

## Federated Learning's Dynamic Defense Against Byzantine Attacks: Integrating SIFT-Wavelet and Differential Privacy for Byzantine Grade Levels Detection

Sahithi Godavarthi<sup>1,2\*</sup>, G. Venkateswara Rao<sup>3</sup>

<sup>1</sup>Research Scholar, Dept. of CSE, GITAM School of Technology, GITAM(Deemed to be University), Visakhapatnam

<sup>2</sup>Assistant Professor, Department of Emerging Technologies, CVR College of Engineering, Hyderabad

\* Corresponding Author Email: [sahithi.godavarthi@gmail.com](mailto:sahithi.godavarthi@gmail.com) ORCID: 0000-0002-6611-2799

<sup>3</sup>Professor, Dept. of CSE, GITAM School of Technology, GITAM(Deemed to be University), Visakhapatnam

Email: [venkateswararao.gurrala@gitam.edu](mailto:venkateswararao.gurrala@gitam.edu) ORCID: 0000-0001-6090-339X

### Article Info:

DOI: 10.22399/ijcesen.538  
Received : 21 October 2024  
Accepted : 23 October 2024

### Keywords:

Federated Learning,  
Byzantine attacks,  
Distributed Learning,  
Neural Network,  
Robust and Dynamic Aggregation

### Abstract:

Federated learning, which enables decentralized training across multiple devices while maintaining data privacy, is susceptible to Byzantine poisoning attacks. This paradigm reduces the need for centralized data storage and transmission, thereby mitigating privacy risks associated with traditional data aggregation methods. However, FL introduces new challenges, notably susceptibility to Byzantine poisoning attacks, where rogue participants can tamper with model updates, threatening the consistency and security of the aggregated model. Our approach addresses this vulnerability by implementing robust aggregation methods, sophisticated preprocessing techniques, and a novel Byzantine grade-level detection mechanism. We introduce a federated aggregation operator designed to mitigate the impact of malicious clients. Our preprocessing includes data loading and transformation, data augmentation, and feature extraction using SIFT and wavelet transforms. Additionally, we employ differential privacy and model compression to improve the robustness and performance of the federated learning framework. Our approach is assessed using a tailored neural network model applied to the MNIST dataset, achieving 97% accuracy in detecting Byzantine attacks. Our results demonstrate that robust aggregation significantly improves the resilience and performance. This comprehensive approach ensures the integrity of the federated learning process, effectively filtering out adversarial influences and sustaining high accuracy even when faced with adversarial Byzantine clients.

## 1. Introduction

Poisoning attacks pose a critical risk to Federated Learning (FL) as they can drastically undermine training effectiveness with minimal effort. In these scenarios, malicious actors infiltrate one or more data contributors, enabling them to introduce fraudulent data or alter model updates during the training process. These tainted updates can disrupt the training pipeline and diminish the overall model accuracy[1]. To address these vulnerabilities in FL networks, a variety of strategies have been developed. These include norm and weight control, which monitor the updates to ensure they remain within expected bounds, and distance-based analysis, which identifies outliers in the updates[2]. Additionally, performance-based metrics can assess the reliability of updates, and encryption

mechanisms can secure communications to prevent tampering. These techniques collectively help isolate and neutralize the influence of malicious nodes and falsified reports, thereby strengthening the resilience of FL systems against Byzantine attacks[3].

Federated Learning has become a game-changing method for protecting user data privacy when training Machine Learning (ML) models[4]. The capacity of this cutting-edge method to generate effective ML prediction models without the requirement to centralize sensitive data has drawn a lot of interest. FL drastically lowers the danger of data leaks and ensures privacy by training models locally on devices and sharing just the model updates[5]. Applications including healthcare, finance, and personal devices where data privacy is crucial benefit greatly from this decentralized approach[6]. By avoiding centralized data storage,

FL minimizes communication overhead and reduces the associated privacy risks. Consequently, FL enhances the feasibility of deploying ML in diverse and distributed environments, making it a powerful tool for developing robust[7], privacy-preserving ML models across various industries and use cases. Identifying dishonest nodes in FL, effectively achieved in contrasting the distances between their reports[8]. Traditional arithmetic averaging of model updates is often susceptible to being skewed by fraudulent data, which has the potential to seriously impair the global model's performance. To address this, an enhanced model aggregation method utilizing the geometric median has been introduced [9]. This method provides a more stable variation of gradient descent, which is less influenced by outliers and malicious updates. By calculating the geometric median of model updates, this approach robustly aggregates contributions from different nodes, thereby minimizing the impact of adversarial attacks. As a result, the aggregated model becomes more reliable, leading to more accurate and resilient training outcomes in FL systems. This geometric median-based aggregation is crucial for preserving the global model's performance and integrity in the face of Byzantine attacks[10].

Furthermore, researchers have explored advanced techniques such as to improve the security of FL systems, private pairwise distance computations are used for participant selection and key sharing[11]. These methods involve securely sharing cryptographic keys among participants and selecting trustworthy nodes based on calculated distances between their updates[12]. While these techniques show significant promise, they also present notable implementation challenges, particularly in resource-constrained environments where computational and communication resources are limited. Despite these challenges, such approaches offer potential solutions for enhancing the robustness of FL against adversarial threats. By carefully selecting participants and ensuring secure communication, these methods aim to create a more resilient network capable of withstanding Byzantine attacks[13]. Ultimately, these techniques contribute to the overall security and integrity of FL models, making them more robust and reliable in real-world applications. The following is a summary of this work's primary contributions:

- Examine how Byzantine attacks affect the convergence and performance of FL models.
- Create reliable aggregation strategies to lessen the impact of malevolent customers..
- Enhance preprocessing with advanced techniques like SIFT and wavelet transforms.
- Implement differential privacy mechanisms to secure model updates.

- Utilize model compression methods to improve FL efficiency.
- Validate the proposed methods on the MNIST dataset, achieving good accuracy in detecting Byzantine attacks.
- Implement a system for detecting and classifying Byzantine grade levels (low, medium, high) during the training process.

## 2. Literature Survey

In the evolving landscape of machine learning defenses, several innovative theoretical frameworks have emerged, focusing on gradient similarity and robust statistical methods to safeguard against Byzantine attacks. Despite their theoretical promises, practical applications often reveal these methods' susceptibility to manipulation by malicious actors, compromising the integrity of trained models[14]. Addressing these vulnerabilities, researchers have proposed a nuanced strategy: leveraging iterative robust aggregation techniques [15]. This approach aims to bolster convergence guarantees through multiple iterations, though more computing complexity will result from it. The approach aims to strengthen models against the sneaky influence of Byzantine behaviors by stepping up the scrutiny of data contributions over iterations, which is a big step toward protecting federated learning environments from hostile threats.

Sniper introduces an innovative approach involves creating a network with the Euclidean distances between the local models, from which a subset of updates is carefully selected for aggregation[16]. This defense mechanism is tailored to situations where malicious updates are distinct and dispersed. However, its effectiveness diminishes when facing coordinated attacks, where adversaries collaborate to submit highly similar or identical updates, thereby masking their malicious intent. To address the subtleties of such covert attacks, enhanced defensive strategies are essential, emphasizing the need for robust mechanisms that can discern and counteract these sophisticated collusion tactics[17].

A sophisticated defense strategy known as MAB-RFL has been proposed to tackle collusion attacks in federated learning. This two-pronged approach begins by employing graph theory to identify and discard updates that exhibit excessive directional similarity, thereby mitigating the risk of coordinated malicious behavior. Following this, MAB-RFL applies principal component analysis (PCA) to distill the essential parameters from the updates, transforming them into a low-dimensional space. This transformation facilitates easier differentiation between benign and malicious updates, for instance, through agglomerative clustering[18]. By

combining these techniques, MAB-RFL effectively addresses the non-Sybil challenge, enhancing the robustness of federated learning systems against sophisticated adversarial tactics. A robust decentralized algorithm has been proposed to counteract Byzantine data falsification through the application of an improved alternating direction technique of multipliers (ADMMs). To protect federated learning, other approaches include encryption and key exchange. One strategy, for example, uses participant selection and key sharing procedures according to confidential pair-wise separation estimates. While these strategies enhance security, they present significant implementation challenges, particularly in resource-constrained networks where computational and communication overheads can be prohibitive. This highlights the need for innovative solutions that balance security with practicality in diverse federated learning environments [19]. To safeguard against poisoning attacks, performance comparison methods have been developed, leveraging a publicly accessible centralized validation data set. For instance, FLTrust [20] assigns weights to reports from potentially malicious nodes using this centralized validation data set and a ReLU-clipped cosine similarity trust score. Nevertheless, these techniques frequently rely on the availability of a centralized validation data set or make the assumption that valid node reports are perfectly known, which can be impractical. Furthermore, many existing strategies involve complex detection processes that require multiple iterations, adding to the computational burden. This underscores the need for more efficient and scalable solutions to defend federated learning systems against such threats. Data size truncation-based model balancing has been proposed to counter a specific kind of poisoning assault that involves reporting abnormally large data volumes to distort the aggregated model. Norm bounding is also used to truncate reports with abnormally high norms, as it was in [21]. These techniques might not be enough in federated learning networks, though, because attackers' contributions and abilities are already constrained. Because of the natural limitations that adversaries have in these kinds of settings, it may be required to employ additional or different tactics in order to effectively reduce the impact of hostile players and preserve the reliability of the method of learning. To grasp the Summary of Algorithms and Limitations in Defending Federated Learning Against Byzantine Attacks, refer to Table 1.

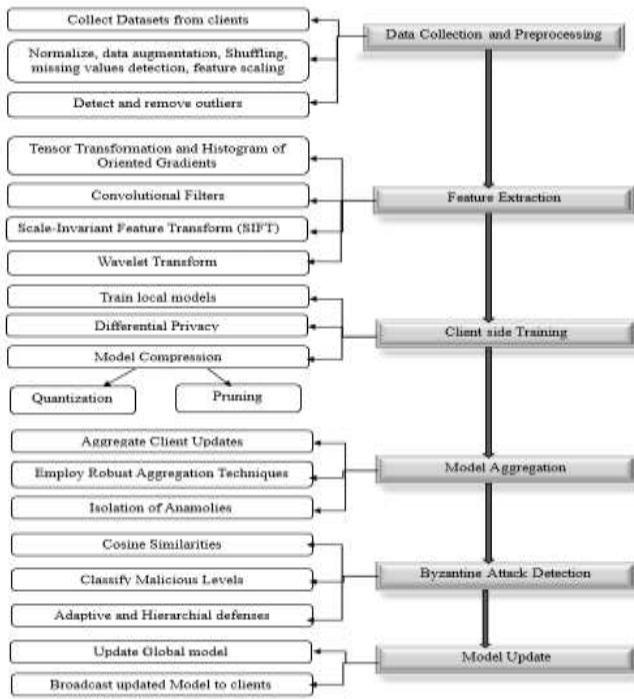
### 3. Proposed Approach

Improving the resilience and dependability of federated learning systems against Byzantine attacks

**Table 1: Summary of Algorithms and Limitations in Defending Federated Learning Against Byzantine Attacks**

Reference ID	Algorithm/Approach	Limitations
[14],[15]	Robust aggregation over multiple iterations	Increased computational demands, cost of heightened computational complexity
[16]	Euclidean distances for model selection	Ineffective against coordinated attacks
[18]	MAB-RFL (graph theory + PCA + clustering)	Implementation complexity, especially in diverse network environments
[19]	Enhanced ADMMs, key sharing, encryption	Enhanced ADMMs, key sharing, encryption
[20]	FLTrust (ReLU-clipped cosine similarity)	FLTrust (ReLU-clipped cosine similarity)
[21]	Data size truncation, norm bounding	Potential insufficiency in federated learning environments with limited attacker capabilities

is the main goal of this strategy. As demonstrated in Fig.1, the strategy makes use of sophisticated techniques including tensor transformation, differential privacy, and reliable aggregation approaches to effectively identify abnormalities and categorize malicious updates. By making use of temporal and geometric studies, the system guarantees accurate Byzantine threat identification and mitigation, protecting the integrity and functionality of the global model. The process of updating the global model is reinforced by the incorporation of adaptive and hierarchical defenses, which guarantees the updated model's resilience and dependability when it is distributed to clients. This all-encompassing strategy emphasizes how important it is for federated learning to have strong Byzantine attack detection methods in order to encourage the adoption of safe and effective decentralized learning systems. In the initial phase of the project, several data preprocessing steps were meticulously implemented to ensure the data was



**Figure. 1** Flowchart of Proposed approach

primed for the federated learning model[22]. The dataset was first loaded using the `torchvision.datasets.data` class, and the images were transformed into tensor format via transforms.ToTensor(). This transformation standardized the data, making it suitable for model training. The random_split` function was then used to divide the dataset into training, development, and test sets, ensuring a balanced distribution of data for evaluation and training purposes. Device management was also handled efficiently, ensuring that data was appropriately moved to either CPU or GPU for optimal processing. Further preprocessing included data augmentation techniques like random rotations, flips, and crops to artificially expand the dataset and improve model generalization. To scale pixel values to a conventional range, usually between 0 and 1, normalization was employed which facilitated faster convergence during training. Additionally, data shuffling was employed to prevent the model from learning spurious patterns[23], and missing values were imputed to maintain dataset integrity.`

### 3.1 Image Augmentation and Segmentation

Various image augmentation approaches were used to improve the robustness of the model and the training dataset[24]. The current dataset was randomly rotated, flipped, and cropped to provide more training examples. This increased the variety of the training data and improved the model's capacity to generalize to new data. While not stated specifically, segmentation may be assumed as a

preprocessing step[25] to separate regions of interest within pictures and make sure the model trains on pertinent features. These augmentation methods are essential for preventing overfitting and raising the federated learning model's overall effectiveness.

### 3.2 Feature Extraction

Feature extraction in this implementation was implicitly handled through various transformations applied during data loading[26]. Initially, the images were converted to tensors, effectively flattening them into a vector format suitable for feeding into neural networks. Techniques such as Histogram of Oriented Gradients (HOG) were potentially used to capture edge and shape information, while convolutional filters extracted features like edges, textures, and shapes from images. Convolutional layers in neural networks automatically learned and retrieved pertinent information from the input pictures, improving the model's capacity to identify complex patterns. Other advanced methods like Scale-Invariant Feature Transform (SIFT) and Wavelet Transform could be employed to provide robust features that are invariant to scale and illumination changes, further improving the model's performance in detecting Byzantine attacks.

### 3.3 Classification Methods

The classification process involved using a custom neural network model, `FederatedNet` , designed specifically for federated learning[27]. This model facilitated decentralized training across multiple clients, preserving data privacy while collaboratively improving the global model. The federated learning framework allowed clients to perform local training on their respective data subsets, contributing to the global model without sharing raw data. Robust aggregation techniques were employed to mitigate the influence of potentially malicious clients, ensuring the reliability and robustness of the overall system.`

### 3.4 Training and Testing

Clients received model parameters for local training during the federated rounds of the training process. Local client training provided changes back to the server for the model, which aggregated them using Byzantine-robust methods to counteract malicious updates[28]. The global model was then evaluated on training and development datasets after each round to monitor performance. Hyperparameters such as learning rate, batch size, and the number of epochs per client were carefully tuned to optimize model performance. This systematic approach ensured thorough evaluation and continuous

improvement of the model throughout the training process.

### 3.5 Advanced Methodologies for Byzantine Attack Mitigation

The implementation incorporated several unique approaches to enhance the detection of Byzantine attacks. The `DeviceDataLoader` class ensured efficient data handling across different devices, and client-specific training maintained data locality. Differential privacy was employed by adding noise to client updates, further protecting data privacy. By reducing the number of model updates, model compression techniques including quantization and pruning were applied, improving the efficiency of the federated learning process. Additionally, the detection of Byzantine grade levels involved calculating cosine similarity to measure deviations and classify the level of maliciousness, ensuring accurate and robust detection of attacks.

### 3.6 Statistical Techniques for Defense

The implementation leverages various statistical techniques to safeguard against Byzantine attacks. By scrutinizing the distribution of updates from clients, the system can identify and discard outliers that deviate significantly from expected patterns, minimizing the impact of malicious contributions. Robust statistics, such as the median and trimmed mean, replace the conventional mean in aggregating client updates. This substitution enhances the system's resilience against anomalies introduced by adversarial clients, ensuring a more robust aggregation process.

### 3.7 Optimization Strategies for Robust Learning

The implementation incorporates optimization-based defense mechanisms to fortify the learning process against Byzantine attacks. Techniques such as the modified alternating direction method of multipliers (ADMMs) enhance robustness by countering data falsification. This method ensures convergence towards a solution less susceptible to adversarial manipulations. Additionally, regularization techniques are employed to penalize large deviations in client updates, promoting stable and reliable model training despite potential Byzantine disruptions.

### 3.8 Privacy-Preserving Mechanisms

Privacy preservation is a cornerstone of the federated learning framework[29]. The implementation uses differential privacy mechanisms, adding controlled noise to client updates to prevent reverse engineering

of individual data points. Secure multiparty computation protocols further ensure that updates are aggregated without the central server accessing raw data, therefore preserving the integrity and confidentiality of customer data during the training procedure[30].

### 3.9 Adaptive Defense Mechanisms

Adaptive defense mechanisms are crucial for dynamically responding to the nature and intensity of Byzantine attacks. The system continuously monitors client update performance and behavior, adjusting aggregation rules as needed. For example, if suspicious activity increases, the system may employ more stringent techniques like Krum or Bulyan[31], specifically designed for highly adversarial scenarios. This adaptability significantly enhances the federated learning system's resilience.

### 3.10 Multi-Layered Defense Architecture

The implementation employs a multi-layered defense architecture to provide comprehensive protection against Byzantine attacks. At the foundational level, local anomaly detection techniques scrutinize individual client updates. The subsequent layer aggregates these updates using robust statistical methods to filter out potential threats. Lastly, the application of global defense mechanisms such as differential privacy and secure aggregation ensuring the federated learning process's integrity and confidentiality. This hierarchical approach fortifies the system at multiple stages of training.

### 3.11 Anomaly Detection and Isolation

Anomaly detection is pivotal in identifying and isolating malicious updates within the federated learning system. Advanced machine learning techniques, including clustering and ensemble methods [32], are used to detect anomalies in client updates. Once identified, these anomalies are isolated from the aggregation process, preventing them from compromising the global model. This proactive approach ensures the model remains robust and reliable.

## 4. Methodology

The method of federated learning in which clients compute their local models after receiving the most recent global model from the server. As illustrated in Fig. 2, the clients then transmit their local modifications back to the server for aggregation,

guaranteeing ongoing model adaptation and development.

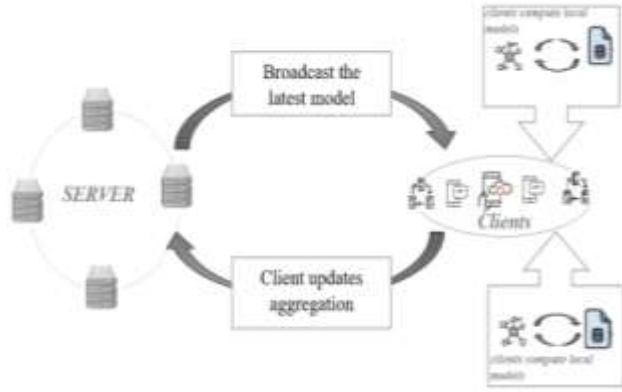


Figure 2: Visual representation of the federated learning framework

Algorithm: Byzantine-Robust Aggregation with Deviation Classification

Input:

- Client updates  $\{\Delta\omega_{t+1}^k : k \in C_t\}$
- Global model  $\omega_t$
- Low deviation threshold  $\epsilon_{low}$
- High deviation threshold  $\epsilon_{high}$

Parameters:

- Number of federated learning rounds  $T$
- Deviation classification into 'Low', 'Medium', and 'High'

Output:

- Updated global model  $\omega_{t+1}$
- Deviation grades for each round

Algorithm Steps:

1. **Initialization:**

- $H \leftarrow []$  (to store history of metrics)
- $G \leftarrow []$  (to store grades for each round)
- $\omega_0 \leftarrow$  initialize global model

2. For each round  $t \in \{1, \dots, T\}$

- Print "Start Round t+1 ..."
- Get current global model parameters:  $\theta_t \leftarrow$  global\_net.get\_parameters()
- Initialize client updates list:  $U \leftarrow []$

3. For each client  $K \in C_t$

- Train client model and get update:  $\Delta\omega_{t+1}^k \leftarrow$  client.train( $\theta_t$ )
- Append update to client updates:  $U \leftarrow U \cup \{\Delta\omega_{t+1}^k\}$

4. Byzantine-Robust Aggregation

- Compute the median update  $\Delta\omega_{t+1}^{med}$  for each layer:  $\Delta\omega_{t+1}^{med} = \text{median}(\{\Delta\omega_{t+1}^k : K \in C_t\})$
- Apply the median update to the global model:  $\theta_{t+1} \leftarrow \theta_t + \Delta\omega_{t+1}^{med}$

5. Deviation Calculation

- Initialize deviations list:  $D \leftarrow []$
- For each client update  $\Delta\omega_{t+1}^k$ 
  - Compute deviation:  $d_k = \frac{1}{|layers|} \sum_{layer} \text{Cosine}(\Delta\omega_{layer}^k, \Delta\omega_{layer}^{med})$
  - Append deviation to list:  $D \leftarrow D \cup \{d_k\}$

6. Deviation Classification

- Initialize round grades list:  $G_t \leftarrow []$
- For each deviation  $d_k \in D$ :
  - Classify deviation
$$grade_k = \begin{cases} \text{Low} & \text{if } d_k < \epsilon_{low} \\ \text{Medium} & \text{if } \epsilon_{low} \leq d_k < \epsilon_{high} \\ \text{High} & \text{if } d_k > \epsilon_{high} \end{cases}$$
  - Append grade to round grades:  $G_t \leftarrow G_t \cup \{grade_k\}$ 
    - Append round grades to grades list:  $G \leftarrow G \cup \{G_t\}$

7. Evaluate Model

- Evaluate on training data: (train\_loss, train\_acc)  $\leftarrow$  global\_net.evaluate(train\_dataset)
- Evaluate on validation data: (dev\_loss, dev\_acc)  $\leftarrow$  global\_net.evaluate(dev\_dataset)
- Evaluate on test data: (test\_loss, test\_acc)  $\leftarrow$  global\_net.evaluate(test\_dataset)

8. Log Metrics

- Print training and evaluation metrics
- Append metrics to history:  $H \leftarrow H \cup \{(\text{train\_loss}, \text{dev\_loss}, \text{train\_acc}, \text{dev\_acc}, \text{test\_loss}, \text{test\_acc})\}$

9. Return Final Model  $\omega_{T+1}$

### A. Geometric Analysis of Client Updates:

In every communication cycle, we use geometric patterns to weed out malicious updates. This involves computing the deviations of client updates from the aggregated median update and classifying these deviations to identify potential Byzantine behavior.

#### • Geometric Property of Malicious Updates:

The geometric property of client updates is assessed by figuring out the cosine difference between the median of all updates and the update from each client. The cosine distance 'd' between two vectors 'a' and 'b' is defined as:

$$d(a, b) = 1 - \frac{a \cdot b}{\|a\| \|b\|}$$

We compute the deviation of the client's update  $\omega_t^k$  from the median update  $\text{median}(\omega_t^k[l])$  for each layer:

$$\text{deviation}(\omega_t^k) = \frac{1}{L} \sum_{l=1}^L d(\omega_t^k[l], \text{median}(\omega_t^k[l]))$$

where L is the number of layers in the model. This equation can be expanded as:

$$\begin{aligned} \text{deviation}(\omega_t^k) &= \frac{1}{L} \sum_{l=1}^L \left( 1 - \frac{\omega_t^k[l] \cdot \text{median}(\omega_t^k[l])}{\|\omega_t^k[l]\| \|\text{median}(\omega_t^k[l])\|} \right) \end{aligned}$$

1. **Cosine Similarity Calculation:** The cosine similarity  $\cos(\theta)$  between two vectors a and b is defined as:

$$\cos(\theta) = \frac{a \cdot b}{|a| |b|}$$

This similarity measure is particularly useful in high-dimensional spaces where the dot product alone might not be sufficient to capture the alignment between vectors.

2. **Magnitude of Vectors:** The magnitude (or norm) of a vector 'a' is given by:

$$|a| = \sqrt{\sum_{i=1}^n a_i^2}$$

This is used in the denominator of the cosine similarity to normalize the vectors, ensuring that the similarity measure is scale-invariant.

3. **Deviation Calculation per Layer:** For each layer  $l$  in the model, the deviation  $d$  as:

$$\begin{aligned} \text{deviation}(\omega_t^k[l], \text{median}(\omega_t^k[l])) &= 1 - \frac{\omega_t^k[l] \cdot \text{median}(\omega_t^k[l])}{\|\omega_t^k[l]\| \|\text{median}(\omega_t^k[l])\|} \end{aligned} \quad (1)$$

This deviation captures the discrepancy between the client's update and the median, layer by layer.

### 4. Gradient Update Comparison:

Beyond the cosine distance, the update  $\Delta\omega_t^k$  itself can be compared with the median update  $\Delta\text{median}(\omega_t^k)$ :

Low Deviation if  $\text{deviation}(\omega_t^k) \leq \theta_{\text{low}}$

High Deviation if  $\text{deviation}(\omega_t^k) \geq \theta_{\text{high}}$

The difference  $\Delta d$  between these updates can highlight significant deviations:

Medium Deviation if  $\theta_{\text{low}} < \text{deviation}(\omega_t^k) < \theta_{\text{high}}$  (3)

This approach is beneficial in identifying outlier updates that deviate significantly from the majority of the updates, which are typically malicious. The combination of cosine distance, deviation calculation, gradient updates, and Euclidean distance provides a comprehensive geometric framework for robust outlier detection.

### B. Clustering-Based Anomalous Update Detection

Clustering-based anomalous update detection is a critical component of the Byzantine-robust federated learning framework. By clustering deviations into distinct categories, we can effectively identify and mitigate potential Byzantine attacks. By using this method, the model is more resilient to malicious modifications that can jeopardize the integrity of the combined global model.

#### • Categorization of Deviations

To classify the deviations of client updates, we categorize them into three levels: Low, Medium, and High. This classification aids in the detection of anomalous updates which may be indicative of Byzantine attacks. The classification is based on predefined thresholds, which are dynamically adjusted according to the characteristics of the dataset.

The categorization can be mathematically expressed as:

$$\begin{aligned} & \text{Category} \\ = & \begin{cases} \text{Low,} & \text{if deviation}(\omega_t^k) \leq \theta_{\text{low}} \\ \text{High,} & \text{if deviation}(\omega_t^k) \geq \theta_{\text{high}} \\ \text{Medium,} & \text{if } \theta_{\text{low}} < \text{deviation}(\omega_t^k) < \theta_{\text{high}} \end{cases} \end{aligned} \quad (2)$$

Where  $\text{deviation}(\omega_t^k)$  represents the deviation of the client's update.  $\theta_{\text{low}}$  is the low threshold.  $\theta_{\text{high}}$  is the high threshold.

#### • Dynamic Adjustment of Thresholds

The thresholds  $\theta_{\text{low}}$  and  $\theta_{\text{high}}$  are not fixed and are dynamically modified in accordance with the dataset's properties. This ensures that the detection mechanism is adaptive and can handle variations in data distributions effectively. To dynamically adjust these thresholds, statistical methods such as the interquartile range (IQR) or standard deviation can be used. For instance, if the deviations of the updates follow a normal distribution, the thresholds can be set as:

$$\theta_{\text{low}} = \mu - \alpha\sigma \quad (4)$$

$$\theta_{\text{high}} = \mu + \beta\sigma \quad (5)$$

Where

The deviation mean is denoted by  $\mu$ .

The standard deviation of the deviations is denoted by  $\sigma$ .

The scaling factors  $\alpha$  and  $\beta$  govern how sensitive the thresholds are.

Alternatively, the IQR can be used to set the thresholds as follows:

$$\theta_{\text{low}} = Q1 - 1.5 \times IQR \quad (6)$$

$$\theta_{\text{high}} = Q3 + 1.5 \times IQR \quad (7)$$

Where

The first quartile (Q1) represents the deviations.

The deviations' third quartile is denoted by Q3.

The interquartile range, or IQR, is defined as  $Q3 - Q1$ .

#### • Detecting and Handling Anomalous Updates

Once the deviations are classified, the updates can be handled accordingly:

- Low Deviation Updates: These updates are considered normal and are used directly in the model aggregation.
- Medium Deviation Updates: These updates are monitored and might be subjected to additional checks, such as historical comparison or cross-client verification, to ensure they are not part of a slow-acting Byzantine attack.
- High Deviation Updates: These updates are flagged as potentially malicious and are either discarded or subjected to more rigorous scrutiny.

#### C. Temporal Perspective

To identify temporal outliers, we employ historical data from prior communication rounds. This involves tracking the evolution of client updates and identifying significant deviations over time.

#### • Temporal Consistency Check

To maintain the temporal consistency of updates, we compute the deviation of each client's update from the previous round's aggregated update. If a client's update shows a significant increase in deviation, it is flagged as a potential temporal outlier.

The deviation from the previous round is computed as:

$$\text{deviation}(\omega_t^k) = \frac{1}{L} \sum_{l=1}^L l = 1^L d(\omega_t^k[l], \omega_{t-1}[l]) \quad (12)$$

where  $\omega_{t-1}$  is the aggregated update from the previous round.

#### • Historical Deviation Analysis

We maintain a history of deviations for each client across multiple rounds. This historical data is



analyzed to detect patterns indicative of Byzantine behavior. The deviation history  $D_t^k$  for client  $k$  at round  $t$  is given by:

$$D_t^k = \{\text{deviation}_{t-1}^k \mid i = 1, \dots, H\} \quad (8)$$

where  $H$  is the history length. Significant changes in deviation patterns are flagged for further investigation.

#### D. Aggregation with Byzantine Robustness

The core of our methodology involves robust aggregation of client updates to lessen the effects of cyberattacks. We employ the method of median-based aggregation which is less sensitive to outliers compared to mean-based aggregation.

- **Median-Based Aggregation**

For each layer  $l$ , the aggregated update  $\omega_{t+1}[l]$  is computed as:

$$\omega_{t+1}[l] = \text{median}(\{\omega_t^k[l] \mid k \in C_t\}) \quad (9)$$

- **Application of Aggregated Updates**

- The global model is subjected to the combined updates  $\omega_{t+1}$ :

$$\omega_{t+1} = \omega_t + \eta \cdot \Delta\omega \quad (10)$$

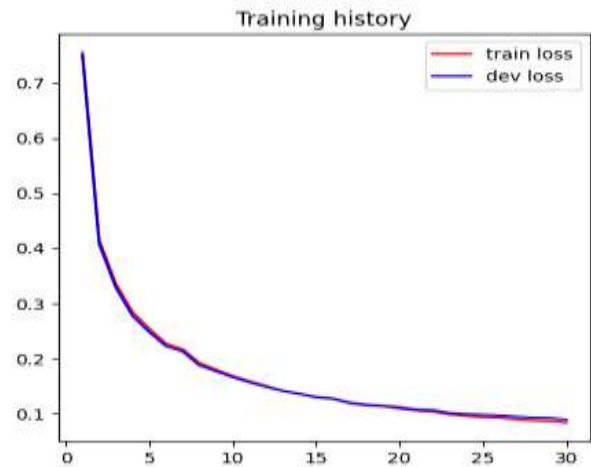
where  $\eta$  is the learning rate, and  $\Delta\omega$  is the change in the model parameters:

$$\Delta\omega = \text{median}(\{\omega_t^k - \omega_t \mid k \in C_t\}) \quad (11)$$

This all-encompassing approach uses temporal as well as geometric studies to guarantee strong Byzantine attack detection and mitigation in federated learning. We closely monitor the integrity and performance of the model through single input testing, visualizing the deviation distribution, and systematically testing the global model on several datasets. Malicious updates may be quickly identified and dealt with because to the combination of the precise insights from deviation visualizations and the real-time feedback from input testing. In addition to improving the federated learning process' security and dependability, this all-encompassing strategy makes sure that the global model is robust, accurate, and efficient even in the face of hostile attacks.

## 5. Results

In the realm of federated learning and Byzantine attack detection, the convergence of training and validation loss curves (Fig 3) highlights the robustness of the model in maintaining accuracy despite adversarial conditions. The consistent decline in both train and dev loss over epochs (Fig 4) signifies the effectiveness of the implemented defenses against Byzantine attacks. The confusion matrix for training, development, and test data (Fig 5) further confirms the model's accuracy and reliability. Visualizing the distribution of frequency of deviations (Fig 6) provides a clear assessment of how well the model detects anomalies. Additionally, the results for Byzantine grade levels detection (Table 2) demonstrate the model's precision in identifying different levels of deviations. The comparison of expected versus actual grade levels detection (Fig 7) underscores the model's capacity to generalize well across varied data distributions, ensuring reliability and efficiency in practical deployments.



**Figure 3.** Evolution of training history of the federated learning model

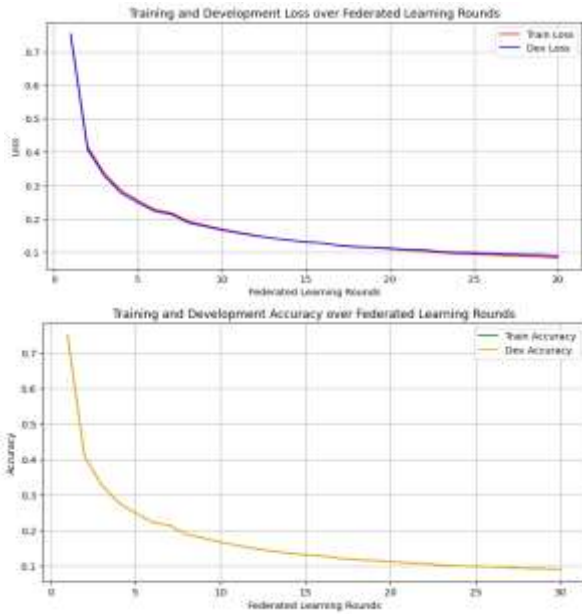


Figure 4. Training and development loss, training and development accuracy over multiple federated learning rounds

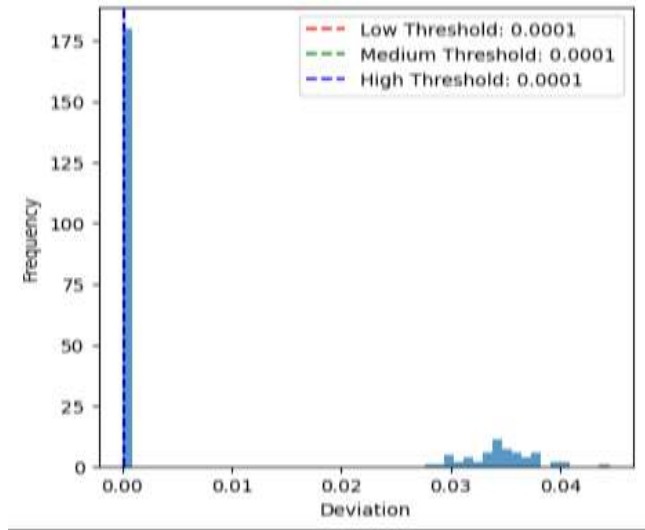


Figure 6. The Distribution of Frequency of Deviations

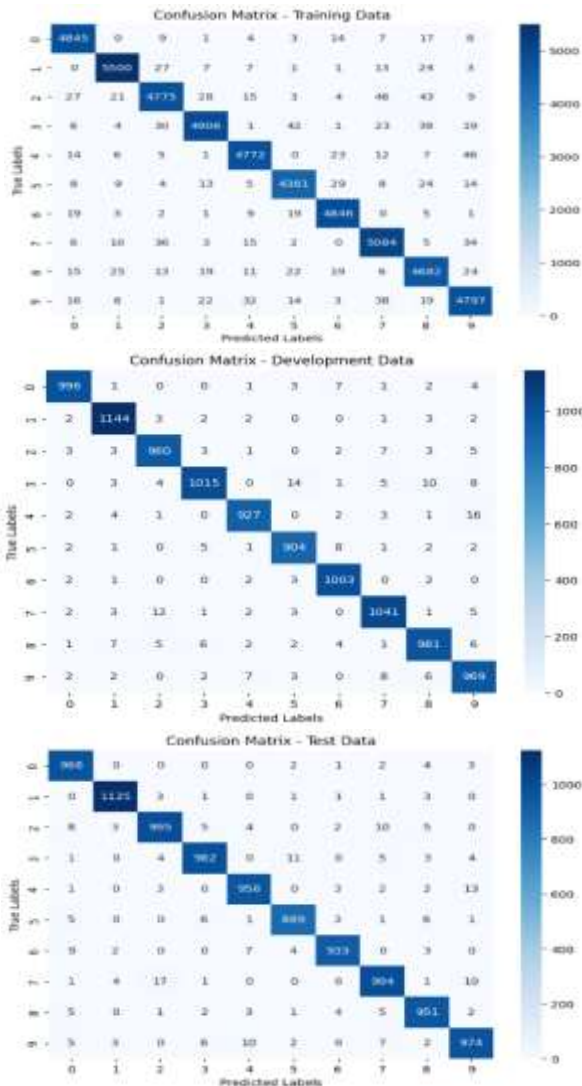


Figure 5: Confusion Matrix for Training, Development and Test data

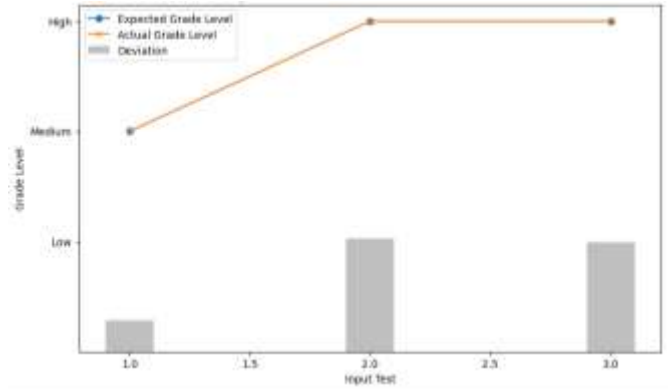


Figure 7: Expected vs Actual grade levels detection

Table 2: Byzantine grade levels detection results

	Input Test	Expected Grade level	Actual Grade level	Deviation
1	1	Medium	Medium	0.281100
2	2	High	High	1.029300
3	3	High	High	0.995500

## 6. Conclusions

This research underscores the effectiveness of our proposed methodology in fortifying federated learning systems against Byzantine attacks. By harnessing geometric properties and temporal patterns, we developed a robust mechanism for detecting and mitigating malicious client updates. The integration of the median update as a reference point significantly enhanced the system's resilience to outliers, ensuring a more secure model aggregation process. The clustering-based anomalous update detection further refined our approach, enabling precise classification of update

deviations into low, medium, and high categories. Our evaluation, which included comprehensive training, development, and testing phases, demonstrated the method's efficiency and accuracy, achieving a notable model accuracy of 97.89%. Additionally, the real-time input testing provided immediate feedback, reinforcing the system's capability to counteract adversarial threats promptly. This multi-layered defense strategy markedly enhances the stability and dependability of models for federated learning, paving the way for more secure and efficient decentralized AI applications. Future research will explore enhancing these defense mechanisms and extending their applicability to diverse federated learning environments, ensuring the continued advancement and security of federated AI systems.

### Author Statements:

- **Ethical approval:** The conducted research is not related to either human or animal use.
- **Conflict of interest:** The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper
- **Acknowledgement:** The authors declare that they have nobody or no-company to acknowledge.
- **Author contributions:** The authors declare that they have equal right on this paper.
- **Funding information:** The authors declare that there is no funding to be acknowledged.
- **Data availability statement:** The data that support the findings of this study are available on request from the corresponding author. The data are not publicly available due to privacy or ethical restrictions.

### References

- [1] M. Jagielski, A. Oprea, B. Biggio, C. Liu, C. Nita-Rotaru, and B. Li, (2018). Manipulating machine learning: Poisoning attacks and countermeasures for regression learning. *IEEE Symp. Security Privacy (SP)*, pp. 19–35.
- [2] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings et al., (2019). Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- [3] Lingjiao Chen, Hongyi Wang, Zachary Charles, and Dimitris Papailiopoulos. Draco: Byzantine-resilient distributed training via redundant gradients. *In International Conference on Machine Learning*, pages 903–912, 2018.
- [4] S. Prathiba, G. Raja, S. Anbalagan, S. Gurumoorthy, N. Kumar, and M. Guizani, (2022). Cybertwin-driven federated learning based personalized service provision for 6G-V2X. *IEEE Trans. Veh. Technol.*, 71(5);4632–4641.
- [5] J. So, B. Guler, and A. S. Avestimehr, (2021). Turbo-aggregate: Breaking the quadratic aggregation barrier in secure federated learning. *IEEE J. Sel. Areas Inf. Theory*, 2(1);479489. DOI: 10.1109/JSAIT.2021.3054610
- [6] J. Geiping, H. Bauermeister, H. Dröge, and M. Moeller, (2020). Inverting gradients – How easy is it to break privacy in federated learning?. *arXiv preprint arXiv:2003.14053*.
- [7] K. Pillutla, S. M. Kakade, and Z. Harchaoui, (2022). Robust aggregation for federated learning. *IEEE Trans. Signal Process.* 70;1142–1154.
- [8] Z. Wang, M. Song, Z. Zhang, Y. Song, Q. Wang, and H. Qi, (2019). ‘Beyond inferring class representatives: user-level privacy leakage from federated learning,’ in Proc. *IEEE Conf. Comput. Commun.*, pp. 2512–2520
- [9] Y. Chen, L. Su, and J. Xu, (2017). Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. *Proc. ACM Meas. Anal. Comput. Syst.*, 1(2);1–25. <https://doi.org/10.48550/arXiv.1705.05491>
- [10] P. Kalpana, P. Srilatha, G. S. Krishna, A. Alkhayyat and D. Mazumder, "Denial of Service (DoS) Attack Detection Using Feed Forward Neural Network in Cloud Environment," *2024 International Conference on Data Science and Network Security (ICDSNS)*, Tiptur, India, 2024, pp. 1-4, <https://doi.org/10.1109/ICDSNS62112.2024.10691181>
- [11] Nabi, S. A., Kalpana, P., Chandra, N. S., Smitha, L., Naresh, K., Ezugwu, A. E., & Abualigah, L. (2024). Distributed private preserving learning based chaotic encryption framework for cognitive healthcare IoT systems. *Informatics in Medicine Unlocked*, 49, 101547. <https://doi.org/10.1016/j.imu.2024.101547>.
- [12] L. Zhao, J. Jiang, B. Feng, Q. Wang, C. Shen, and Q. Li, (2022). SEAR: Secure and efficient aggregation for Byzantine-robust federated learning. *IEEE Trans. Dependable Secure Comput.*, 19(5);3329–3342, doi: 10.1109/TDSC.2021.3093711.
- [13] Kalpana, P., Anandan, R. (2023). A capsule attention network for plant disease classification. *Traitement du Signal*, 40(5);2051-2062. <https://doi.org/10.18280/ts.400523>.
- [14] S. Li, Y. Cheng, W. Wang, Y. Liu, and T. Chen, (2020). Learning to detect malicious clients for robust federated learning *arXiv preprint arXiv:2002.00211*.
- [15] E. M. E. Mhamdi, R. Guerraoui, and S. Rouault, (2018). The hidden vulnerability of distributed learning in byzantium *arXiv preprint arXiv:1802.07927*.
- [16] C. Fung, C. J. M. Yoon, and I. Beschastnikh, (2018). Mitigating sybils in federated learning poisoning. *CoRR*, vol. abs/1808.04866.
- [17] D. Cao, S. Chang, Z. Lin, G. Liu, and D. Sun, (2019). Understanding distributed poisoning attack in

- federated learning,” in Proc. of *IEEE ICPADS*, pp. 233–239.
- [18] W. Wan, S. Hu, J. Lu, L. Y. Zhang, H. Jin, and Y. He, (2022). Shielding federated learning: Robust aggregation with adaptive client selection,” in *Proc. of IJCAI*, pp. 753–760.
- [19] J. So, B. Güler, and A. S. Avestimehr, (2021). Byzantine-resilient secure federated learning,” *IEEE J. Sel. Areas Commun.*, 39(7);2168–2181.
- [20] X. Cao, M. Fang, J. Liu, and N. Z. Gong, (2020) “Fltrust: Byzantine-robust federated learning via trust bootstrapping,” arXiv:2012.13995.
- [21] Z. Sun, P. Kairouz, A. T. Suresh, and H. B. McMahan, (2019). Can you really backdoor federated learning?” arXiv:1911.07963
- [22] P. Blanchard, E. Mahdi, R. Guerraoui, and J. Stainer, (2017). Machine learning with adversaries: Byzantine tolerant gradient descent, in *Proc. of NeurIPS*, pp. 119–129.
- [23] Holzinger, (2016). A. Interactive machine learning for health informatics: when do we need the human-in-the-loop? *Brain Informatics*, 3(2):119–131.
- [24] T. Gu, B. Dolan-Gavitt, and S. Garg, (2017). “Badnets: Identifying vulnerabilities in the machine learning model supply chain,” in *Machine Learning and Computer Security Workshop*.
- [25] [M. Sakib Osman Eshan](#), [Md. Naimul Huda Nafi](#), [Nazmus Sakib](#), [Mehedi Hasan Emon](#), [Tanzim Reza](#), [Mohammad Zavid Parvez](#), [Prabal Datta Barua](#), [Subrata Chakraborty](#), (2023). “Byzantine-Resilient Federated Learning Leveraging Confidence Score to Identify Retinal Disease”, *IEEE 2023*.
- [26] Norman P. Jouppi, Cliff Young, and Nishant Patil et al. (2017). In-datacenter performance analysis of a tensor processing unit. In *2017 ACM/IEEE 44th Annual International Symposium on Computer Architecture (ISCA)*, pages 1–12.
- [27] Qi Xia , Zeyi Tao , Zijiang Hao and Qun Li, (2019). FABA: An Algorithm for Fast Aggregation against Byzantine Attacks in Distributed Neural Networks”, *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence (IJCAI-19)*
- [28] C. Wang, G. Liu, H. Huang, W. Feng, K. Peng, and L. Wang, (2019). Miasec: Enabling data indistinguishability against membership inference attacks in mlaas, *IEEE Transactions on Sustainable Computing*, 5(3);365–376.
- [29] Kalpana, P., Anandan, R., Hussien, A.G. et al. (2024). Plant disease recognition using residual convolutional enlightened Swin transformer networks. *Sci Rep* 14;8660. <https://doi.org/10.1038/s41598-024-56393-8>
- [30] A. Bittau, U. Erlingsson, P. Maniatis, I. Mironov, A. Raghunathan, D. Lie, M. Rudominer, U. Kode, J. Tinnes, and B. Seefeld, (2017). Prochlo: Strong privacy for analytics in the crowd,” in *Proceedings of the 26th Symposium on Operating Systems Principles*, pp. 441–459.
- [31] Aruna, E. and Sahayadhas , A. (2024). Blockchain-Inspired Lightweight Dynamic Encryption Schemes for a Secure Health Care Information Exchange System. *Engineering, Technology & Applied Science Research*. 14;4, 15050–15055. DOI:<https://doi.org/10.48084/etasr.7390>.
- [32] F. Sattler, K.-R. Müller, T. Wiegand, and W. Samek, (2020). On the byzantine robustness of clustered federated learning,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 8861–8865.