

# Report on

## LeadFL: Client Self-Defense against Model Poisoning in Federated Learning

Claus Guthmann, Ramazan Tan

**Abstract**—Existing defense mechanisms in federated learning often assume an approximately constant proportion of malicious clients throughout training. However, when only a subset of clients participates in each communication round, the fraction of malicious clients can vary substantially across rounds. Over a sufficient number of iterations, this variability makes it likely that some rounds will contain a majority of malicious participants. Under this premise, the paper *LeadFL: Client Self-Defense against Model Poisoning in Federated Learning* [1] demonstrates that poisoning attacks occurring in such rounds can have long-lasting effects on the global model. To address this issue, the authors propose a client-side defense mechanism that significantly improves the model’s ability to recover from these transient but severe poisoning events.

### I. INTRODUCTION

FEDERATED learning enables the collaborative training of machine learning models across distributed clients while keeping local data private. However, this decentralized setting also exposes the training process to model poisoning attacks, where malicious clients manipulate their updates to degrade the global model. Many existing defense mechanisms implicitly assume that the proportion of malicious clients remains approximately constant throughout training. In practice, federated learning systems typically select only a subset of clients to participate in each communication round, causing the fraction of malicious clients to fluctuate over time. As a result, it becomes increasingly likely that some rounds contain an unusually high—or even majority—share of malicious participants, especially over long training horizons.

The paper *LeadFL: Client Self-Defense against Model Poisoning in Federated Learning* [1] addresses this often-overlooked scenario. It shows that poisoning attacks occurring during such adversarial spikes can have long-lasting effects on the global model, even when the system later returns to benign conditions. To mitigate this issue, LeadFL introduces a client-side defense mechanism designed to limit the impact of these transient but severe attacks and to accelerate recovery in subsequent rounds.

In this report, we begin by presenting the preliminary concepts and results required to understand the LeadFL approach. We then provide an overview of the method itself and discuss its potential advantages in federated learning environments subject to spiking adversarial participation, followed by an examination of its limitations. To further assess its effectiveness, we implemented the LeadFL algorithm and extended the original evaluation to cover previously unaddressed cases. Finally, we present the results of our own experiments, which corroborate that LeadFL significantly improves robustness against the attack patterns considered in the paper.

### II. BACKGROUND

A central idea is to reduce the lasting influence of a burstial attack. One way to formally characterize this lingering effect was introduced in *FL-WBC: Enhancing Robustness against Model Poisoning Attacks in Federated Learning from a Client Perspective* [2] through the notion of the **Attack Effect on Parameters**. For a given communication round  $t$ , let  $W_t$  denote the global model parameters resulting from potentially malicious client updates, and let  $W_t^M$  denote the parameters that would have been obtained in the absence of malicious behavior. The attack effect is then defined as

$$\delta_t := W_t - W_t^M. \quad (1)$$

The authors further derive a recursive approximation of this quantity,

$$\delta_t \approx \hat{\delta}_t := \frac{1}{|C_t|} \left[ \sum_{c \in C_t} \prod_{i=0}^{I-1} (I - \eta_t H_{t,i}^c) \right] \hat{\delta}_{t-1},$$

where  $C_t$  denotes the set of clients selected in round  $t$  (assumed to be equally weighted),  $\eta_t$  is the learning rate,  $I$  is the number of local training steps, and  $H_{t,i}^c$  represents the Hessian matrix of the local loss function.

Empirically, the Hessian matrices  $H_{t,i}^c$  are often highly sparse, which causes the aggregated transformation

$$\frac{1}{|C_t|} \left[ \sum_{c \in C_t} \prod_{i=0}^{I-1} (I - \eta_t H_{t,i}^c) \right] \quad (2)$$

to remain close to the identity matrix. As a consequence, the attack effect  $\delta_t$  decays only slowly over time, leading to the observed persistence of poisoning attacks.

The FL-WBC [2] client-side defense seeks to counteract this by injecting random noise into local updates, thereby accelerating the decay of the attack effect. Building on this insight, LeadFL [1] proposes a more targeted approach: instead of relying on random perturbations, it introduces a specifically designed regularization term.

Directly computing the full Hessian matrix  $H_{t,i}^c$  for neural networks is, however, computationally prohibitive. To make the approach tractable in a federated learning setting, LeadFL relies on an efficient approximation of the Hessian inspired by the *Optimal Brain Damage* framework [3].

Specifically, the method restricts attention to the diagonal elements of the Hessian, thereby ignoring inter-parameter dependencies. While this simplification sacrifices second-order interactions, it significantly reduces computational complexity and has been shown to provide useful curvature information in large-scale models. The diagonal entries capture the sensitivity of the gradient with respect to individual parameters and can be approximated via finite differences of successive gradients.

Formally, the Hessian matrix is approximated as

$$H_{t,i}^c \approx \tilde{H}_{t,i}^c := \text{diag}(\nabla L(W_{t,i+1}^c) - \nabla L(W_{t,i}^c)), \quad (3)$$

where  $W_{t,i}^c$  denotes the local model parameters of client  $c$  at global round  $t$  and local step  $i$ , and  $L(\cdot)$  is the local loss function.

### III. LEADFL

The core idea of LeadFL [1] is to reduce the persistence of poisoning effects by locally minimizing the term  $I - H_{t,i}^c$  on each client, thereby accelerating the decay of the attack effect described in the previous section. To achieve this, LeadFL augments the standard local optimization objective with an additional regularization term whose gradient is incorporated into each update step.

Concretely, each local update is computed in two stages. First, a standard gradient descent step is performed,

$$\tilde{W}_{t,i+1}^c := W_{t,i}^c - \eta_t \nabla L(W_{t,i}^c), \quad (4)$$

yielding intermediate model parameters  $\tilde{W}_{t,i+1}^c$ . In a second step, the regularization update is applied,

$$W_{t,i+1}^c := \tilde{W}_{t,i+1}^c - \eta_t \alpha \text{clip}(\nabla(I - \eta_t \tilde{H}_{t,i}^c), q), \quad (5)$$

where  $\alpha$  controls the strength of the regularization and  $\text{clip}$  limits the absolute value of each entry to the clipping threshold  $q$ .

This two-stage formulation allows the intermediate parameters  $\tilde{W}_{t,i+1}^c$  to be reused in constructing a computationally efficient approximation of the Hessian. However, the derivation of this approximation in the original paper contains errors. We therefore adopt the following formulation, which we believe to be a mathematically consistent approximation closely aligned with the intended method:

$$\tilde{H}_{t,i}^c \approx (W_{t,i}^c - \tilde{W}_{t,i+1}^c + \Delta W_{t,i}^c) / \eta_t, \quad (6)$$

where  $\Delta W_{t,i}^c := W_{t,i}^c - W_{t,i-1}^c$  denotes the parameter change between successive local steps. It is important to note that this Hessian approximation is only accurate up to the contribution of the regularization term applied in the previous local update.

Combining the previous steps, the local client iteration of LeadFL can be summarized as follows (cf. Algorithm 1 in [1]):

---

#### Algorithm 1 *LeadFL* Local Client Iteration

---

- 1: Compute gradients and update intermediate weights:  
 $\tilde{W}_{t,i+1}^c \leftarrow W_{t,i}^c - \eta_t \nabla L(W_{t,i}^c)$
  - 2: Estimate the Hessian matrix:  
 $\tilde{H}_{t,i}^c \leftarrow (W_{t,i}^c - \tilde{W}_{t,i+1}^c + \Delta W_{t,i}^c) / \eta_t$
  - 3: Compute the regularization term:  
 $R_{t,i}^c \leftarrow \text{clip}(\nabla(I - \eta_t \tilde{H}_{t,i}^c), q)$
  - 4: Update the local weights with the regularization term:  
 $W_{t,i+1}^c \leftarrow \tilde{W}_{t,i+1}^c - \eta_t \alpha R_{t,i}^c$
- 

The authors of LeadFL subsequently implement and evaluate the proposed defense under a range of experimental settings. Their primary setup consists of 100 clients, of which

25 are malicious. Training is performed over 80 global rounds, with each client executing 10 local training steps per round. In each global round, 10 clients are selected uniformly at random; this client selection is kept consistent across different experimental runs to ensure comparability (additional results for alternative client selection strategies are reported in the appendix of [1]).

The evaluation considers several server-side aggregation defenses, namely SparseFed, Multi-Krum, and Bulyan. For the LeadFL client-side defense, the clipping threshold is fixed at  $q = 0.2$ . The regularization strength  $\alpha$  is varied across datasets, with  $\alpha = 0.4$  for FashionMNIST and  $\alpha = 0.25$  for CIFAR10. Due to the randomized client selection process, the number of malicious participants fluctuates across rounds, leading to significant variability in backdoor success rates. To account for this effect, the authors report both the final backdoor accuracy and the average backdoor accuracy over the entire training process.

TABLE I  
COMPARISON OF DEFENSES UNDER 9-PIXEL PATTERN BACKDOOR ATTACK ON IID FASHIONMNIST DATASET. (CF. TABLE 1 [1])

Server-side Defense	Multi-Krum			
Client-side Defense	None	LDP	FL-WBC	LeadFL
Maintask Accuracy	89.3	87	87.2	87.9
Backdoor Accuracy Avg.	82.6	76.0	77.5	32.9
Backdoor Accuracy Final	93.2	79.6	80.6	0.0

TABLE II  
COMPARISON OF DEFENSES UNDER 9-PIXEL PATTERN BACKDOOR ATTACK ON IID CIFAR10 DATASET. (CF. TABLE 2 [1])

Server-side Defense	Multi-Krum			
Client-side Defense	None	LDP	FL-WBC	LeadFL
Maintask Accuracy	76.3	48.0	43.3	56.9
Backdoor Accuracy Avg.	77.5	53.1	56.9	35.6
Backdoor Accuracy Final	80.5	43.8	40.5	25.6

Tables I and II present a selection of results for the IID FashionMNIST and CIFAR10 datasets, respectively, under a 9-pixel pattern backdoor attack using the Multi-Krum server-side defense.

For the FashionMNIST dataset (Table I), LeadFL consistently reduces the final backdoor accuracy to nearly zero while maintaining the main task accuracy within a 10% reduction compared to the undefended baseline. In contrast, the alternative client-side defenses, LDP and FL-WBC, achieve only moderate reductions in backdoor accuracy.

For the more complex CIFAR10 dataset (Table II), the impact of LeadFL on the main task accuracy is more pronounced, reducing it to approximately 50–60%. Nevertheless, LeadFL still substantially lowers the final backdoor accuracy to around 25%, outperforming the other client-side defenses in terms of robustness.

Across both datasets and all evaluated server-side defenses, the absence of a client-side defense results in consistently high average backdoor accuracies—exceeding 75%—highlighting the severity of bursty poisoning attacks and the necessity of additional client-side protection mechanisms.

Fig. 1. 9-pixel back door attack on FashionMNIST when the attack is periodic(cf. Figure 3 [1])

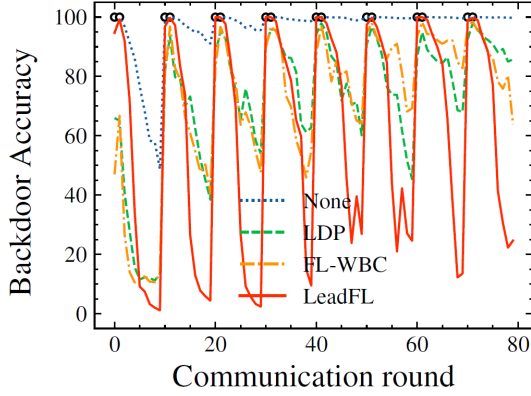
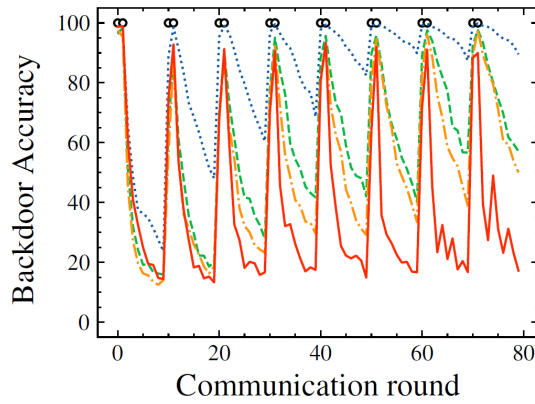


Fig. 2. 9-pixel back door attack on CIFAR10 when the attack is periodic (cf. Figure 3 [1])



The authors further investigate the behavior of LeadFL under explicitly bursty attack patterns. In this setting, client selection is no longer random. Instead, training proceeds in cycles of ten rounds, where the first two rounds contain six malicious clients out of ten selected participants, while the remaining eight rounds involve exclusively benign clients.

The results of this experiment are illustrated in Figure 1 for the FashionMNIST dataset and in Figure 2 for the CIFAR10 dataset. While LeadFL exhibits limited impact during the rounds dominated by malicious clients, its effect becomes apparent in the subsequent benign rounds. In particular, the recovery from each attack spike is significantly faster, and the achieved backdoor accuracy stabilizes at a lower level compared to the undefended baseline.

In summary, LeadFL introduces a novel client-side defense mechanism that is particularly effective at mitigating the long-lasting effects of spikes in malicious client participation. Although the method incurs a measurable reduction in main task accuracy, it demonstrates that client-side defenses can play a crucial role in complementing existing server-side aggregation strategies. These findings motivate a closer examination of LeadFL’s practical limitations and empirical behavior, which we address in the remainder of this report through an extended experimental evaluation.

#### IV. LIMITATIONS AND EXTENDED EVALUATION

Before presenting our own experimental evaluation, we first discuss several limitations and ambiguities encountered in the original LeadFL paper. These issues are relevant both for the correct interpretation of the reported results and for the reproducibility of the proposed method.

First, while the regularization strength  $\alpha$  plays a crucial role in balancing robustness against performance degradation, the paper does not provide a principled justification or systematic tuning strategy for the selected values. Instead,  $\alpha$  is varied empirically across experiments without clear guidance, limiting the interpretability and generalizability of the reported outcomes.

Fig. 3. 9-pixel back door attack on FashionMNIST when the malicious client selection is periodic

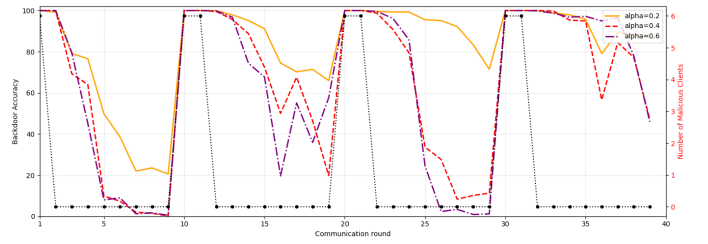
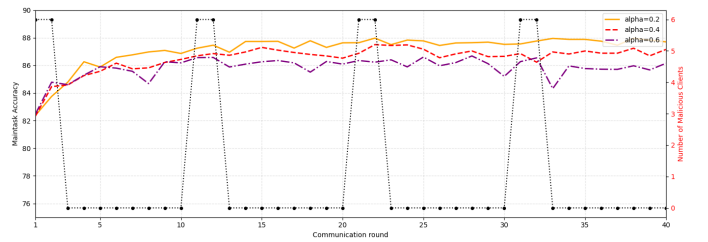


Fig. 4. 9-pixel back door attack on FashionMNIST when the malicious client selection is periodic



To investigate the sensitivity of LeadFL to the choice of  $\alpha$ , we conducted additional experiments using the authors’ original code with modified regularization strengths. Due to computational constraints, these experiments were limited to 40 global rounds. Figures 3 and 4 show the resulting backdoor accuracy and main task accuracy, respectively. All of the evaluations in this section use the bulky server-side defense.

As shown in Figure 3, reducing the regularization strength to  $\alpha = 0.2$  significantly weakens the defense: the decline in backdoor accuracy is slower and can not reach the same lower bound compared to the default setting of  $\alpha = 0.4$ . Increasing the regularization strength to  $\alpha = 0.6$  yields a slight improvement on average, but the effect is not consistent across all rounds. Regarding the main task performance, illustrated in Figure 4,  $\alpha = 0.2$  provides no observable benefit, while  $\alpha = 0.6$  introduces occasional drops in accuracy.

Finally, as discussed in the previous sections, the Hessian approximation used in LeadFL appears to contain fundamen-

tal inconsistencies between the theoretical proposal and the reference implementation.

The manuscript describes the defense mechanism as a direct regularization of the client-side Hessian matrix, explicitly deriving an approximation based on the finite difference of gradients. Under this formulation, the client is expected to perform a lookahead step to estimate the local curvature  $\tilde{H}$  and essentially perturb the loss landscape to eliminate the null-space blind spots where the attack hides.

However, our analysis of the authors' official source code reveals a different approach. As detailed in **Algorithm 2**, the reference implementation does not compute the finite difference of gradients or estimate the Hessian diagonal. Rather than computing the Hessian approximation, the code employs a heuristic that minimizes the  $L_2$  norm of the divergence between the current update and a theoretical "unit step." This effectively constrains the weight trajectory through geometry rather than the curvature-based regularization proposed in the text. This raises the question of whether the method's reported efficacy stems from this norm-based clipping rather than the intended Hessian regularization.

---

**Algorithm 2** *LeadFL* Local Client Iteration (Original Implementation)

---

- 1: Compute standard gradients:  
 $g_{t,i}^c \leftarrow \nabla L(W_{t,i}^c)$
  - 2: Perform standard SGD step:  
 $\tilde{W}_{t,i+1}^c \leftarrow W_{t,i}^c - \eta_t g_{t,i}^c$
  - 3: Compute regularization loss ( $L_2$  Norm):  
 $L_{reg} \leftarrow \alpha \cdot \|\tilde{W}_{t,i+1}^c - W_{t,i}^c - g_{old}^c - \eta_t \cdot \mathbf{1}\|_2$
  - 4: Compute gradients of the regularization loss:  
 $g_{reg}^c \leftarrow \nabla L_{reg}(\tilde{W}_{t,i+1}^c)$
  - 5: Clip the regularization gradients:  
 $\tilde{g}_{reg}^c \leftarrow \text{clip\_grad\_value}(g_{reg}^c, q)$
  - 6: Secondary update step:  
 $W_{t,i+1}^c \leftarrow \tilde{W}_{t,i+1}^c - \eta_t \tilde{g}_{reg}^c$
- 

To address this ambiguity and rigorously evaluate the paper's theoretical claims, we developed a corrected implementation, detailed in **Algorithm 3**. Unlike the original code, our version strictly adheres to the mathematical derivation provided in the manuscript, explicitly computing the gradient difference  $(\tilde{g} - g)/\eta$  to approximate the Hessian diagonal. By testing both versions, we aim to determine if the defense's robustness relies on the intended Hessian regularization or merely on the norm constraint found in the code.

To evaluate this, we reproduced the periodic attack setup shown in Figure 1 and augmented it with results obtained from our own implementation. The corresponding main task accuracy and backdoor accuracy are shown in Figures 5 and 6, respectively. We observe that the original implementation acts as a strict geometric constraint, whereas our version relies on a noisy finite-difference approximation. This distinction likely contributes to the performance variance, suggesting that the strict norm constraint in the reference code may offer additional stability distinct from the proposed Hessian regularization. Although our implementation performs slightly

---

**Algorithm 3** *LeadFL* Local Client Iteration (Our Implementation)

---

- 1: Compute gradients at current parameters:  
 $g_{t,i}^c \leftarrow \nabla L(W_{t,i}^c)$
  - 2: Perform temporary SGD step (lookahead):  
 $\tilde{W}_{t,i+1}^c \leftarrow W_{t,i}^c - \eta_t g_{t,i}^c$
  - 3: Compute gradients at the lookahead position:  
 $\tilde{g}_{t,i+1}^c \leftarrow \nabla L(\tilde{W}_{t,i+1}^c)$
  - 4: Estimate the Hessian diagonal via Finite Difference:  
 $\tilde{H}_{t,i}^c \leftarrow (\tilde{g}_{t,i+1}^c - g_{t,i}^c)/\eta_t$
  - 5: Compute and clip the regularization term:  
 $R_{raw} \leftarrow 1 - \eta_t \tilde{H}_{t,i}^c$   
 $R_{t,i}^c \leftarrow \text{clamp}(R_{raw}, -q, q)$
  - 6: Final update with regularization:  
 $W_{t,i+1}^c \leftarrow \tilde{W}_{t,i+1}^c - \eta_t \alpha R_{t,i}^c$
- 

Fig. 5. 9-pixel back door attack on FashionMNIST when the malicious client selection is periodic

TODO: Replace with proper client selection

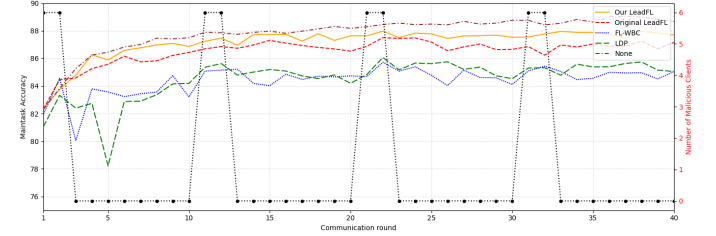
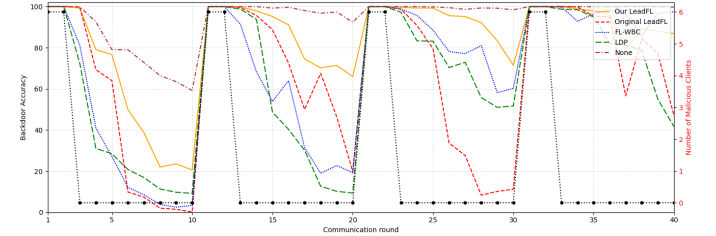


Fig. 6. 9-pixel back door attack on FashionMNIST when the malicious client selection is periodic

TODO: Replace with proper client selection



worse than the original results reported by the authors, the central qualitative behavior of LeadFL remains intact. In particular, the defense induces only a moderate reduction in main task accuracy while substantially improving the recovery rate following bursty poisoning attacks.

Fig. 7. 9-pixel back door attack on FashionMNIST when the malicious client selection is random

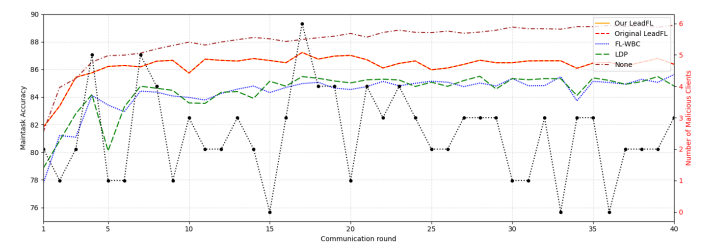
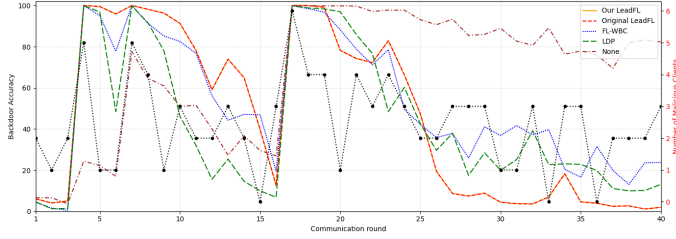


Fig. 8. 9-pixel back door attack on FashionMNIST when the malicious client selection is random



We also repeated the same evaluation using random client selection instead of the periodic pattern. Interestingly, the main-task accuracy, shown in Figure 7, remains largely consistent with the periodic case. In contrast, the backdoor accuracy exhibits more pronounced differences. While LeadFL still provides the strongest defense toward the end of training, it is occasionally outperformed by other methods—including the baseline with no client-side defense—during the initial rounds.

In summary, our evaluation confirms that LeadFL is effective at its core: it accelerates recovery from bursty poisoning attacks and lowers the long-term backdoor accuracy, though at a cost to main-task performance. Crucially, the observation that our corrected Hessian-based implementation exhibits comparable qualitative behavior suggests that the defense’s efficacy is primarily driven by the theoretical regularization, though the implicit norm constraint in the original reference code appears to contribute additional stability. Notably, even with this trade-off, LeadFL maintains higher main-task accuracy compared to alternative client-side defenses. The periodic bursty attack pattern considered in this report represents an optimal scenario for LeadFL, whereas random client selection reduces its overall impact and introduces greater variability in backdoor mitigation. These findings suggest that while LeadFL shows promise, further evaluation is necessary before deployment in real-world federated learning systems. In particular, large-scale experiments and a more systematic investigation of the regularization parameter  $\alpha$  in complex models would be essential.

## REFERENCES

- [1] C. Zhu, S. Roos, and L. Chen. LeadFL: Client Self-Defense against Model Poisoning in Federated Learning. *Proceedings of the 40th International Conference on Machine Learning*, vol. 202, pp. 43158–43180, Jul. 2023.
- [2] J. Sun, A. Li, L. DiValentin, A. Hassanzadeh, Y. Chen and H. Li. FL-WBC: Enhancing Robustness against Model Poisoning Attacks in Federated Learning from a Client Perspective *Advances in Neural Information Processing Systems*, vol. 34, pp. 12613–12624, Dec. 2021.
- [3] Y. LeCun, J. Denker and S. Solla. Optimal Brain Damage *Advances in Neural Information Processing Systems*, vol. 2, pp. 598–605, Nov. 1989.

## V. SUMMARY: MODEL POISONING ATTACKS TO FEDERATED LEARNING VIA MULTI-ROUND CONSISTENCY

Many existing attacks on federated learning rely on assumptions that are often unrealistic in practice, such as knowledge of other clients' updates or awareness of the aggregation rule used by the server. The paper *Model Poisoning Attacks to Federated Learning via Multi-Round Consistency* [2] proposes an attack strategy that relaxes these assumptions and operates effectively without direct knowledge of other clients or the aggregation mechanism.

The principal idea behind the attack is to overcome the natural cancellation of updates caused by random noise added by benign clients. To achieve this, the malicious clients first select a consistent attack direction  $s$ . They then estimate the honest update by removing their own previous malicious contribution from the global model update.

$$\tilde{h} = g^{t-1} - \frac{\|g^{t-1}\|}{\|k^{t-1} \odot s\|} (k^{t-1} \odot s) \quad (7)$$

where  $g^{t-1}$  is the last global model update and  $k^{t-1} \odot s$  the last model updates on the fake clients. The factor  $\frac{\|g^{t-1}\|}{\|k^{t-1} \odot s\|}$  normalizes the malicious update to have the same magnitude as the global one. This estimate is normalized to form a template of a benign update.

$$v^t = \frac{|\tilde{h}|}{\|\tilde{h}\|_2} \quad (8)$$

$$\lambda_t = c_t \|g^{t-1}\|_2 \quad (9)$$

A scaling factor  $\lambda_t$  is determined relative to the norm of the global update, with an additional check: if the global model does not drift in the intended attack direction, it is assumed that the malicious updates are being filtered, and the scaling factor is reduced accordingly:  $c_t = \begin{cases} c_{t-1} & \text{Alignment} > 50\% \\ \beta c_{t-1} & \text{else} \end{cases}$

with  $\beta < 1$ . The final malicious update  $g_i^t$  is computed as the product of the normalized template, the scaling factor, and the attack direction.

$$g_i^t = (v^t \cdot \lambda_t) \odot s \quad (10)$$

The presenting group [1] conducted their own evaluation of this attack. They observed that the attack requires several rounds to become effective. Nevertheless, with at least 9% of selected clients acting maliciously, the attack successfully degrades the global model's accuracy.

The students further investigated the performance of a novel server-side defense, PRoDIGY, which was not included in the original paper. PRoDIGY proved extremely effective against this attack because the malicious updates are highly similar, making them easily detectable and filtered by the dissimilarity-based criteria.

To examine potential circumvention strategies, the group adapted the attack to counter PRoDIGY. They tested two techniques proposed in the original paper: noise adaptation, where random perturbations are added to the malicious update, and sign flipping, where each element of the malicious direction

vector is flipped with a small probability. For both techniques, different scaling factors and probabilities were evaluated. The noise adaptation strategy was over all more effective than sign flipping. However, in all cases, the main task accuracy remained above 80%.

The group also noted that the original paper assumes a fixed, very small learning rate. They experimented with a learning rate scheduler, which accelerates the increase of main task accuracy and causes the learning rate to decay before the attack can fully take effect. Under this condition, the attack's effectiveness is further reduced.

Building on the work of the group, we believe that an additional direction for future research is to investigate the impact of client selection on attack effectiveness. In the paper, the percentage of malicious clients is fixed after selection, remaining consistent across rounds. Since one of the main observations is that the attack's effect accumulates over time, it would be valuable to study the impact of a fully random client selection process, where the fraction of malicious participants varies from round to round.

Further, if attackers are assumed to coordinate, an alternative strategy would be to partition the dimensions of the malicious direction vector among participants, with each attacker contributing only to their assigned subset of parameters. Against defenses such as PRoDIGY, which rely on detecting similarity across malicious updates, this approach could reduce the artificial correlation among adversarial updates and potentially increase the attack's effectiveness.

## REFERENCES

- [1] Y. Xie, M. Fang and N. Z. Gong. Model Poisoning Attacks to Federated Learning via Multi-Round Consistency. *Trustworthy Distributed Learning, TUM*, Jan. 2026.
- [2] L. Kunze, F. Durchdewald and N. Schenk. Presentation on: PoisonedFL via Multi-Round Consistency. *Tri*, vpp. 5454-15463, Jun. 2025.

## VI. SUMMARY: TOWARDS THE ROBUSTNESS OF DIFFERENTIALLY PRIVATE FEDERATED LEARNING

A common intuition in Federated Learning (FL) is that Differential Privacy (DP) noise naturally enhances robustness by "drowning out" malicious inputs. The paper *Towards the Robustness of Differentially Private Federated Learning* [2] challenges this belief, demonstrating that DP-FL is not inherently robust. In fact, standard outlier detection mechanisms (like Krum or Median) fail in DP-FL because the heavy DP noise increases the variance of benign updates, masking the malicious ones.

To exploit this, the authors propose **Attack-DPFL**. The core insight is that while benign updates are perturbed by noise with a large norm ( $\|n\|_2 \approx \sigma\sqrt{d}$ ), malicious clients can upload unperturbed, "clean" poisoned gradients that are amplified to match this large norm.

$$A_u = A \cdot G_u, \quad \text{where } A = \frac{\|S_{benign}\|}{\|G_u\|} \quad (11)$$

By scaling the poisoned gradient  $G_u$  by factor  $A$ , the attack hides the malicious update within the magnitude of the noise, allowing it to dominate the aggregation without being detected as an outlier.

To counter this, the paper introduces **Robust-DPFL**. The defense relies on the observation that DP noise vectors are independent and zero-mean, meaning their element-wise sum tends toward zero. In contrast, the amplified poisoned updates introduce a systematic bias. The defense calculates a *Detection Score*  $Z(S)$  for each update:

$$Z(S) = \left| \frac{1}{d} \sum_{i=1}^d S[i] \right| \quad (12)$$

Updates with high detection scores are identified as malicious and filtered out via K-means clustering.

The presenting group [1] reproduced the experiments on MNIST, FEMNIST, and CIFAR-10. Their evaluation confirmed that Attack-DPFL achieves an Attack Success Rate (ASR) of nearly 100% against standard defenses like Krum and Median, which struggle to distinguish noise from poison.

They further analyzed the "Three-way Trade-off" between Privacy, Accuracy, and Robustness using the Robust-DPFL defense. They observed that while Robust-DPFL successfully lowers ASR, increasing the privacy level (lowering the privacy budget  $\epsilon$ ) introduces a side effect: the intense noise distorts gradient distributions so heavily that the score-based detection becomes slightly less precise. Despite this, they concluded that Robust-DPFL achieves the optimal balance, maintaining high main task accuracy where other defenses fail.

The presenters have outlined critical areas for future investigation, specifically emphasizing **Adaptive Privacy** strategies. While the current approach uses static noise levels, a dynamic mechanism that adjusts in real-time to the estimated threat severity could better preserve model accuracy. Additionally, the presenters highlight the risk of "Smarter Attackers" capable of circumventing the defense by optimizing their updates to

yield a low Detection Score  $Z(S)$ . Countering such adversaries will necessitate the development of more sophisticated defense metrics that go beyond simple element-wise averaging.

## REFERENCES

- [1] K. Ise, W. Tang and R. Yagi. Group Project: Towards the Robustness of Differentially Private Federated Learning. *Trustworthy Distributed Learning, TUM*, Jan. 2026.
- [2] T. Qi, H. Wang and Y. Huang. Towards the Robustness of Differentially Private Federated Learning. *AAAI Conference on Artificial Intelligence (AAAI-24)*, Vol. 38, No. 19, pp. 19911-19919, Mar. 2024.