bar{f}_{\boldsymbol{\theta}^{j,K}}(\mathbf{x}_{\text{rand}})\|$ is the Euclidean distance of clients’ output proba-
533 bilities for the global random vectors.

534 After get three clustering result vectors $Re_{\text{bias}}, Re_{\text{conv}}, Re_{\text{prob}} \in \mathbb{R}^K$ (For example, $Re_{\text{bias}}[i]$ is the i -th
535 local model’s cluster label based on D_{bias}), DeepSight defines a new distance matrix D_{final} as:

$$D_{\text{final}}[i, j] = \sum_{r \in Res} \mathbb{1}(r[i] = r[j]) \text{ where } Res = \{Re_{\text{bias}}, Re_{\text{conv}}, Re_{\text{prob}}\}$$

536 DeepSight performs DBSCAN the last time according to $D_{\text{final}}[i, j]$ and get Re_{final} . Based on Re_{final} ,
537 the cluster with potentially malicious clients more than a given proportion threshold would be
538 excluded from the next round of aggregation.

539 **CRFL**[34] clips the training-time global model parameters $Clip_{\rho_t}(\boldsymbol{\theta}_t) = \boldsymbol{\theta}_t / \max(1, \frac{\|\boldsymbol{\theta}_t\|}{\rho_t})$ so that
540 its norm is bounded by ρ_t , and then add isotropic Gaussian noise $\tilde{\boldsymbol{\theta}}_t \leftarrow Clip_{\rho_t}(\boldsymbol{\theta}_t) + \epsilon_t$, where
541 $\epsilon_t \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})$. Aligned with the training time Gaussian noise (perturbing), CRFL adopts the same
542 Gaussian smoothing measures $\mu(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}, \sigma_T^2 \mathbf{I})$ M times independently on the tested model, to
543 get M sets of noisy model parameters, such that $\tilde{\boldsymbol{\theta}}_T^k \leftarrow Clip_{\rho_t}(\boldsymbol{\theta}_t) + \epsilon_t^k$, runs the classifier with each
544 set of noisy model parameters $\tilde{\boldsymbol{\theta}}_T^k$ for one test sample x_{test} to returns its class counts, with which
545 take the voted most probable class and its probability.

546 F Additional Experimental Details

547 F.1 General Experimental Settings

548 We evaluate the attacks on two classification datasets with non-i.i.d. data distributions: CIFAR-10 [15]
549 and Tiny-ImageNet [16]. To simulate non-i.i.d. training data and supply the server with unbalanced
550 samples from each class for model refinement, we divide the training images (in both datasets) using
551 a Dirichlet distribution [23] with a concentration hyperparameter h (a larger h means a more i.i.d.
552 data distribution). A shared global model is trained by all participants each round for aggregation. We
553 evaluate CIFAR-10 and Tiny Imagenet dataset separately with a simple CNN (2 convolutional layers
554 and 2 fully connected layers) and Resnet-18. Each participating client selected in one round will train
555 for local epochs using SGD with the learning rate of $\eta = 0.001$ for both CIFAR and Tiny ImageNet.
556 To ensure that the backdoor trigger is practical and hard to notice by human eyes, we limit the size
557 of the trigger to a small 3×3 square on the CIFAR-10 dataset, and 4×4 for the Tiny ImageNet
558 dataset. For F3BA, we set the trigger optimization iteration $P = 10$ and $\eta = 0.1$. We apply different
559 candidate parameter selection proportions 1% and 0.1% respectively for convolutional layers and
560 fully-connected layers.

561 F.2 Compare with the DBA’s Experimental Setting

562 Our goal is to use a realistic experimental settings to fairly and accurately evaluate the real-world
563 performance of various federated backdoor defenses in the face of advanced backdoor attacks. We
564 believe that it is unrealistic to specify a fixed number of malicious clients in each training round as in
565 the experimental setup of DBA. In reality, due to the server’s lack of knowledge on the clients (the
566 server cannot know in advance whether a client is benign or malicious), it can only randomly select a
567 subset of clients to participate in each training round (the number of malicious clients is unknown).
568 In this case, it is more practical to randomly select participating clients among all the clients without
569 distinguishing between benign and malicious clients as the experimental setting in F3BA.

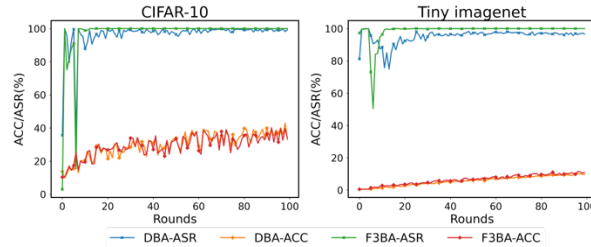


Figure 8: Attack success rate and accuracy of F3BA and DBA against plain FedAvg on CIFAR-10 and Tiny-ImageNet datasets under the setting of DBA.

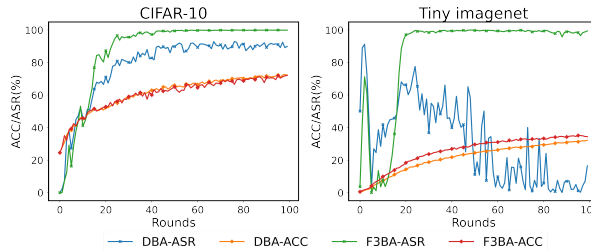


Figure 9: Attack success rate and accuracy of F3BA and DBA against plain FedAvg on CIFAR-10 and Tiny-ImageNet datasets under the setting of F3BA.

570 We attack plain FedAvg with F3BA and DBA in the pre-tuned (in DBA) and post-tuned (in F3BA)
571 experimental settings. Figure.8 and Figure.9 show that even without any defense, DBA could not
572 evade FedAvg on Tiny-Imagenet dataset in the post-tuned setting, while F3BA succeed in both
573 datasets for the two settings. Our chosen setting is actually harder and more practical than the setting

574 of DBA. Since all the clients are randomly selected, as Figure.10 shows, the proportion of malicious
 575 clients among all the participating clients is much less than that in DBA’s setting.

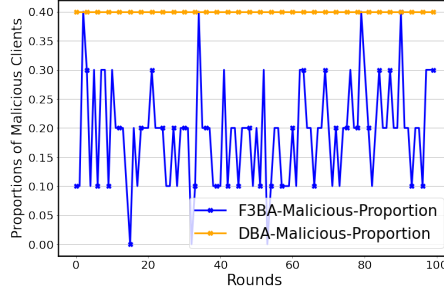


Figure 10: The proportion of malicious clients of each training round in DBA’s and F3BA’s settings.

576 **F.3 Additional Details on the Defense Baselines**

577 **FedDF** [20] leverages unlabeled data or artificially generated samples from a GAN’s generator to
 578 achieve robust server-side model fusion, aggregate knowledge from all received (heterogeneous)
 579 client models. Formally, in the t -th round, the server first obtains a global model θ_{t+1} with standard
 580 FedAvg as equation 2.1

581 FedDF distills the global model with the ensembled logit predictions on clean unlabeled samples from
 582 each client’s model, hence overcoming the limitation of the number of data samples and eliminating
 583 the effect of data drift.

584 **FedRAD** [29] utilizes median-based scoring along with knowledge distillation for ensemble dis-
 585 tillation. In the t -th round with K participating clients, the median-based scoring assigns the i -th
 586 client a score s_i by counting how many times it gave the prediction for a server-side dataset \mathbf{x}_s . For
 587 a classification task with M probable classes $\mathcal{C} = [c_1, c_2, \dots, c_M]$ and a sample \mathbf{x} in \mathbf{x}_s , the i -th
 588 client’s logits output for c_m is denoted as $f_{\theta_t^{i,K}}(\mathbf{x})[m]$, and we get its score:

$$s_i = \sum_{\mathbf{x} \in \mathbf{x}_s} \sum_{c_m \in \mathcal{C}} \mathbb{1}(f_{\theta_t^{i,K}}(\mathbf{x})[m] = \text{median}(F_{\theta_t^K}(\mathbf{x})[m]))$$

$$\text{where } F_{\theta_t^K}(\mathbf{x})[m] = [f_{\theta_t^{1,K}}(\mathbf{x})[m], \dots, f_{\theta_t^{K,K}}(\mathbf{x})[m]]$$

589 Then FedRAD normalizes $s_i \leftarrow s_i / \sum_{i=1}^K (s_i)$ to let all the clients’ scores adding up to 1, and
 590 aggregates local models. The major difference between FedDF and FedRAD during distillation is that
 591 the median $F_{\theta_t^K}[m]$ is used instead of $\frac{1}{K} \sum_{i \in [K]} f_{\theta_t^{i,K}}(\mathbf{x}_j)[m]$ to generate the m -th logit of teacher
 592 model’s soft labels.

593 **FedMV Pruning** [33] lets each client provide a ranking of all p filters in its last convolutional layer
 594 based on their averaged activation values and decide which filters would ultimately be pruned in the
 595 global model. Suppose each client has a test dataset \mathbf{x}_{test} , and the activations of its last convolutional
 596 layer is $\sigma(\mathbf{x}_{\text{test}}) = [\sigma_1(\mathbf{x}_{\text{test}}), \sigma_2(\mathbf{x}_{\text{test}}), \dots, \sigma_p(\mathbf{x}_{\text{test}})]$. The i -th client get the ranks of the p filter $\mathcal{R}^{i,K}$
 597 based on the ascending order of $\sigma(\mathbf{x}_{\text{test}})^{i,K}$ (the smallest $\sigma(\cdot)$ means the smallest rank), and sends
 598 $\mathcal{R}^{[i,K]}$ with its local model $\theta^{[i,K]}$ for each round. The server averages the received rankings by filters
 599 $\mathcal{R}^K = \text{avg}(\{\mathcal{R}^{i,K}\}_{1 \leq i \leq K})$. As a result, the server obtains a global ranking for all the filters in the
 600 last convolutional layers.

601 **Bulyan** [9] argues that any current gradient aggregation rules (e.g. Krum, GeoMED, BRUTE), where
 602 the aggregated vector is the result of a distance minimization scheme, can not defend the proposed
 603 malicious updates that is highly-divergent on only a few coordinates while keeping the others closed.
 604 In view of this, more than a combination of these rules, Bulyan bounds the aggregated updates around
 605 the median of all proposed updates at each coordinate and excludes those potentially malicious
 606 updates that disagree a lot. The two steps of Bulyan can be formalized as follows:

607 (1). Choose the updates closest to other updates among the proposed local updates. For pairwise
 608 distance, this would be the sum of Euclidean distance to other models. Move the chosen update from
 609 the “received set” to the “selection set”, noted \mathcal{S} . Repeat the procedure until get $n - 2f$ updates in \mathcal{S} .

610 (2). Aggregate $n-4f$ updates closest to the median by coordinate. Hence for each $i \in [1\dots d]$ (Suppose
 611 the uploaded models has d dimensions), The resulting update $\mathbf{G} = \boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1} = (\mathbf{G}[1]\dots\mathbf{G}[d])$, so
 612 that for each of its coordinates $\mathbf{G}[\cdot]$:

$$\forall i \in [1, \dots, d], \mathbf{G}[i] = \frac{1}{n-4f} \sum_{\mathbf{X} \in \mathcal{M}[i]} \mathbf{X}[i]$$

$$\text{where } \mathcal{M}[i] = \arg \min_{\mathcal{R} \subset \mathcal{S}, |\mathcal{R}|=n-4f} \left(\sum_{\mathbf{X} \in \mathcal{R}} |\mathbf{X}[i] - \text{median}[i]| \right) \quad (\text{F.1})$$

$$\text{median}[i] = \arg \min_{m=\mathbf{Y}_i, \mathbf{Y} \in \mathcal{S}} \left(\sum_{\mathbf{Z} \in \mathcal{S}} |m - \mathbf{Z}_i| \right)$$

613 Simply stated: each i -th coordinate of \mathbf{G} equals to the average of the $n-4f$ closest i -th coordinates
 614 to the median of the $n-2f$ selected updates.

615 **Robust LR** requires a sufficient number of proposed updates with the same signs to decide the
 616 global optimization direction. It assumes that the direction of proposed updates for benign clients
 617 and malicious clients is in most cases inconsistent, hence the presence of malicious clients would
 618 change the distribution of the signs of all proposed local updates. For each dimension with the sum of
 619 signs of updates fewer than a pre-defined threshold β , the learning rate is multiplied by -1 . With the
 620 number of adversarial agents sufficiently below β , Robust LR is expected to move the global model
 621 from the backdoored model to the benign one. Since Robust LR only adjusts the learning rate, the
 622 approach is agnostic to the aggregation rules. For example, it can trivially work with update clipping
 623 and noise addition.

624 **DeepSight** [26] inspects the output probabilities $f_{\boldsymbol{\theta}^{i,K}}(\mathbf{x}_{\text{rand}}) \in \mathbb{R}^{d_0 \times d_c}$ (d_c is the number of classes)
 625 of the i -th local model $\boldsymbol{\theta}^{i,K}$ on d_0 given random d_1 -dimensional inputs $\mathbf{x}_{\text{rand}} \in \mathbb{R}^{d_0 \times d_1}$. After inspect
 626 the model's sample-wise average ($\bar{f}_{\boldsymbol{\theta}^{i,K}}(\mathbf{X}_{\text{rand}}) = \sum_{p=1}^{n_0} (f_{\boldsymbol{\theta}^{i,K}}(\mathbf{X}_{\text{rand}}[p])) \in \mathbb{R}^{d_c}$ and label the
 627 potentially malicious clients, DeepSight applies DBSCAN on participant clients [6] three times
 628 with distance matrices $D_{\text{bias}}, D_{\text{conv}}, D_{\text{prob}} \in \mathbb{R}^{K \times K}$ (assume all the local models have the same
 629 architecture with L layers).

- 630 • $D_{\text{bias}}[i, j] = 1 - \text{cosine}(\mathbf{u}_{t+1}^{i,K,[L]}, \mathbf{u}_t^{j,K,[L]})$, $\mathbf{u}_{t+1}^{i,K,[L]}$ is the update of the i -th local update in
 631 their last layers.
- 632 • $D_{\text{conv}}[i, j] = \|\mathbf{w}_{t+1}^{i,K,[L]} - \mathbf{w}_t^{j,K,[L]}\|$ is the Euclidean distance of the i -th and j -th local
 633 models' last layers $\mathbf{w}_{t+1}^{i,K,[L]}$ and $\mathbf{w}_t^{j,K,[L]}$
- 634 • $D_{\text{prob}}[i, j] = \|\bar{f}_{\boldsymbol{\theta}^{i,K}}(\mathbf{x}_{\text{rand}}) - \bar{f}_{\boldsymbol{\theta}^{j,K}}(\mathbf{x}_{\text{rand}})\|$ is the Euclidean distance of clients' output
 635 probabilities for the global random vectors.

636 After getting the three clustering result vectors $Re_{\text{bias}}, Re_{\text{conv}}, Re_{\text{prob}} \in \mathbb{R}^K$ (For example, $Re_{\text{bias}}[i]$ is
 637 the i -th local model's cluster label based on D_{bias}), DeepSight defines a new distance matrix D_{final} as:

$$D_{\text{final}}[i, j] = \sum_{r \in Res} \mathbb{1}(r[i] = r[j]), \quad (\text{F.2})$$

$$\text{where } Res = \{Re_{\text{bias}}, Re_{\text{conv}}, Re_{\text{prob}}\}$$

638 DeepSight performs DBSCAN the last time according to $D_{\text{final}}[i, j]$ and get Re_{final} . Based on Re_{final} ,
 639 the cluster with potentially malicious clients more than a given proportion threshold would be
 640 excluded from the next round of aggregation.

641 **CRFL** [34] clips the training-time global model parameters $Clip_{\rho_t}(\boldsymbol{\theta}_t) = \boldsymbol{\theta}_t / \max(1, \frac{\|\boldsymbol{\theta}_t\|}{\rho_t})$ so that
 642 its norm is bounded by ρ_t , and then add isotropic Gaussian noise $\tilde{\boldsymbol{\theta}}_t \leftarrow Clip_{\rho_t}(\boldsymbol{\theta}_t) + \epsilon_t$, where
 643 $\epsilon_t \sim \mathcal{N}(0, \sigma_t^2 \mathbf{I})$. Aligned with the training time Gaussian noise (perturbing), CRFL adopts the same
 644 Gaussian smoothing measures $\mu(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}, \sigma_T^2 \mathbf{I})$ M times independently on the tested model, to
 645 get M sets of noisy model parameters, such that $\tilde{\boldsymbol{\theta}}_T^k \leftarrow Clip_{\rho_t}(\boldsymbol{\theta}_t) + \epsilon_t^k$, runs the classifier with each
 646 set of noisy model parameters $\tilde{\boldsymbol{\theta}}_T^k$ for one test sample x_{test} to returns its class counts, with which
 647 take the voted most probable class and its probability.

648 **G Ablation Study**

649 In this section, we provide ablation studies towards our proposed F3BA method and study how
 650 various factors affect the ASR and ACC of our proposed attack. To eliminate the influence of random
 651 client participation, in this, section, we set a total of 10 participant clients with only 1 malicious client,
 652 and all the participants would be selected in each round. The data heterogeneity hyperparameter
 653 $h = 1.0$ (Except when we investigate the effects of data heterogeneity). We do not set more clients or
 654 lower data heterogeneity because it would make it easier to achieve high ASR for F3BA with various
 655 parameter settings, which is not conducive for us to investigate the effects of various factors on the
 656 backdoor attack.

657 **G.1 Attack with Different Trigger**

658 In our attack design, the malicious clients do not need to share the optimized trigger during training,
 659 instead, they optimize their own triggers when conducting the F3BA attack. During the test phase,
 660 we simply average all the optimized triggers as a global trigger and attached it to the test dataset for
 661 ASR evaluation. Note that even if we do not perform averaging but directly use one of the optimized
 662 triggers, the attack still works.

663 Table.1 shows the ASR of different triggers on CIFAR-10 and Tiny-Imagenet datasets. Note that
 664 whether using the average or local triggers does not have a major impact on the performance of F3BA
 665 on both datasets.

Tested Tigger	Round	CIFAR-10	Tiny-imagenet
Averaged Trigger	50	97.97%	97.01%
	100	99.23%	99.15%
Local Trigger #1	50	98.03%	96.36%
	100	99.49%	99.11%
Local Trigger #2	50	98.57%	96.57%
	100	99.49%	98.95%
Local Trigger #3	50	97.87%	98.15%
	100	98.35%	99.62%
Local Trigger #4	50	98.84%	96.57%
	100	98.41%	98.48%

Table 1: ASR of F3BA and DBA attacks with the averaged global trigger and clients’ local triggers against plain FedAvg.

666 **G.2 Attack Other Network Architecture**

667 As mentioned in Section.B, the Focused Flip operations can be generally applied to any network
 668 structure that rely on the dot product, and be applied independently to boost the traditional training-
 669 based backdoor attack. We show the ASR/ACC of F3BA on attacking plain FedAvg with MLP
 670 models to verify the applicability of the F3BA attack on architectures beyond CNN. Table.2 suggests
 671 that F3BA is still highly effective on MLPs (and still better than the DBA baseline) without loss on
 the performance of its main task.

Round	CIFAR-ACC		CIFAR-ASR		TinyImagenet -ACC		TinyImagenet -ASR	
	F3BA	DBA	F3BA	DBA	F3BA	DBA	F3BA	DBA
25	44.13%	44.57%	98.76%	87.21%	7.06%	7.02%	98.52%	89.54%
50	47.17%	47.26%	99.78%	93.30%	8.58%	8.87%	97.37%	86.35%
75	48.53%	49.13%	99.82%	91.55%	9.40%	9.50%	98.39%	90.39%
100	51.10%	50.99%	99.88%	93.68%	10.04%	10.00%	98.94%	93.54%

Table 2: ASR/ACC of F3BA and DBA with MLP network architecture against plain FedAvg.

673 **G.3 Attack with More Benign Clients**

674 We test F3BA on FedAvg with in total 100 clients, of which 4 are malicious. 40 clients are randomly
 675 selected for each round. When the number of malicious clients is fixed, the benign clients becomes
 676 more thus to some extent deterring the attack from both F3BA and DBA. As Table.3, we can observe
 677 that the advantages of F3BA over DBA still holds. When all the clients are randomly selected, the
 678 decreasing chances for malicious clients being selected do partially slows down the process of F3BA’s
 679 evasion but not able to remove the injected backdoor. In this circumstance, F3BA still performs better
 680 than DBA. Stopping F3BA entirely by the number of benign clients would require potentially much
 more benign clients and fewer malicious clients.

Round	CIFAR-ACC		CIFAR-ASR		TinyImagenet -ACC		TinyImagenet -ASR	
	F3BA	DBA	F3BA	DBA	F3BA	DBA	F3BA	DBA
25	39.01%	38.10%	12.39%	10.12%	9.25%	9.64%	7.15%	4.72%
50	46.54%	45.93%	34.26%	20.27%	17.50%	17.23%	18.10%	9.27%
75	51.21%	50.95%	56.20%	25.55%	20.05%	20.66%	42.06%	15.20%
100	55.38%	55.60%	75.25%	24.36%	23.84%	24.20%	60.11%	30.05%

Table 3: ASR/ACC with more benign clients against plain FedAvg. 100 clients, 4 malicious.

681 **G.4 Attack Sparsification-based Defense**

683 Besides 3 Model Refinement defenses, 3 Robust Aggregation defenses, and 1 Certified Robustness
 684 defenses in Section.3, we also explore the effect of SparseFed, a theoretical framework for analyzing
 685 the robustness of defenses against poisoning attacks. Since it is not specifically designed for the
 686 robustness towards backdoor attacks, we put the result in the supplementary as Table.4.

687 Based on the model architecture and task complexity, we set the number of accepted parameters
 688 $K = 1e4$ for CNN and $K = 4e5$ for ResNet-18 respectively. From the Table.4, we can observe that
 689 although SparseFed only allows a small fraction of aggregated parameters for global model updates
 at each round, it still can be backdoored by our F3BA.

Round	CIFAR		TinyImagenet	
	ACC	ASR	ACC	ASR
25	44.70%	76.75%	1.75%	9.64%
50	49.63%	90.51%	4.78%	49.56%
75	54.43%	99.30%	8.05%	87.21%
100	57.71%	99.75%	11.10%	82.77%

Table 4: ASR/ACC with against sparseFed FedAvg. 20 clients, 4 malicious.

690 We conjecture that the sparsity criterion cannot rule out all the backdoor-related model parameters as
 691 our attack does not necessarily lead to an update that is small in magnitude.

693 **G.5 Effect of Data Heterogeneity**

694 We study how data heterogeneity would affect the attack of F3BA. We manually adjust the concentra-
 695 tion hyperparameter h to split non-i.i.d dataset. (The larger the h , the more i.i.d the data is distributed).
 696 On the two datasets, the ASR both grows when the h becomes smaller, and the ACC decreases at the
 697 same time. The result shows that lower h strongly hurts the accuracy of the global model but gives
 698 convenience to the backdoor attack.

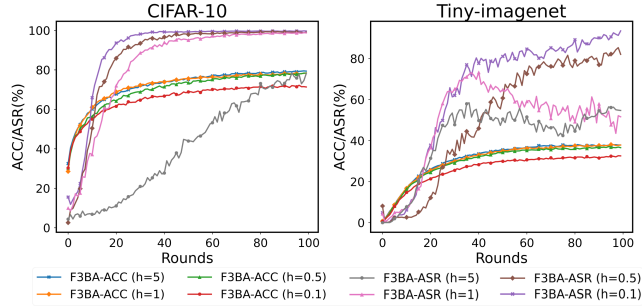


Figure 11: Attack success rate and accuracy against FedAvg under different data heterogeneity h .

699 **G.6 Effect of The Proportion of Candidate Parameters**

700 In our experiments, we find that different proportions of candidate parameters should be used for the
 701 different parts of a neural network in order to achieve a better attack. To valid this, we conducted
 702 a grid search on CIFAR-10 for ACC and ASR at the 20th round. Based on the result in Figure 12,
 703 moderately scaling up the candidate parameters rate (e.g. from 0.02 to 0.05) in fully-connected
 704 layers would increase the ASR without the loss of ACC. Meanwhile, too high candidate parameters
 705 proportion for fully-connected layers can cause an obvious loss of ASR. Intuitively, the parameters of
 706 the convolutional part are more sparse than the fully-connected part and therefore can be selected
 707 with a higher candidate parameter proportion. Similarly, fully-connected layers indeed involve more
 708 parameters, and thus only require a smaller candidate parameter proportion. According to our practice,
 709 a sound choice for attackers is to set the proportions for convolutional layers and fully-connected
 710 layers respectively below 5% and 1%.

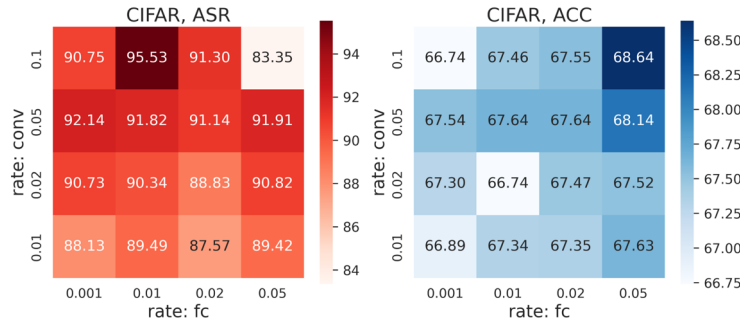


Figure 12: Attack success rate and accuracy with different candidate parameter proportion for convolutional layers and fully-connected layers on CIFAR-10 dataset.

711 **G.7 Effect of Local Training**

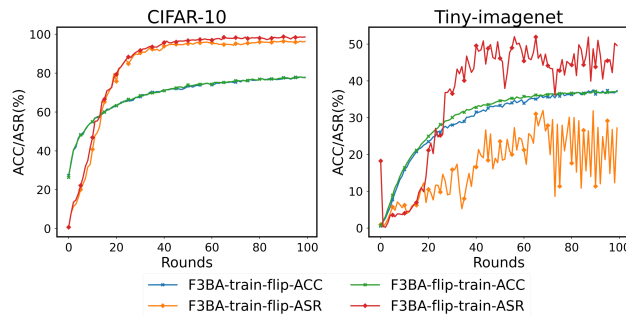


Figure 13: Attack success rate and accuracy against FedAvg under different data heterogeneity h .

712 As discussed in Section 2.2, local training is still a must for the effects of F3BA, while simply
713 flipping without training can not induce activation led by the trigger to the targeted label. We further
714 discover how the order of local training and Focused Flip would affect our attack. Based on the
715 results on CIFAR-10 and Tiny Imagenet, the flipping-training-pipeline can achieve better ASR than
716 the flip-training one. It also ensures a slightly higher ACC. The significant difference of ASR on
717 attacking Tiny Imagenet dataset for two pipelines also suggests that training a local model to bridge
718 the sudden changes in model weights caused by focused flip can be of benefit to the effectiveness of
719 the attack.