

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), η_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 δ_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 α controls the strength of the regularization and 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 α 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])

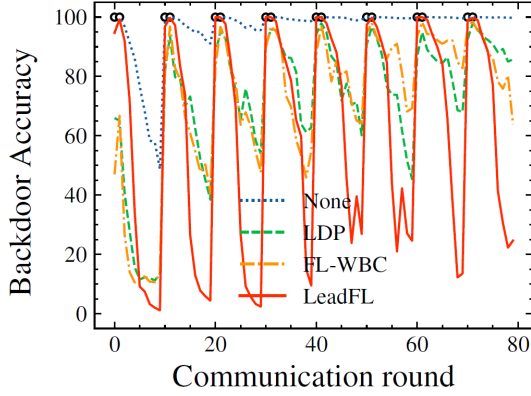
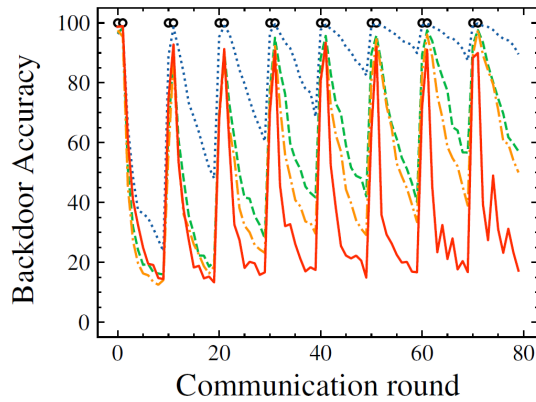


Fig. 2. 9-pixel back door attack on CIFAR10 when the attack is periodic (cf. Figure 3 [1])



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, as discussed in the previous sections, the Hessian approximation used in LeadFL is not derived rigorously and appears to contain minor mathematical inconsistencies. As a result, any practical implementation—including our own—necessarily relies on an approximation that deviates from the formulation presented in the paper.

Second, we observed discrepancies between the algorithmic description in the paper and the behavior required for a functioning implementation.

Finally, while the regularization strength α 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, α is varied empirically across experiments without clear guidance, limiting the interpretability and generalizability of the reported outcomes.

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: