

Sentinel: An Aggregation Function to Secure Decentralized Federated Learning

Chao Feng^{a,*}, Alberto Huertas Celdrán^a, Janosch Baltensperger^a, Enrique Tomás Martínez Beltrán^b, Pedro Miguel Sánchez Sánchez^b, G r me Bovet^c and Burkhard Stiller^a

^aCommunication Systems Group CSG, Department of Informatics, University of Zurich UZH, CH–8050 Z rich, Switzerland

^bDepartment of Information and Communications Engineering, University of Murcia, 30100–Murcia

^cCyber-Defence Campus, armasuisse Science & Technology, CH–3602 Thun, Switzerland

Abstract. Decentralized Federated Learning (DFL) emerges as an innovative paradigm to train collaborative models, addressing the single point of failure limitation. However, the security and trustworthiness of FL and DFL are compromised by poisoning attacks, negatively impacting its performance. Existing defense mechanisms have been designed for centralized FL and they do not adequately exploit the particularities of DFL. Thus, this work introduces Sentinel, a defense strategy to counteract poisoning attacks in DFL. Sentinel leverages the accessibility of local data and defines a three-step aggregation protocol consisting of similarity filtering, bootstrap validation, and normalization to safeguard against malicious model updates. Sentinel has been evaluated with diverse datasets and data distributions. Besides, various poisoning attack types and threat levels have been verified. The results improve the state-of-the-art performance against both untargeted and targeted poisoning attacks when data follows an IID (Independent and Identically Distributed) configuration. Besides, under non-IID configuration, it is analyzed how performance degrades both for Sentinel and other state-of-the-art robust aggregation methods.

1 Introduction

Federated Learning (FL) emerged in 2016 as a promising paradigm able to train Machine and Deep Learning (ML & DL) models in a privacy-preserving and collaborative manner [1]. Nowadays, FL is used in many different scenarios such as healthcare, cybersecurity, networking, and so on. More in detail, FL combines the capabilities of communication systems with ML to solve data privacy concerns and increase scalability while distributing the learning process throughout a federation of participants. However, the original concept of FL relies on a central entity aggregating weights to create global models, which implicates drawbacks such as communication bottlenecks and single points of failure. In this context, Decentralized FL (DFL) has been recently proposed as an alternative to improve these limitations [11]. In DFL, communication systems get much more relevance than in CFL, since the aggregation task is decentralized and different network topologies, communication protocols, and network optimization techniques play a key role.

Despite the merits of both centralized and decentralized FL, recent research has unveiled their vulnerabilities to adversarial attacks, particularly poisoning attacks that manipulate training data [16]. FL is highly susceptible to such attacks because traditional data inspection techniques cannot be used to effectively detect altered data, creating a challenge in balancing privacy protection with trust and robustness establishment [17]. Successful poisoning in FL can yield significant repercussions, including inaccurate predictions, misclassification of inputs, and the introduction of exploitable backdoors within compromised models [21]. Furthermore, these attacks can profoundly disrupt the collaborative learning process, eroding trust among participants and undermining the core privacy and security advantages of FL.

To defend against poisoning attacks in Centralized Federated Learning (CFL), researchers have actively explored various robust and secure techniques dealing with Byzantine-robust aggregation, anomaly detection, and hybrid approaches [19]. These defenses primarily concentrate on detecting and mitigating the detrimental effects of poisoning attacks, striving to uphold the integrity of FL and sustain the performance of individual participants' models. Nevertheless, some challenges emerge in the transition to DFL, where defense techniques are unexplored. When comparing DFL with CFL, DFL has the advantage of leveraging more information during the aggregation process, such as local data, weights, and models (apart from network topology information), which could be used to enhance the defenses effectiveness. In this sense, by comparing the similarity between the local model and the shared models and calculating the loss of the shared models on the local dataset, the DFL participant could improve the security of the model aggregation.

To improve the aforementioned limitations, this work introduces Sentinel, a hybrid defense strategy designed to enhance the resilience of DFL against poisoning attacks. Sentinel proposes a novel three-phase aggregation protocol to bolster security. In the initial phase, it employs a layer-wise average cosine similarity metric to identify and filter highly suspicious model updates. Subsequently, the remaining updates are aggregated based on local bootstrap validation loss, with aggregation weights determined through adaptive loss distance mapping. In the final phase, trusted models are normalized using the local model norm as a threshold to mitigate the impact of potential stealth attacks.

* Corresponding Author. Email: cfeng@uzh.ch.

Sentinel performance has been evaluated in a well-established DFL framework, Fedstellar [2], utilizing Computer Version (CV) datasets: MNIST, Fashion-MNIST, CIFAR10, Natural Language Processing (NLP) dataset: Sentiment140, and tabular dataset: Purchase, following IID (Independent and Identically Distributed) and non-IID data distributions. It undergoes evaluation under various data and model poisoning attacks in both targeted and untargeted forms, with varying numbers of attackers modifying different data portions. Results show high performance on IID, improving other state-of-the-art methods, while non-IID results remark the drawbacks of extensive filtering, something common to other robust aggregation methods.

The remainder of this document is organized as follows. First, Section 2 reviews previously proposed defense mechanisms for CFL and DFL. Section 3 introduces the design details of Sentinel, the proposed robust aggregation mechanism. Then, Section 4 evaluates and compares the performance of Sentinel with related work. Section 5 provides the discussion of the achieved results. Lastly, Section 6 summarizes the contributions of this work and proposes future opportunities.

2 Related Work

The most prominent methods to mitigate poisoning attacks are designed for CFL, where the model aggregation is done by a central server to reduce complexity. As indicated in [19], existing mitigation approaches are categorized into: Byzantine-robust aggregation, anomaly detection, and hybrid approaches.

Table 1: Techniques Against Poisoning Attacks. U Stands for Untargeted Attacks and T Stands for Targeted Attacks

Category	Type	Method	Technique	Schema	Obj.	
					U	T
Robust Aggregation	Geometry	COMED [24]	Coordinate-wise median	CFL	✓	x
		RFA [15]	Geometric median	CFL	✓	x
		TrimmedMean [23]	Filtered mean	CFL	✓	x
		Krum [3]	Euclidean distance	CFL	✓	x
		Multi-Krum [3]	Euclidean distance	CFL	✓	x
		Bulyan [12]	Krum and TrimmedMean	CFL	✓	x
Anomaly Detection	Validation	ERR, LFR [7]	Global validation	CFL	✓	✓
		PDGAN [26]	Model accuracy auditing	CFL	✓	✓
	Gradient	FLDetector [25]	Hessian-based gradient consistency	CFL	✓	✓
Hybrid		FLTrust [4]	ReLU-clipped cosine similarity, norm thresholding	CFL	✓	✓
		Gholami et al. [8]	Trusted aggregation	DFL	✓	x
		FLAME [13]	Clustering (cosine similarity), adaptive clipping, noising	CFL	✓	✓
This work		Sentinel	Model similarity, bootstrap validation and normalization	DFL	✓	✓

Byzantine-robust aggregation mechanisms try to prevent malicious updates deteriorating the model performance. Among the existing alternatives, Coordinate-wise Median is a defense method that applies dimension-wise filtering [24]. This approach is a generalization of the median in higher dimensions, is insensitive

to skewed distributions, and effectively defends against a model replacement attack. TrimmedMean [23] is similar to the coordinate-wise median but excludes the lowest and largest values in each dimension of the sorted model updates. RFA [15] is a simple alteration of FedAvg, where the mean aggregation is replaced with the geometric median. Krum [3] is a more sophisticated approach that assigns a score to each model update based on the similarity to other updates. Client scores are calculated by summing the Euclidean distance to the closest client updates. Bulyan [12] was proposed to address the limitations of Krum concerning the curse of dimensionality. However, experiments have demonstrated that in some scenarios, Bulyan performs worse than Krum due to potentially trimming benign updates.

Defense mechanisms through anomaly detection, also referred to as Byzantine detection, aim at identifying and removing potentially malicious updates. In contrast to Byzantine robustness, these anomaly detection schemes do not implement the defense strategy into the aggregation rules. [7] proposed two adapted defense strategies for CFL: Error Rate based Rejection (ERR) and Loss Function based Rejection (LFR). The key concept of these strategies is to evaluate the performance of the collected client models on a server-side validation dataset. Before averaging the received models, a predefined number of updates that have the largest impact on either the loss or the validation error are rejected. As a drawback, the server is required to collect a clean data set, which may violate the privacy-preserving concept of FL. To overcome that limitation, the authors of [26] proposed PDGAN. Instead of requiring a clean validation dataset, the system trains a generative adversarial network (GAN) concurrently to the original federated learning task. Once trained, PDGAN reconstructs a client's data from their model updates received on the server side. The generated data can be used to audit individual updates and label malicious clients as attackers. However, PDGAN can only defend from attackers after a certain number of iterations when the GAN is trained. Another approach is to investigate the gradient consistency of client updates. In this direction, FLDetector [25] predicts a client's update based on past contributions. Specifically, the server applies a quasi-newton approach to estimate the Hessian of an update and compare it to the integrated Hessian. FLDetector is effective against various adaptive attacks, such as distributed backdoors, and untargeted model poisoning. However, predicting consistency is computationally expensive.

Finally, hybrid defenses employ a combination of robust aggregation and anomaly detection mechanisms. [4] proposed FLTrust, a trusted aggregation mechanism for CFL based on comparing local model updates to a trusted bootstrapping model. Therefore, the server acquires a clean root dataset to train a reference model, and in each iteration, the received model updates are compared to the reference model. The result is then used as a trust score, and updates that deviate too much from the reference model receive a score of zero. [8] implemented the concept of trusted aggregation in the decentralized architecture. Each node is judged based on a behavioral score reflecting a participant's performance contribution and consistency. Computationally, this is represented by cluster-based and distance-based metrics. Each neighboring node's local trust score is then broadcast so that each node can compute a global trust score based on their neighbors' opinions. Finally, [13] proposed FLAME, a hybrid approach based on dynamic clustering, adaptive clipping, and noising. In the first step, the client updates are clustered based on pair-wise cosine similarity to capture angular deviation. The authors argue that existing clustering-based

approaches often group the models into malicious and benign. This leads to the issue that benign updates become removed when no adversaries are present.

In conclusion, as can be seen in Table 1, a noticeable research gap exists concerning defense mechanisms in the domain of DFL. Previous attempts prominently explored defenses against poisoning attacks in CFL, but only a few works investigated defense mechanisms in DFL. Therefore, while DFL has gained significant attention as a promising approach without a central entity, the security aspects of poisoning attacks still need to be addressed.

3 Sentinel

Sentinel is a defense mechanism that aims to defend against poisoning attacks in DFL. It considers local data availability and relies on the cosine similarity and bootstrap loss to evaluate the trust performance of neighbors' models. As can be seen in Figure 1, Sentinel is a three-phase aggregation protocol consisting of similarity filtering, bootstrap validation, and layer normalization.

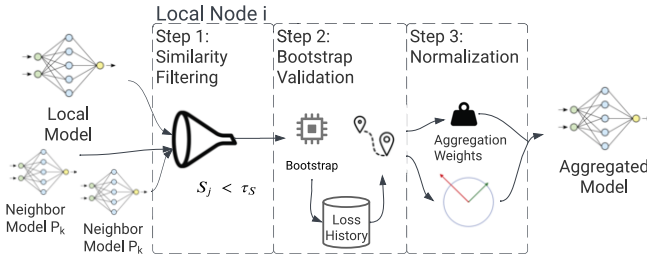


Figure 1: Overview of the Sentinel Aggregation Process.

The formal overview of Sentinel is given in Algorithm 1. The local model M is defined as the model trained on the individual node available in any round. Additionally, a model received from an adjacent node n_i is referred to as a neighbor model P_i . More information about the three steps of Sentinel are elucidated in the following paragraphs.

Algorithm 1 Sentinel Aggregation Algorithm

Require: P : Neighbor parameters, D_{bs} : Bootstrap dataset, H : Loss history, τ_S : Similarity threshold, τ_L : Loss distance threshold.

- 1: $r \leftarrow$ current round
- 2: $M \leftarrow$ local model
- 3: $w \leftarrow 0$
- 4: **for** j in P , $j \neq i$ **do** ▷ (1) Similarity filtering
- 5: $S_j \leftarrow \text{CosineSimilarity}(P_j, M)$
- 6: **if** $S_j < \tau_S$ **then** Remove P_j from P
- 7: $H_i[r] \leftarrow \text{ComputeBootstrapLoss}(M, D_{bs})$ ▷ (2) Bootstrap validation
- 8: **for** j in P , $j \neq i$ **do**
- 9: $l_j \leftarrow \text{ComputeBootstrapLoss}(P_j, D_{bs})$
- 10: $H_j[r] \leftarrow l_j$ ▷ Update loss history
- 11: $w_j \leftarrow \text{MapLossDistance}(H_i, H_j, \tau_L)$
- 12: **for** j in N , $j \neq i$ **do** ▷ (3) Layer Normalization
- 13: $\tilde{P}_j \leftarrow \text{NormaliseModel}(P_j, P_j)$
- 14: $\theta \leftarrow \frac{1}{\sum_{j \in N} w_j} \left(\sum_{j \in N} w_j \tilde{P}_j \right)$ ▷ Aggregate
- 15: **return** θ

Step 1: Similarity Filtering. At the first stage, Sentinel computes the layer-wise average cosine similarity between the current local model and all neighbor models received. This similarity is denoted as

S_j . The computation of the cosine similarity is defined in Algorithm 2. The cosine similarity is chosen due to its advantage over the Euclidean distance. In contrast to existing methods, Sentinel does not cluster or weigh models based on cosine similarity, but rather filters them instead. This step serves as a preliminary similarity evaluation to exclude potentially malicious models and reduce the computation of the second stage. The decision boundary is defined by τ_S , such that all models exhibiting a lower cosine similarity to the local model than τ_S will be filtered.

Algorithm 2 Cosine Similarity

Require: P : Neighbor parameters, M : Local model parameters

Ensure: $|M| = |P|$

- 1: **for** l_m, l_p in M, P **do**
- 2: $S_P \leftarrow S_P + \phi\left(\frac{\mathbf{v}_1 \cdot \mathbf{v}_2}{\|\mathbf{v}_1\| \|\mathbf{v}_2\|}\right)$ ▷ row-wise average
- 3: $S_P \leftarrow S_P / |M|$ ▷ layer-wise average
- 4: **return** S_P

Step 2: Bootstrap Validation. At the second step, the remaining models are evaluated on the basis of their loss. To do so, each node holds a small bootstrap dataset D_{bs} , which is randomly sampled from the validation dataset of the same node. The size of this bootstrap dataset is set to a third of the validation dataset or at least 300 samples, for example $|D_{bs}| = \max(|D_{val}|, 300)$. The number of collected samples required to effectively measure the loss has been investigated in related studies and is therefore not within the scope of this work [4]. In general, choosing the bootstrap dataset size is a trade-off between computational overhead and performance. Therefore, the chosen size should be suitable for the device specifications of each participating node. The computation of the bootstrap validation loss is defined in Algorithm 3. Note that the local node also computes the bootstrap validation loss for its own local model for a more accurate comparison.

Algorithm 3 Compute Bootstrap Loss

Require: P : Model parameters, D_{bs} : Bootstrap dataset

- 1: **for** each (x, y) in D_{bs} **do**
- 2: $y_{pred} \leftarrow$ Predict with P on x
- 3: $L \leftarrow \frac{1}{|D_{bs}|} \sum_{i=1}^{|D_{bs}|} l(y, y_{pred})$ ▷ Compute loss
- 4: **return** L

After computing the performance of each received model, Sentinel computes the average of the loss history up to the current aggregation round. This procedure simply takes all previously computed loss values into account. Thus, if the bootstrap loss was not computed in a previous round, since the model was filtered in step (1), there will be consequently fewer values to be averaged. Generally, the averaging serves for better computational stability. Finally, all average loss values are then compared to the local average bootstrap loss according to Algorithm 4, which results in an aggregation weight w_i .

Algorithm 4 Map Loss Distance

- 1: $\bar{l}_i \leftarrow \phi(H_i)$, $\bar{l}_j \leftarrow \phi(H_j)$
- 2: $\kappa \leftarrow \max(\bar{l}_i, l_{min})^{-1}$ ▷ Damping factor
- 3: $\bar{d}_i \leftarrow \max(\bar{l}_j - \bar{l}_i, 0)$
- 4: $w \leftarrow \exp(-k * \bar{d}_i)$
- 5: **if** $w < \tau_L$ **then** $w = 0$
- 6: **return** w

Sentinel maps the average loss distance according to a damped decay function. The damping factor κ is defined as the inverse of the average local loss. In theory, the local loss could be 0. Hence, a minimum average loss l_{min} must be defined for numerical stability, for example $l_{min} = 0.001$. Consequently, the local model and any neighbor model that presents a loss lower than the local model will receive $w = 1$. With this approach, Sentinel becomes more defensive as the average local loss decreases, which is expected during the FL process. However, if the local loss increases, which means the node is not learning, Sentinel becomes more exploratory.

Step 3: Layer Normalization. Inspired by recent approaches [4], the last step of Sentinel is the normalization of models to defend against potential stealth attacks that were able to pass defense layers (1) and (2). This procedure is defined in Algorithm 5. The normalization reduces the magnitude of potentially scaled attacks, which are commonly used to introduce backdoors. Such that Sentinel does not rely on a threshold, the ratio of the local model norm and the neighbor model norm is used as a scaling factor ρ . Sentinel only decreases the norm of neighbor models, hence $\rho \leq 1$. Subsequently, the normalized models are aggregated according to Algorithm 1 line 14.

Algorithm 5 Normalize Model

Require: P : Neighbor parameters, M : local model parameters

```

1:  $L \leftarrow$  number of layers
2: for  $l$  in  $|M|$  do
3:    $\rho \leftarrow \min(1, \frac{\|P[l]\|}{\|M[l]\|})$ 
4:    $\tilde{P}[l] \leftarrow \rho * P[l]$ 
5: return  $\tilde{P}$ 

```

4 Deployment and Experiments

Fedstellar [2] was chosen as the underlying platform to deploy and test Sentinel in DFL scenarios. Fedstellar allows users to train FL models in a decentralized, semi-decentralized, and centralized manner. It deploys the desired scenario to each specified client and manages the network of participants in terms of node connectivity. Users can run FL scenarios on selected physical devices or in a containerized simulation on Docker. Fedstellar is fully extensible, allowing users to implement custom models, load new datasets, and define the preferred aggregation algorithms.

4.1 Deployment Configuration

In this work, the Sentinel aggregation method (see Algorithm 1) was implemented in Python. PyTorch Lightning was used to implement the DL models. Then, for testing the Sentinel performance, a set of experiments was executed with the following configurations:

- The federation is composed of 10 fully connected nodes, virtualized with Docker, and acting as trainers and aggregators.
- Three types of datasets are employed to conduct experiments in this work. (i) Commonly used CV datasets MNIST [5], Fashion-MNIST [22], and CIFAR10 [10]; (ii) NLP sentiment analysis dataset Sentiment140 [9]; and (iii) tabular dataset Purchase [6].
- For the IID evaluation, data are equally distributed among all nodes. For MNIST and Fashion-MNIST, each node receives 6 000 training and 1 000 test samples; for CIFAR10, 5 000 training and 1 000 test samples; for Sentiment140, 8 500 training and 1 500

test samples; for Purchase, 7 384 training and 1 300 test samples; and the validation datasets are 10% of the training sets.

- For the non-IID evaluation, the same datasets are distributed across the participants by following a Dirichlet function [20] with $\alpha = 0.5$, and using the same ratio for training, validation, and testing.
- The chosen model topology for MNIST and Fashion-MNIST datasets is a Multi-Layer Perceptron (MLP) containing two fully-connected hidden layers with 256 and 128 neurons.
- For CIFAR10, the model topology is a Convolutional Neural Network (CNN), namely SimpleMobileNet [18], which contains several convolutional layers followed by pooling layers and fully connected layers towards the end. The first layer uses a set of 32 filters with a kernel size of 3x3. This is followed by a max pooling layer to reduce the spatial dimensions. Then, five depthwise separable convolutional layers with different input and output channel sizes and strides. These layers help reduce the model computational cost while preserving expressive power. Finally, global pooling in the form of an adaptive average pooling layer is applied to reduce the spatial dimensions of the output feature map to a size of 1×1 . This vector feeds into a dense layer of 512 units.
- For Sentiment140, a bidirectional Long Short-Term Memory (LSTM) network is used for the model topology. In the preprocessing of the data, the original text data is firstly tokenized using the lexicon of the pre-trained word embedding of Glove's Twitter.27B [14], and then the tokens are transformed into a 100-dimensional word vector. In order to make the length of each input vector the same, all samples are pad to 64 dimensions, i.e., after preprocessing, all the text is transformed into a 64*100 tensor. In terms of the model topology, after the input layer, a bi-directional LSTM with a hidden layer of 256 is utilized, and finally, a fully connected layer with an output layer of 2 is connected.
- MLP is employed as a model for the Purchase dataset, which contains only a 256-neuron hidden layer.
- In all model topologies, *ReLU* is employed as an activation function in the hidden layers and *softmax* in the output layer. *Adam* is employed as an optimizer with default hyperparameters.
- The duration of the learning process is configured to 10 rounds of 3 epochs. To make experiments reproducible, synchronous aggregation is implemented in Fedstellar.
- The federation network bandwidth is configured to 1 Mbps, the loss to 0%, and the delay to 0 ms.
- The ratio of poisoned nodes (PNR) is established to 10%, 50%, and 80%. Malicious nodes are randomly selected, and they launch model poisoning, label flipping (targeted and untargeted), and backdoor attacks.
- FedAvg, FLTrust, Krum, TrimmedMean, and Sentinel are implemented in Fedstellar, evaluated, and compared. FedAvg acts as a baseline, and the other methods act as defense mechanisms based on robust aggregation.

4.2 Experiments

To compare the performance of Sentinel with existing defense mechanisms, a baseline performance reference under benign settings is established. Subsequently, a set of experiments evaluates the performance of each defense for the previous datasets and attacks.

Baseline. The goal of this experiment is to establish a baseline performance that serves as a reference for subsequent experiments

in malicious environments. The evaluation of baseline, model poisoning, and untargeted label flipping attacks employs the F1-Score.

As can be seen in Table 2, Sentinel achieves a similar F1-Score to related work for the five datasets when the dataset follows an IID distribution. In the case of non-IID, the F1-Score performance degrades compared to FedAvg in ≈ 0.05 , 0.10, and 0.15 for MNIST, Fashion-MNIST and CIFAR10, and a decrement of 0.10 for both Sentiment140 and Purchase. However, this decrease is consistent with other robust aggregation mechanisms such as Krum. These results demonstrate that Sentinel is aligned with other robust mechanisms and is a valid option when no attacks affect the federation nodes.

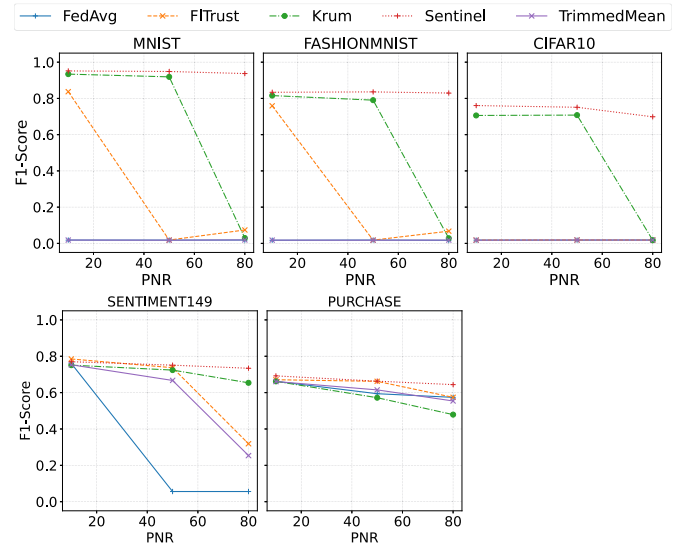
Table 2: F1-Score for MNIST, Fashion-MNIST, CIFAR10, Sentiment140, and Purchase without Attacks

	MNIST	Fashion-MNIST	CIFAR10	Sentiment140	Purchase
Balanced data (IID)					
FedAvg	0.953 \pm 0.037	0.838 \pm 0.027	0.764 \pm 0.022	0.784 \pm 0.012	0.760 \pm 0.092
FLTrust	0.952 \pm 0.037	0.837 \pm 0.022	0.764 \pm 0.025	0.785 \pm 0.011	0.761 \pm 0.061
Krum	0.935 \pm 0.046	0.822 \pm 0.021	0.671 \pm 0.029	0.744 \pm 0.016	0.781 \pm 0.022
TrimmedMean	0.952 \pm 0.039	0.835 \pm 0.025	0.762 \pm 0.018	0.789 \pm 0.009	0.708 \pm 0.050
Sentinel	0.951 \pm 0.041	0.834 \pm 0.025	0.754 \pm 0.018	0.763 \pm 0.031	0.750 \pm 0.082
Unbalanced data (Non-IID Dirichlet($\alpha = 0.5$))					
FedAvg	0.941 \pm 0.014	0.836 \pm 0.006	0.599 \pm 0.044	0.660 \pm 0.012	0.653 \pm 0.009
FLTrust	0.884 \pm 0.045	0.770 \pm 0.039	0.459 \pm 0.053	0.573 \pm 0.017	0.636 \pm 0.013
Krum	0.908 \pm 0.009	0.741 \pm 0.031	0.375 \pm 0.014	0.547 \pm 0.018	0.661 \pm 0.010
TrimmedMean	0.963 \pm 0.009	0.833 \pm 0.011	0.636 \pm 0.013	0.668 \pm 0.013	0.661 \pm 0.010
Sentinel	0.894 \pm 0.072	0.727 \pm 0.042	0.448 \pm 0.051	0.666 \pm 0.012	0.652 \pm 0.035

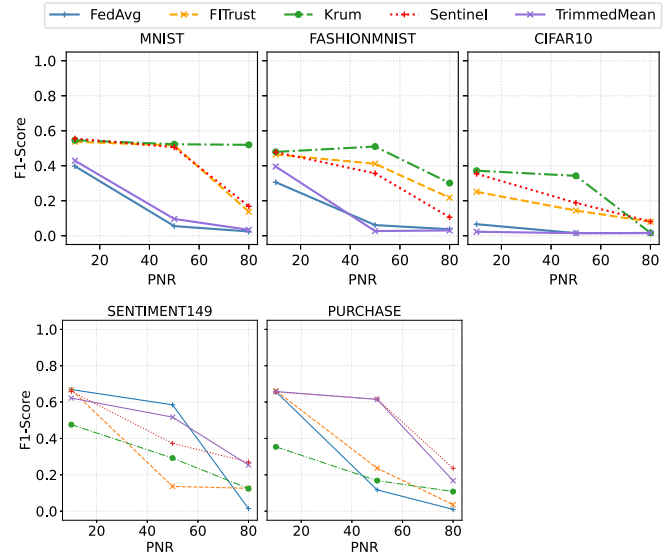
Model Poisoning. In this attack, the weights are randomly altered with salt noise before the model is sent to their neighbors at each round. The noise ratio (NR) remained fixed at 80% for all configurations. Instead, the poisoned node ratio (PNR) was varied with the ratios 10%, 50%, and 80% for all datasets. Figure 2a shows the F1-Score for the five datasets and each PNR configuration under IID distribution. As can be seen, Sentinel performs the best, with a maximum F1-Score reduction of 0.04 from the baseline performance with FedAvg, observed on CIFAR10 and a PNR of 80%. Sentinel shows the same high performance not only in the CV datasets but also in the Sentiment140 and Purchase datasets. In contrast, Figure 2b shows that Krum performs better when the distribution is non-IID. This is caused by a higher filtering rate from Sentinel compared to other methods, which could affect non-IID clients that are legitimate. Still, the results prove that Sentinel performance in non-IID is close to FLTrust and improves FedAvg and TrimmedMean aggregation.

Untargeted Label Flipping. In this experiment, adversaries modify all the local training data by randomly altering the labels of their data samples. The aim is to degrade the performance of the model from benign actors, provoking random mispredictions. As can be seen in Figure 3a, Sentinel is the unique defense able to maintain its baseline F1-Score for any attack configuration and dataset when the distribution is IID. More in detail, with 80% of malicious participants, the performance of the rest of the defense mechanisms is destroyed. However, in the case of non-IID distribution, the robust aggregation mechanisms are the ones performing the worse, especially in the case of Sentinel, as Figure 3b shows. This can be caused by the model filtering in robust aggregation mechanisms, which could filter honest clients that generate more varied models due to their local data distribution.

To conclude, in untargeted attacks, Sentinel is effective against many types of attacks, including model poisoning and label flipping attacks. This effectiveness is well validated on different types of



(a) F1-Score Under Model Poisoning Attacks on IID distributions



(b) F1-Score Under Model Poisoning Attacks on non-IID distributions (Dirichlet $\alpha = 0.5$)

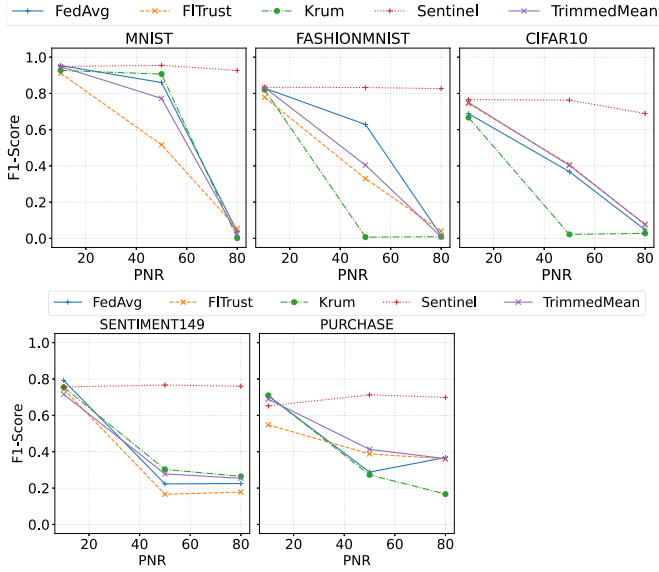
Figure 2: F1-Scores Under Model Poisoning Attacks

datasets, including CV, NLP, and tabula datasets.

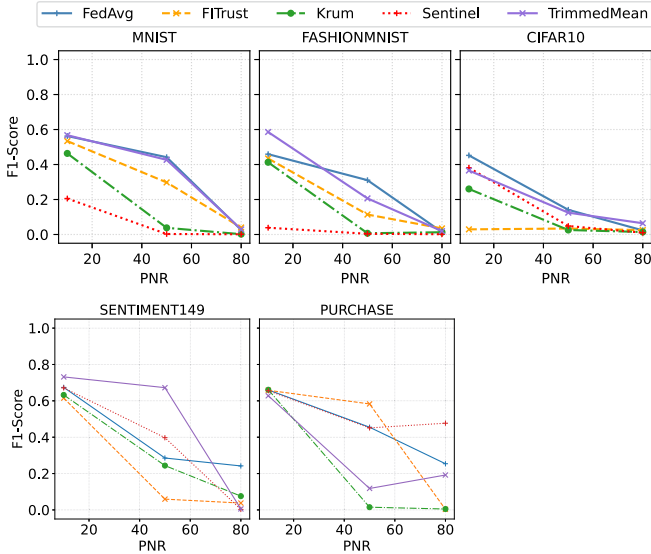
Targeted Label Flipping. The attack causes a misclassification of a selected source label l_{src} to a target label l_t . The success of this poisoning attack is measured by the amount of test data samples labeled with the source label classified as the target label by the final models of benign participants. This ratio is represented by the Attack Success Rate for label flipping (ASR_{lf}), as shown in (1), where c_{ij} is the number of samples having true label y_i and predicted label \hat{y}_j . L represents the set of labels in the corresponding dataset.

This experiment was performed only in MNIST, Fashion-MNIST, and CIFAR10. These three datasets are ten classification tasks and, thus, can be attacked by targeting specific labels. However, Sentiment140 and Purchase are classification tasks, so the targeted label flipping attack is not applicable in this case.

$$ASR_{lf} = \frac{c_{src,t}}{\sum_{j=0}^{|L|} c_{src,j}} \quad (1)$$



(a) F1-Score Under Untargeted Label Flipping Attacks on IID distribution

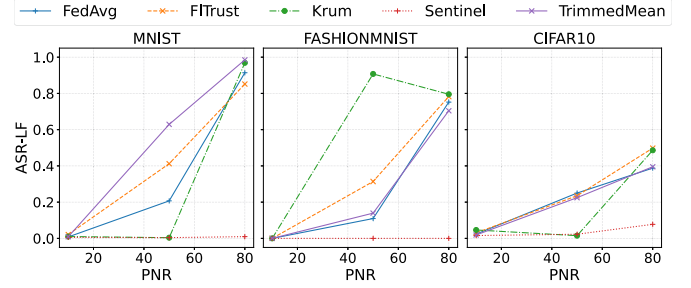


(b) F1-Score Under Untargeted Label Flipping Attacks on non-IID distribution

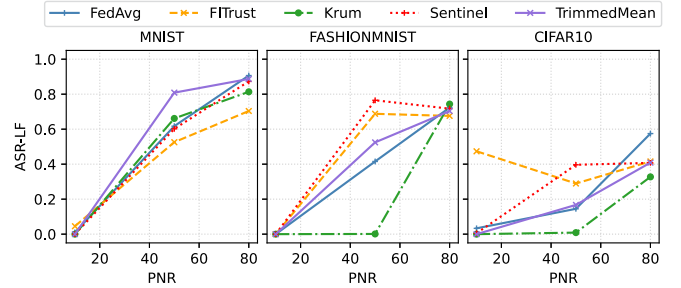
Figure 3: F1-Scores Under Untargeted Label Flipping Attack

As can be seen in Figure 4a, Sentinel performs the best for all datasets and PNR when the distribution is IID. This can be explained by the exclusion of untrusted nodes and thereby preventing an aggregation of malicious, but highly similar, poisoned models. However, as Figure 4b shows, and similarly to the Untargeted setup, the robust aggregation methods do not perform well when the distribution is non-IID, again especially for Sentinel. The reason behind this lack of performance is the extensive filtering of different models based on distance, which can degrade performance under non-IID. Note that this problem is present in all robust aggregation mechanisms, not only for Sentinel.

Backdoor. In this attack, malicious participants try to provoke a misclassification to a predefined target label l_t by using an artificial trigger. In this work, this trigger is represented by the form of an “X” on any image of the three CV datasets (MNIST, Fashion-MNIST, and CIFAR10). In Sentiment140, 20% of the training data with a



(a) ASR Under Targeted Label Flipping Attacks on IID distribution



(b) ASR Under Targeted Label Flipping Attacks on non-IID distribution

Figure 4: ASR-LF (Attack Success Rate - Label Flipping) Under Targeted Label Flipping Attack

sentiment of positive are randomly selected, and their first three-word embeddings are changed to 0 vectors. In Purchase, 20% of the data in labeling ‘1’ are selected and their first seven dimensions of the data are changed to 1.

Backdoor Accuracy (BA) is used for evaluation of defenses, as shown in (2), where c_{ij} is the number of samples having true label y_i and predicted label \hat{y}_j . L represents the set of labels in the corresponding dataset, and B is the backdoor dataset.

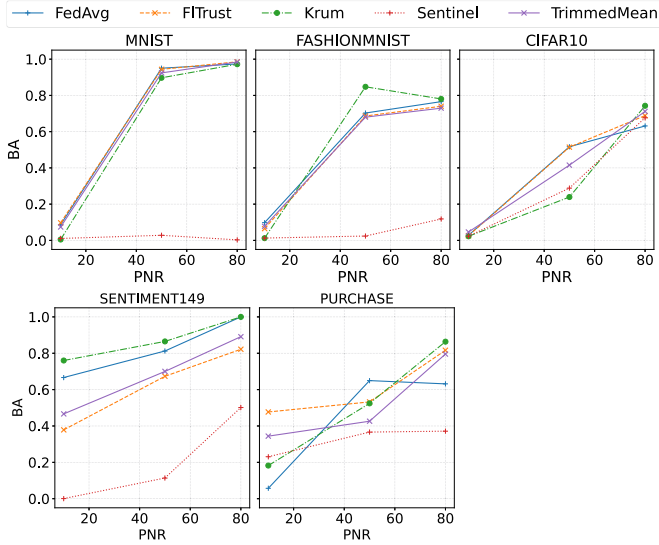
$$BA = \frac{\sum_{j=0}^{|L|} c_{j,t} - c_{t,t}}{|B| - c_{t,t}} \quad (2)$$

As can be seen in Figure 5a, Sentinel is the most robust solution when the distribution is IID, performing almost perfectly for MNIST, Fashion-MNIST, Sentiment140, and Purchase (regardless of the PNR). In the CIFAR10 case, Sentinel achieves similar performance to the state of the art. In the case of non-IID, the results in Figure 5b show that the backdoor attack is not effective in the non-IID setup, probably due to the variations in the distributions in each local client dataset.

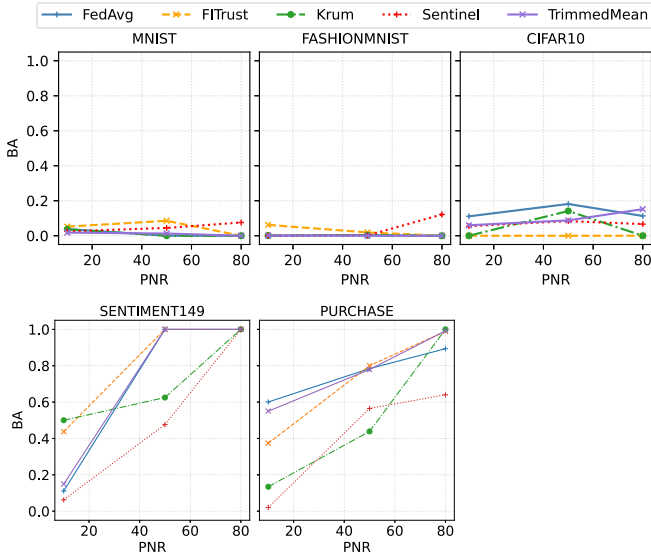
5 Discussion

As experiments demonstrated, the resilience of Sentinel against a range of attacks in IID scenarios improved state-of-the-art methods. Therefore, Sentinel is a reliable defense mechanism in controlled environments. However, the fluctuating performance in non-IID settings suggests that further refinements are needed to enhance its adaptability and discrimination capabilities in more complex data distributions.

In IID, Sentinel, which sequentially integrates three distinct defensive mechanisms, effectively filters out malicious updates before aggregation. This multi-phase approach minimizes the impact of poisoning and ensures that only trustworthy model updates contribute to the aggregated global model. The experimental results



(a) BA Under Backdoor Attacks on IID distribution



(b) BA Under Backdoor Attacks on non-IID distribution

Figure 5: BA (Backdoor Accuracy) Under Backdoor Attacks on Different Distributions

validate that Sentinel can maintain integrity and performance even under aggressive attack scenarios, such as model poisoning and targeted label flipping.

While Sentinel effectively secures federated learning networks in IID scenarios by filtering out anomalous updates, its performance under non-IID conditions requires careful consideration. The distinct data characteristics inherent to non-IID distributions can lead to high variability in local model updates. This variability could be misinterpreted by Sentinel similarity filtering phase as potentially malicious, leading to the unwarranted exclusion of legitimate model updates. Such exclusions reduce the overall model quality and affect learning efficiency and convergence rates.

Furthermore, the bootstrap validation step in Sentinel, which assesses model updates based on their loss performance on a small, locally held validation dataset, can penalize nodes that genuinely represent minority patterns within the data. Since these nodes naturally exhibit higher loss values due to their divergent data

characteristics, their contributions might be undervalued in the aggregation process. This scenario underscores a critical trade-off in Sentinel design between robustness to attacks and sensitivity to the legitimate diversity of participant data.

To better accommodate non-IID distributions, future iterations of Sentinel could benefit from incorporating adaptive mechanisms that adjust the thresholds for similarity and loss validation based on the observed data distribution variance among nodes. For instance, dynamic adjustment of the similarity threshold (τ_S) and the loss distance threshold (τ_L) could help mitigate the risk of excluding non-malicious yet statistically unusual data patterns. Additionally, implementing anomaly detection techniques that consider the context of data distribution or employing unsupervised learning to better understand the typical behavior of node updates could enhance Sentinel capability to discern between attacks and naturally occurring data anomalies.

6 Conclusion and Future Work

This work has proposed Sentinel, a sophisticated defense strategy that applies a multi-level defense protocol composed of similarity filtering, bootstrap validation, and model normalization to mitigate poisoning attacks. By taking advantage of the local data availability in DFL, Sentinel demonstrates the transferability of promising defense strategies proposed for CFL. Extensive evaluations on the datasets MNIST, Fashion-MNIST, CIFAR10, Sentiment140, and Purchase datasets with various attack configurations have demonstrated their effectiveness. Compared to other state-of-the-art aggregation algorithms, the aggregation protocol proposed in this work reliably defended against poisoning attacks under IID configurations, especially in highly malicious environments. Furthermore, this work demonstrated the usability of Fedstellar for more comparable security benchmarks in DFL. In contrast, under non-IID configurations Sentinel shows some disadvantages compared to traditional non-robust aggregation methods, drawbacks that are also present in other robust methods.

As future work, it is planned to optimize the Sentinel algorithm better to handle the variability and unpredictability of non-IID data. Exploring hybrid models that combine Sentinel robust aggregation with other techniques could potentially yield methods that are robust to attacks and sensitive to legitimate data diversity. Other relevant directions are evaluating scenarios involving alternative network topologies, such as ring or random graphs, with or without clustered components. Finally, evaluating the impact of different Dirichlet values in the non-IID is another experiment for the future.

Acknowledgments

This work has been partially supported by (a) the Swiss Federal Office for Defense Procurement (armasuisse) with the CyberForce (CYD-C-2020003) project, (b) the University of Zürich UZH, (c) 21629/FPI/21, Fundación Séneca, Región de Murcia (Spain), and (d) the strategic project DEFENDER from the Spanish National Institute of Cybersecurity (INCIBE) and by the Recovery, Transformation and Resilience Plan, Next Generation EU.

References

- [1] M. Alazab, S. P. Rm, P. M, P. K. R. Maddikunta, T. R. Gadekallu, and Q.-V. Pham. Federated Learning for Cybersecurity: Concepts, Challenges, and Future Directions. *IEEE Transactions on Industrial Informatics*, 18(5):3501–3509, May 2022. ISSN 1551-3203, 1941-0050. doi: 10.1109/TII.2021.3119038.

- [2] E. T. M. Beltrán, Á. L. P. Gómez, C. Feng, P. M. S. Sánchez, S. L. Bernal, G. Bovet, M. G. Pérez, G. M. Pérez, and A. H. Celdrán. Fedstellar: A platform for decentralized federated learning. *Expert Systems with Applications*, 242:122861, 2024.
- [3] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer. Machine learning with adversaries: Byzantine tolerant gradient descent. In *Proceedings of the 31st international conference on neural information processing systems*, NIPS'17, pages 118–128, NY, USA, 2017. ISBN 978-1-5108-6096-4.
- [4] X. Cao, M. Fang, J. Liu, and N. Z. Gong. FLTrust: Byzantine-robust Federated Learning via Trust Bootstrapping. In *Proceedings 2021 Network and Distributed System Security Symposium*, Virtual, 2021. ISBN 978-1-891562-66-2. doi: 10.14722/ndss.2021.24434.
- [5] L. Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- [6] W. C. DMDave, Todd B. Acquire valued shoppers challenge, 2014. URL <https://kaggle.com/competitions/acquire-valued-shoppers-challenge>.
- [7] M. Fang, X. Cao, J. Jia, and N. Z. Gong. Local Model Poisoning Attacks to Byzantine-Robust Federated Learning, Nov. 2021. arXiv:1911.11815.
- [8] A. Gholami, N. Torkzaban, and J. S. Baras. Trusted Decentralized Federated Learning. In *2022 IEEE 19th Annual Consumer Communications & Networking Conference (CCNC)*, pages 1–6, NV, USA, Jan. 2022. IEEE. ISBN 978-1-66543-161-3. doi: 10.1109/CCNC49033.2022.9700624.
- [9] A. Go, R. Bhayani, and L. Huang. Twitter sentiment classification using distant supervision. *CS224N project report, Stanford*, 1(12):2009, 2009.
- [10] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [11] E. T. Martínez Beltrán, M. Quiles Pérez, P. M. Sánchez Sánchez, S. López Bernal, G. Bovet, M. Gil Pérez, G. Martínez Pérez, and A. Huertas Celdrán. Decentralized federated learning: Fundamentals, state of the art, frameworks, trends, and challenges. *IEEE Communications Surveys & Tutorials*, pages 1–1, In press. doi: 10.1109/COMST.2023.3315746.
- [12] E. M. E. Mhamdi, R. Guerraoui, and S. Rouault. The hidden vulnerability of distributed learning in byzantium. *arXiv preprint arXiv:1802.07927*, 2018.
- [13] T. D. Nguyen, P. Rieger, H. Chen, H. Yalame, H. Möllering, H. Fereidooni, S. Marchal, M. Miettinen, A. Mirhoseini, S. Zeitouni, F. Koushanfar, A.-R. Sadeghi, and T. Schneider. FLAME: Taming Backdoors in Federated Learning, 2021. Report Number: 025.
- [14] J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [15] K. Pillutla, S. M. Kakade, and Z. Harchaoui. Robust Aggregation for Federated Learning, Jan. 2022. arXiv:1912.13445.
- [16] N. Rodríguez-Barroso, D. Jiménez-López, M. V. Luzón, F. Herrera, and E. Martínez-Cámara. Survey on federated learning threats: Concepts, taxonomy on attacks and defences, experimental study and challenges. *Information Fusion*, 90:148–173, Feb. 2023. ISSN 15662535. doi: 10.1016/j.inffus.2022.09.011.
- [17] P. M. S. Sánchez, A. H. Celdrán, N. Xie, G. Bovet, G. M. Pérez, and B. Stiller. Federatedtrust: A solution for trustworthy federated learning. *Future Generation Computer Systems*, 152:83–98, 2024.
- [18] D. Sinha and M. El-Sharkawy. Thin mobilenet: An enhanced mobilenet architecture. In *2019 IEEE 10th annual ubiquitous computing, electronics & mobile communication conference (UEMCON)*, pages 0280–0285. IEEE, 2019.
- [19] Z. Tian, L. Cui, J. Liang, and S. Yu. A Comprehensive Survey on Poisoning Attacks and Countermeasures in Machine Learning. *ACM Computing Surveys*, 55(8):1–35, Aug. 2023. ISSN 0360-0300, 1557-7341. doi: 10.1145/3551636.
- [20] Y. Wang, Y. Tong, and D. Shi. Federated latent dirichlet allocation: A local differential privacy based framework. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 6283–6290, 2020.
- [21] G. Xia, J. Chen, C. Yu, and J. Ma. Poisoning Attacks in Federated Learning: A Survey. *IEEE Access*, 11:10708–10722, 2023. ISSN 2169-3536. doi: 10.1109/ACCESS.2023.3238823. Conference Name: IEEE Access.
- [22] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [23] D. Yin, Y. Chen, R. Kannan, and P. Bartlett. Byzantine-Robust distributed learning: Towards optimal statistical rates. In *Proceedings of the 35th international conference on machine learning*, volume 80, pages 5650–5659, 2018.
- [24] J. Yin, X. Cui, and K. Li. A Reputation-Based Resilient and Recoverable P2P Botnet. In *2017 IEEE Second International Conference on Data Science in Cyberspace (DSC)*, pages 275–282, June 2017. doi: 10.1109/DSC.2017.20.
- [25] Z. Zhang, X. Cao, J. Jia, and N. Z. Gong. FLDetector: Defending Federated Learning Against Model Poisoning Attacks via Detecting Malicious Clients, Oct. 2022. arXiv:2207.09209.
- [26] Y. Zhao, J. Chen, J. Zhang, D. Wu, J. Teng, and S. Yu. PDGAN: A Novel Poisoning Defense Method in Federated Learning using Generative Adversarial Network. In *Algorithms and Architectures for Parallel Processing: 19th International Conference*, pages 595–609, 2020.