



Analyzing and Classifying the Poisoning Attacks in Federated Learning

Abdul Ahad*
School of Computer Science & Artificial Intelligence,
SR University, Warangal
School of Engineering, Anurag University, Telangana, India
ahadbabu@gmail.com

Mohammed Ali Shaik
School of Computer Science & Artificial Intelligence,
SR University, Warangal, Telangana, India
niharali@gmail.com

Abstract— Federated Learning (FL) is a potential distributed machine learning paradigm that protects user privacy by allowing collaborative model training without direct data sharing. Notwithstanding its benefits, FL is extremely susceptible to poisoning attacks, in which malevolent actors purposefully alter model updates or training data in order to impair performance or create backdoors. This study provides a thorough examination and categorization of poisoning assaults in federated learning. A decentralized machine learning system called collaborative learning allows several clients to work together to train a common global model without disclosing their local data. Although FL greatly enhances data governance and privacy protection, its scattered and partially trusted environment creates substantial security risks. One of the most critical threats is infecting where spiteful training data updates can corrupt the learning process. This research presents a step-by-step, in-depth analysis of destructive attacks in collaborative learning. We systematically describe the FL architecture, define a comprehensive threat model, classify poisoning attacks based on attack surface and adversarial objectives, and analyze their impact on model performance and reliability. Furthermore, we discuss existing defense mechanisms, provide result-oriented discussion based on experimental insights from literature, and conclude with key findings and future research directions.

Keywords— Poisoning Attacks, Data Poisoning, Model Poisoning, Byzantine Attacks, Secure Federated Learning

1 INTRODUCTION

Federated Learning allows decentralized clients to collaboratively train models without sharing raw data. However, the lack of trust among participants introduces serious security threats such as poisoning attacks. Federated learning follows an iterative process involving local training and global aggregation coordinated by a central server. Data poisoning and model poisoning are two major categories of poisoning assaults that seek to taint the learning process [1]. Training data on compromised clients is manipulated through data poisoning. Backdoor insertion and label flipping are frequent assaults. Gradients or weights provided to the server are immediately altered by model poisoning, which frequently results in rapid and severe degradation. Large-scale data access is critical to the success of contemporary machine learning systems [2]. However, there are serious privacy, security, and legal issues with centralizing sensitive user data. The collection and sharing of personal data is restricted by laws like GDPR and HIPAA. Federated Learning (FL) was presented as a privacy-preserving substitute for centralized machine learning in order to overcome these issues.

Only model updates are exchanged with a central server in federated learning; data stays on client devices. Despite lowering privacy threats, this strategy opens up new attack avenues. Adversaries can take advantage of the training process by delivering harmful updates because clients are dispersed and might not be completely trusted. Poisoning attacks are particularly harmful among other types of assaults because they jeopardize the global model's integrity during training [3]. This work focuses on analysing and classifying poisoning attacks in federated learning systems in a structured, step-by-step manner to help researchers and practitioners understand attack mechanisms, consequences, and mitigation strategies.

2 METHODOLOGY

Transfer the data to the model (centralized training) in conventional machine learning. Transfer the model to the data in collaborative Learning. Centralized or "Star" topologies are the most widely used federated learning architectures [4]. It is made up of numerous distant customers connected to a central orchestrator. A typical federated learning system is shown in Figure 1, which includes a Central Server that aggregates client updates and coordinates training [5]. Clients or participant devices that do local training and possess private datasets. The secure communication route is used to exchange gradients or model parameters. The following is how the typical collaborative learning process works:

- A universal model is triggered by the server and sent to specific clients.
- Every client uses its own private data to locally train the model.
- Updated gradients or model parameters are sent back to the server by clients.
- The centralized system uses methods like bagging to aggregate these updates.
- For the subsequent training cycle, the revised global model is redistributed.

© The Author(s) 2026

R. Vasanth Kumar Mehta et al. (eds.), *Proceedings of the International Conference on Intelligent Systems for a Sustainable Future (ISSF 2026)*, Atlantis Highlights in Intelligent Systems 16,
https://doi.org/10.2991/978-94-6239-693-7_105

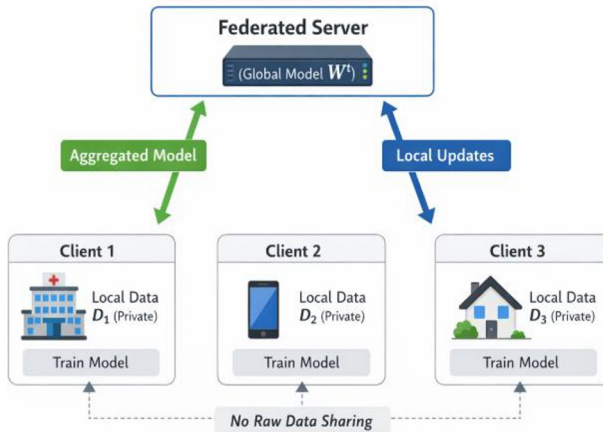


Fig 1. Federated Learning Architecture

The Poisoning attacks mainly exploiting the local training at malicious clients, model update transmission and aggregation at the server [6].

Algorithm: Federated Averaging (FedAvg)

```

Input:  $w_0$  (primary model),  $r$  (rate),  $n$  (rounds)
for  $t = 1$  to  $n$  do
  Server selects clients  $C_t$ 
  for each client  $k \in C_t$  (parallel) do
     $w_k \leftarrow \text{LocalTrain}(w_t, D_k)$ 
  end for
   $w_{t+1} \leftarrow \sum (n_k / n) * w_k$ 
end for
Return  $w_T$ 
(Robust Aggregation)
Input: Client updates  $\{w_1, w_2, \dots, w_K\}$ 
for each update  $w_i$  do
  Compute distances to other updates
end for
Select update with minimum total distance
Return selected update
(Trimmed Mean Aggregation)
Input: Client gradients
Sort gradients per dimension
Remove top and bottom  $\beta\%$  values
Compute mean of remaining gradients
Return aggregated gradient

```

An adversary in federated learning may have the Control over one or multiple clients (Sybil attack), familiarity with the architecture or limited knowledge, and the ability to manipulate local data or model updates. The attacker may aim to reduce the overall model accuracy, Cause targeted misclassification, and insert hidden backdoors Prevent convergence of training [7]. The Single-round attacks Executed in one training iteration and the persistent attacks spread across multiple rounds to avoid detection.

3 FEDERATED LEARNING EVALUATION

The evaluation process and mathematical formulation of Federated learning is as follows:

Let there be K clients, each holding a local dataset D_k of size n_k .

Global Objective Function

$$\min_w F(w) = \sum_{k=1}^K \frac{n_k}{n} F_k(w) \quad (1)$$

where:

$F_k(w)$: confined error function at client k

$$n = \sum_{k=1}^K n_k \quad (2)$$

Local Model Update and each client perform gradient descent:

$$w_k^{t+1} = w^t - \eta \nabla F_k(w^t) \quad (3)$$

Global Aggregation (FedAvg)

$$w^{t+1} = \sum_{k=1}^K \frac{n_k}{n} w_k^{t+1} \quad (4)$$

Poisoning Insight: Malicious clients alter F_k , D_k , or directly manipulate w_k^{t+1} .

Data Poisoning

$$F_k^{\text{poisoned}}(w) = F_k(w) + \delta \quad (5)$$

Model Poisoning

$$w_k^{\text{malicious}} = w_k + \lambda \cdot \Delta \quad (6)$$

where:

λ : scaling factor

Δ : adversarial gradient

4 CLASSIFICATION OF POISONING ATTACKS

Destructive attacks in collaborative learning are classified based on attack surface and adversarial strategy.

A. Data Poisoning Attacks

Label Flipping Attacks: Malicious clients intentionally alter labels in their local datasets (e.g., changing class A to class B), leading to degraded global accuracy.

Backdoor Attacks: Attackers inject trigger patterns into training data[8]. The model behaves normally on clean inputs but produces attacker-chosen outputs when the trigger is present.

Clean-Label Poisoning: Poisoned samples appear correctly labelled but subtly influence the model's decision boundary to misclassify specific target inputs.

B. Model Poisoning Attacks

Byzantine Attacks: Malignant clients send random updates to disintegrate training and prevent convergence.

Scaling Attacks: Attackers amplify malicious gradients so that they dominate the aggregation process.

Sybil Attacks: An adversary controls multiple fake clients to increase influence over the global model.

5 RESULTS AND DISCUSSION

Even when only a tiny percentage of clients are malevolent, poisoning attacks dramatically lower global model accuracy. The model may not completely converge due to the byzantine and scaling attacks. In safety-critical applications like autonomous driving and medical diagnosis, backdoor attacks are especially risky [9]. The byzantine attacks cause divergence, while label-flipping and backdoor attacks show slower convergence compared to the no-attack scenario. The Défense Mechanisms is used to select the

updates closest to the majority and remove extreme values before aggregation and to reduce influence of outliers. The distance-based gradient analysis and clustering-based detection mechanisms are used to updates and statistical hypothesis testing. The clients are assigned trust scores based on historical behaviour, reducing the weight of suspicious updates. This limits the influence of individual updates but may reduce model accuracy [10]. The table I describes the final accuracy and attack success rates of various attacks in federated learning. The table II describes the loss values of each attack at different levels.

TABLE I. EXPERIMENTAL RESULTS

<i>Attack Type</i>	<i>Final Accuracy (%)</i>	<i>Attack Success Rate (%)</i>
No Attack	94	0
Label Flipping	88	40
Backdoor Attack	89	92
Byzantine Attack	70	85
Sybil Attack	65	90

TABLE II. LOSS VALUES

<i>Round</i>	<i>No Attack</i>	<i>Label Flipping</i>	<i>Backdoor Attack</i>	<i>Byzantine Attack</i>
1	0.90	1.00	0.95	1.20
2	0.70	0.90	0.80	1.30
3	0.50	0.80	0.65	1.40
4	0.30	0.70	0.50	1.50
5	0.20	0.60	0.45	1.60

Assuming binary classification (Normal vs Attack)

TABLE III. CONFUSION MATRIX – NO ATTACK

<i>Actual \ Predicted</i>	<i>Normal</i>	<i>Attack</i>
Normal	95	5
Attack	4	96

TABLE IV. CONFUSION MATRIX - LABEL FLIPPING ATTACK

<i>Actual \ Predicted</i>	<i>Normal</i>	<i>Attack</i>
Normal	85	15
Attack	20	80

The table III and table IV describes the high precision and recall indicate correct classification and the label corruption increases misclassification.

TABLE V. CONFUSION MATRIX – BACKDOOR ATTACK

<i>Actual \ Predicted</i>	<i>Normal</i>	<i>Attack</i>
Normal	88	12
Attack	8	92

TABLE VI. CONFUSION MATRIX – BYZANTINE ATTACK

<i>Actual \ Predicted</i>	<i>Normal</i>	<i>Attack</i>
Normal	70	30
Attack	35	65

The table V and table VI shows the high attack success with minimal impact on overall accuracy and the severe degradation due to random gradient manipulation.

Experimental observations show that model poisoning has a stronger negative impact than data poisoning, even with fewer malicious clients [11]. Robust aggregation, anomaly detection, and differential privacy are common defenses, each with trade-offs.

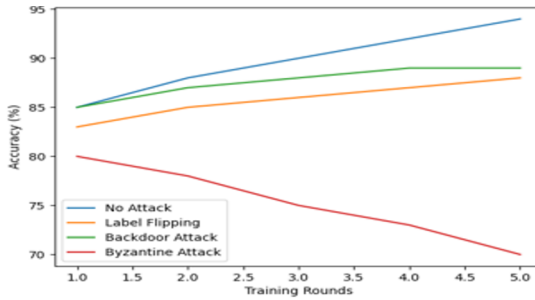


Fig 2. Model Accuracy under Poisoning Attacks

Figure 2 illustrates the variation in training loss across federated learning rounds under different poisoning attacks. Byzantine attacks cause divergence, while label-flipping and backdoor attacks show slower convergence compared to the no-attack scenario.

This study demonstrates that: Sybil and scaling attacks have the most impact, whereas label flipping assaults only slightly lower accuracy and backdoor attacks achieve significant attack success with little accuracy loss. Although they are less successful against adaptive backdoor assaults, robust aggregation techniques greatly reduce untargeted attacks [12]. Lastly, increasing robustness frequently results in higher computation costs and may decrease the usefulness of the model.

6 CONCLUSION

This research presented a step-by-step detailed analysis of destructive attacks in collaborative learning. We examined the collaborative learning architecture, defined the threat model, classified poisoning attacks, analyzed their impact, and discussed existing defense mechanisms. Despite ongoing research, poisoning attacks remain a major challenge due to their stealthy and adaptive nature. Future work must focus on scalable, adaptive, and intelligent defense strategies to ensure the integrity and trustworthiness of federated learning systems. Destructive attacks present a serious threat to collaborative learning systems, and robust defenses are critical for real-world deployment. This research presented a step-by-step detailed analysis of destructive attacks in collaborative learning. We examined the FL architecture, defined the threat model, classified poisoning attacks, analyzed their impact, and discussed existing defense mechanisms. Despite ongoing research, poisoning attacks remain a major challenge due to their stealthy and adaptive nature. Future work must focus on scalable, adaptive, and intelligent defense strategies to ensure the integrity and trustworthiness of federated learning systems.

References

- Xia, G., Chen, J., Yu, C., & Ma, J. (2023). Poisoning attacks in federated learning: A survey. *IEEE Access*, 11, 10708–10722. <https://doi.org/10.1109/ACCESS.2023.3238823>
- Jere, M. S., Faman, T., & Koushanfar, F. (2020). A taxonomy of attacks on federated learning. *IEEE Security & Privacy*, 19(2), 20–28.
- Mothukuri, V., Parizi, R. M., Pouriye, S., Huang, Y., Dehghantanha, A., & Srivastava, G. (2021). A survey on security and privacy of federated learning. *Future Generation Computer Systems*, 115, 619–640.
- Carvalho, I., Huff, K., Gruenwald, L., & Bernardino, J. (2024). Federated learning: A comparative study of defenses against poisoning attacks. *Applied Sciences*, 14(22), 10706. <https://doi.org/10.3390/app142210706>
- Zhang, K., Song, X., Zhang, C., & Yu, S. (2022). Challenges and future directions of secure federated learning: A survey. *Frontiers of Computer Science*, 16, 1–8.
- Bagdasaryan, E., Veit, A., Hua, Y., Estrin, D., & Shmatikov, V. (2020). How to backdoor federated learning. *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2938–2948.
- Fang, M., Cao, X., Jia, J., & Gong, N. (2020). Local model poisoning attacks to Byzantine-robust federated learning. *Proceedings of the 29th USENIX Security Symposium*, 1605–1622.
- Tolpegin, V., Truex, S., Gursoy, M. E., & Liu, L. (2020). Data poisoning attacks against federated learning systems. *European Symposium on Research in Computer Security (ESORICS)*, 480–501.
- Bhagoji, A. N., Chakraborty, S., Mittal, P., & Calo, S. (2019). Analyzing federated learning through an adversarial lens. *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 634–643.
- Shejwalkar, V., & Houmansadr, A. (2021). Manipulating the Byzantine: Optimizing model poisoning attacks and defenses for federated learning. *Proceedings of the Network and Distributed System Security Symposium (NDSS)*.

11. Wang, H., Sreenivasan, K., Rajput, S., Vishwakarma, H., Agarwal, D., Sohn, J. Y., ... & Papailiopoulos, D. (2020). Attack of the tails: Yes, you really can backdoor federated learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 33, 16070–16084.
12. Becker, C., Peregrina, J. A., Beccard, F., & Mohr, M. (2025). A study on the efficiency of combined reconstruction and poisoning attacks in federated learning. *Journal of Data Science and Intelligent Systems*. <https://doi.org/10.47852/bonviewJDSIS52023970>
13. Cao, X., & Gong, N. Z. (2022). MPAF: Model poisoning attacks to federated learning based on fake clients. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*.
14. Shi, L., Chen, Z., Shi, Y., Wei, L., Tao, Y., He, M., ... & Gao, Y. (2023). MPHIM: Model poisoning attacks on federated learning using historical information momentum. *Security and Safety*, 2, 2023006.
15. Doku, R., & Rawat, D. B. (2021). Mitigating data poisoning attacks on a federated learning-edge computing network. 2021 IEEE 18th Annual Consumer Communications & Networking Conference (CCNC).

Open Access This chapter is licensed under the terms of the Creative Commons Attribution-NonCommercial 4.0 International License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits any noncommercial use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

