



ARFED: Attack-Resistant Federated averaging based on outlier elimination^{☆,☆☆}

Ece Isik-Polat^{*}, Gorkem Polat, Altan Kocyigit

Graduate School of Informatics, Middle East Technical University, Universiteler, Dumlupinar Blv. 1/6 D:133, Cankaya, Ankara, 06800, Turkey



ARTICLE INFO

Article history:

Received 20 April 2022

Received in revised form 9 October 2022

Accepted 2 December 2022

Available online 8 December 2022

Dataset link: <https://github.com/eceisik/ARFED>

Keywords:

Federated learning

Data poisoning

Model poisoning

Label flipping attacks

Byzantine attacks

Adaptive attacks

ABSTRACT

In federated learning, each participant trains its local model with its own data and a global model is formed at a trusted server by aggregating model updates coming from these participants. Since the server has no effect and visibility on the training procedure of the participants to ensure privacy, the global model becomes vulnerable to attacks such as data poisoning and model poisoning. Although many defense algorithms have recently been proposed to address these attacks, they often make strong assumptions that do not agree with the nature of federated learning, such as assuming Non-IID datasets. Moreover, they mostly lack comprehensive experimental analyses. In this work, we propose a defense algorithm called ARFED that does not make any assumptions about data distribution, update similarity of participants, or the ratio of the malicious participants. ARFED mainly considers the outlier status of participant updates for each layer of the model architecture based on the distance to the global model. Hence, the participants that do not have any outlier layer are involved in model aggregation. We have performed extensive experiments on diverse scenarios and shown that the proposed approach provides a robust defense against different attacks. To test the defense capability of the ARFED in different conditions, we considered label flipping, Byzantine, and partial knowledge attacks for both IID and Non-IID settings in our experimental evaluations. Moreover, we proposed a new attack, called organized partial knowledge attack, where malicious participants use their training statistics collaboratively to define a common poisoned model. We have shown that organized partial knowledge attacks are more effective than independent attacks.

© 2022 Elsevier B.V. All rights reserved.

1. Introduction

Digitalization has been reshaping the economy, organizations, and society [1]. It is driving the extensive use of artificial intelligence techniques to increase efficiency and productivity, lower costs, and improve the quality of products and services in all industries and sectors [2]. Recent advances in computing and communications technologies, widespread deployment of smart and connected devices, and the proliferation of cloud computing enable the collection and cost-effective storage and processing of massive data that is essentially the fuel for successful organizations today. Hence, there has been a great interest in the Internet of Things (IoT) and Big Data concepts [3]. In this context,

machine learning is one of the primary techniques to extract non-obvious and useful patterns and actionable insight from data. Specifically, deep learning offers enormous potential and has achieved remarkable success [4]. Machine learning in general, and deep learning in particular, is a computation-intensive task that processes data usually collected from many sources and stored in a central location. However, with big data, collecting, storing, and processing data in a scalable and efficient manner is a fundamental challenge [5]. Besides, the performance and adequacy of a trained model mostly rely on the availability of sufficiently large data relevant to the task of interest. Hence, data collected from various sources such as online transaction systems, IoT devices, smartphones, and social media are integrated to have a rich dataset to derive high-quality models. Such data generally contain sensitive information of which collection and use may cause violation of regulations such as the General Data Protection Regulation (GDPR) [6–8]. Hence, data privacy turns out to be a fundamental challenge.

Federated Learning (FL) [9] is a distributed approach to training a machine learning model without requiring training data to be available in a central place. In FL, participants with relevant data and processing resources collaborate to train a machine

[☆] You can visit <https://github.com/eceisik/ARFED> to see all implemented methods and designed experiments.

^{☆☆} The numerical calculations reported in this paper were partially performed using TUBITAK ULAKBIM, High Performance and Grid Computing Center (TRUBA resources).

^{*} Corresponding author.

E-mail addresses: eceisik@metu.edu.tr (E. Isik-Polat), gorkem.polat@metu.edu.tr (G. Polat), kocyigit@metu.edu.tr (A. Kocyigit).

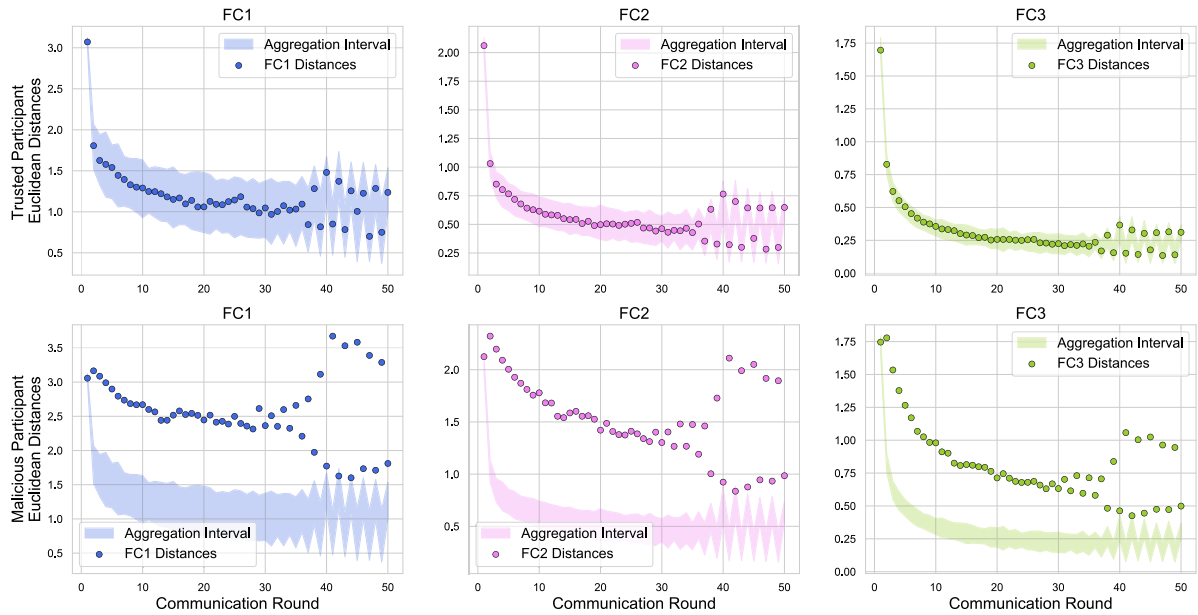


Fig. 1. The layer-wise distances of a randomly selected trusted participant (top row) and a randomly selected malicious participant (bottom row) to the global model under label flipping attack on the MNIST dataset. The points refer to the distances of the participants to the same layer of the global model in corresponding rounds. The shaded regions indicate the upper and lower bounds determined for each round according to the IQR outlier identification technique.

learning model without revealing their data. In a typical FL setting, a trusted server communicates with the participants to jointly train a common model in several iterations. The server chooses a model architecture, determines training parameters, and initializes a global model, which is iteratively improved by performing local training on participants' devices. In each iteration, the server sends the global model to the participants, which use their training data to improve the model for a while and send the resulting local models back to the server. Then the server aggregates received models to have a better global model. These iterations can be continued until some convergence criteria are met. FL is a viable approach to overcome data privacy issues as the participants do not need to disclose their training data; instead, they only share locally trained model parameters. Hence, FL is a promising approach to large-scale application of machine learning in domains where protecting user's privacy is of great concern, such as healthcare [10], smart city [11], and others widely employing IoT and big data technologies [12].

FL has some unique characteristics [13]. Training data may be massively distributed onto a large number of devices with heterogeneous resources, and the sizes and distributions of data on different devices may vary considerably. FL aims to train a single model that can perform well on all participants' data. However, this may not be possible when participants have heterogeneous data and processing resources [14]. Not independent and identically distributed (IID) data on devices lead to severe issues in FL, and improving the performance of FL on non-IID data is an active research area [15]. This statistical heterogeneity affects the convergence behavior of FL and may lead to biased models toward the participants having larger training data. The heterogeneity of storage, computation, and communication resources of participants, expensive communication (especially when there are a massive number of participants), and ensuring privacy against inference attacks are other core challenges in FL [16]. The presence of adversaries manipulating their data or locally trained models exacerbates the problems caused by the heterogeneity of data and resources [14]. Therefore, FL is also vulnerable to security attacks as the central server has no control over the

participants' data and local training processes. There are several vulnerabilities that an attacker might exploit to manipulate the learning process, manipulate the global model, or gain access to participants' private data [17]. For instance, an attacker may pretend like an ordinary participant(s) or the central server or gain control of one or more participants or the central server to target the learning process without getting noticed. The attacker's goal may be to slow down or impede the convergence of model training, degrade the trained model's performance, manipulate the global model to get wrong inferences under specific cases, lead to an ineffective global model, or extract participants' local data from the parameters exchanged during training rounds. Data poisoning [18] and model poisoning [19] are two significant security threats that attackers can pose in FL. In data poisoning attacks, malicious participants manipulate their training data by adding noise or flipping target labels. In model poisoning attacks, participants alter their models before sending them to the server.

Several aggregation approaches and optimization algorithms, mainly variations of the gradient descent algorithm, have been proposed for model training with FL [13,20–22]. Federated Averaging (FedAvg) [9] is one of the most commonly used FL algorithms. In each iteration, FedAvg aggregates the locally trained models returned by the participants to form the new global model by averaging. Each parameter of the new global model is set to the weighted arithmetic mean of the corresponding parameters of the participants' models. In this process, weights are determined according to the number of training examples in the participants to give each training example an equal weight on the global model update. This feature represents a vital vulnerability if there are malicious participants sending arbitrary or deliberately manipulated model parameters, which are likely to be very different from the updates received from reliable participants. Hence, checking the models before aggregation to identify anomalies and dropping (or lowering the weights of) the models coming from suspicious participants in the aggregation could be a promising approach to deal with poisoning attacks.

In this paper, we propose the Attack-Resistant Federated Learning (ARFED) algorithm, which is an extended version of

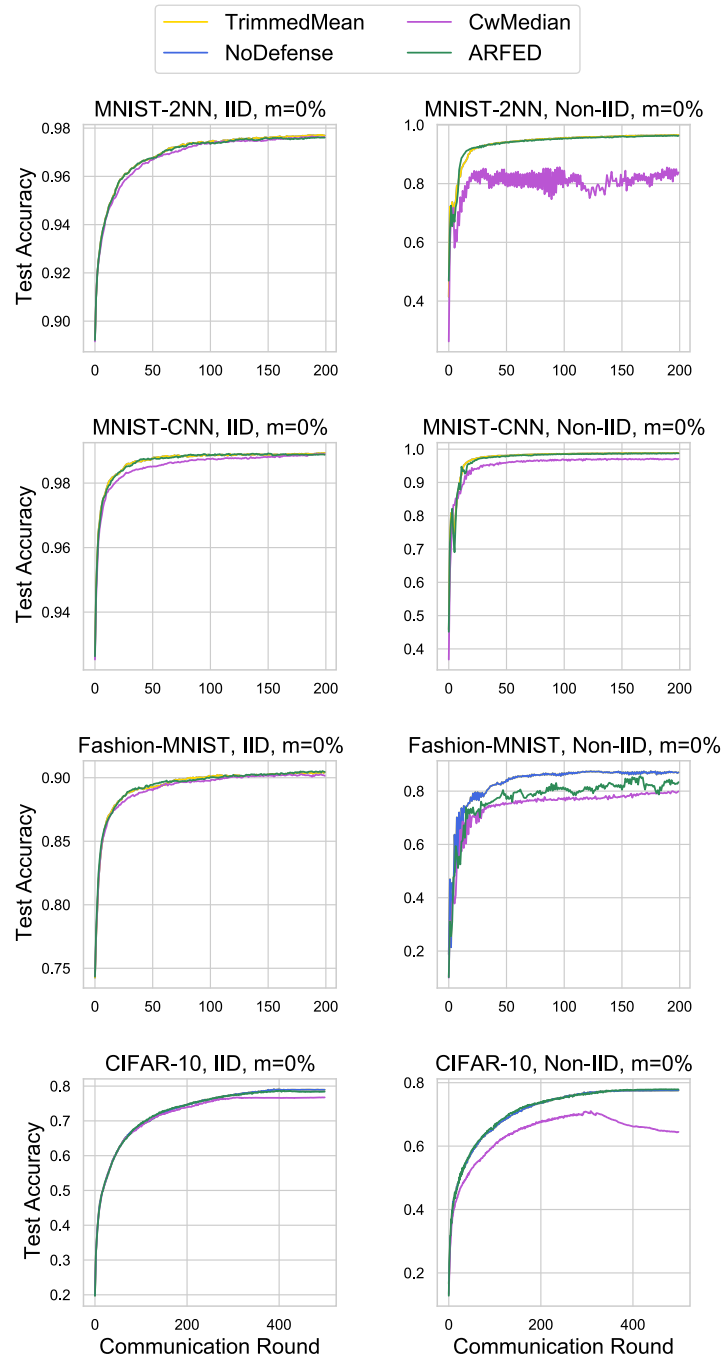


Fig. 2. Accuracy curves of different strategies when all participants are trusted.

FedAvg. The primary objective of the algorithm is to defend against poisoning attacks in FL. In ARFED, the parameters of the models received from the participants in each iteration are examined using a statistical outlier detection technique to identify potentially malicious participants. Accordingly, such participants are discarded in the model aggregation step to mitigate poisoning threats. Many defense methods for poisoning attacks are proposed in the literature, as summarized in Section 2. These methods usually require some knowledge about the attacks, such as malicious participant ratio, examining local datasets, which may compromise the privacy of participants, assuming IID data, which is not valid in typical FL settings or introducing

too much computational overhead, which restricts the practical implementations. Unlike the other defense methods, ARFED employs a relatively simple malicious participant identification technique that does not require making unrealistic assumptions inconsistent with typical FL settings.

The main contributions of this paper are as follows:

- (a) We propose an extension to FedAvg called ARFED to defend against poisoning attacks. Unlike the other similar algorithms proposed in the literature, ARFED does not make unrealistic assumptions about data distributions or participants' update similarities. Nor does it require information

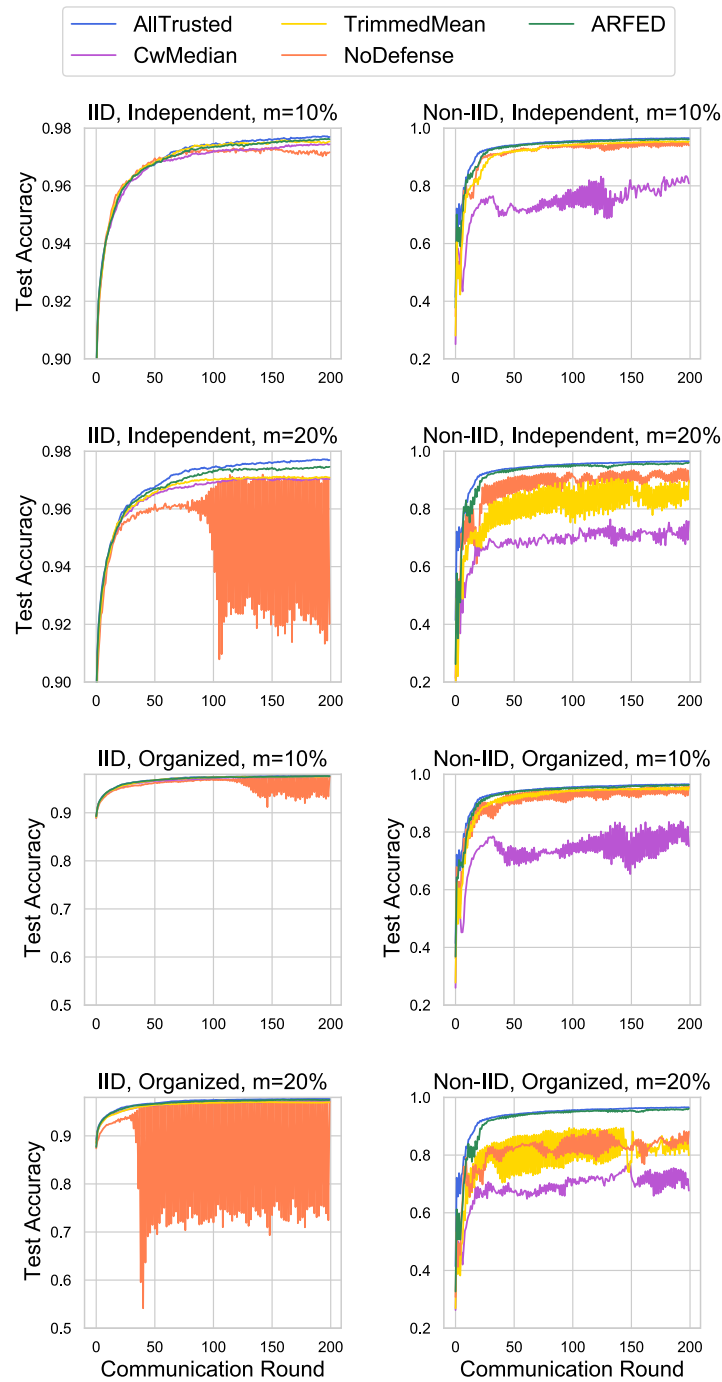


Fig. 3. Accuracy curves of different strategies for MNIST-2NN under label flipping attacks at different attacker ratios.

about attacks, such as the malicious participant ratios, in advance. As shown in Section 3, the computational complexity of the extension made to FedAvg is lower compared to other similar defense algorithms.

- (b) We evaluated the performance of the vanilla FedAvg, ARFED, and two similar defense approaches in various FL scenarios, such as IID data, non-IID data, and under various kinds of organized and independent poisoning attacks, as well as the no-attack case. The results show that attacks in non-IID cases are more severe than in IID cases. Moreover, the attacks committed by an organized group of attackers

can be more detrimental than those implemented by a group of independent attackers.

- (c) The experimental results show that ARFED can mitigate the effects of independent and organized attacks in IID and non-IID data cases. It outperforms the evaluated alternatives, especially in non-IID data and organized attack scenarios. Moreover, it does not cause significant performance loss under no-attack cases where the evaluated alternative defense methods cause some performance loss, especially in non-IID cases. Hence, ARFED can be used to defend against various kinds of poisoning attacks without

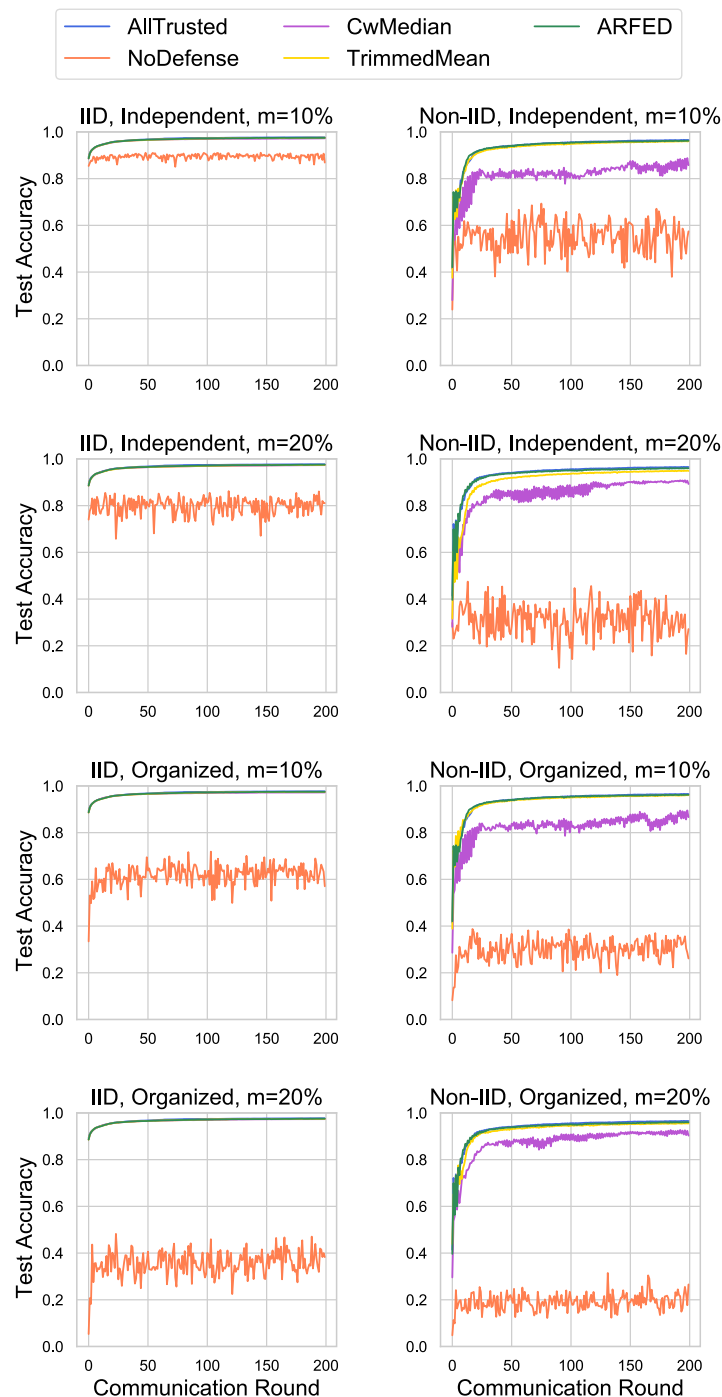


Fig. 4. Accuracy curves of different strategies for MNIST-2NN under Byzantine attacks at different attacker ratios.

worrying about significant performance degradation under no-attack cases.

The rest of the paper is organized as follows. The related work in the literature is reviewed in Section 2. Section 3 describes the proposed outlier-detection-based malicious participant identification technique and delineates the ARFED algorithm. The details of the experimental design used to evaluate the performance and compare it with the other approaches are presented in Section 4. The experimental results are presented and discussed in Section 5. Section 6 gives concluding remarks and directions

for future work. Finally, the supplementary experimental results are given in [Appendix](#).

2. Related work

Security is a critical issue in FL as it is vulnerable to several attacks, such as poisoning, backdoor, free-riding, inference, and eavesdropping [17]. This paper focuses on poisoning attacks in FL and proposes a defense mechanism called ARFED. There are two kinds of poisoning attacks: data poisoning [18] and model

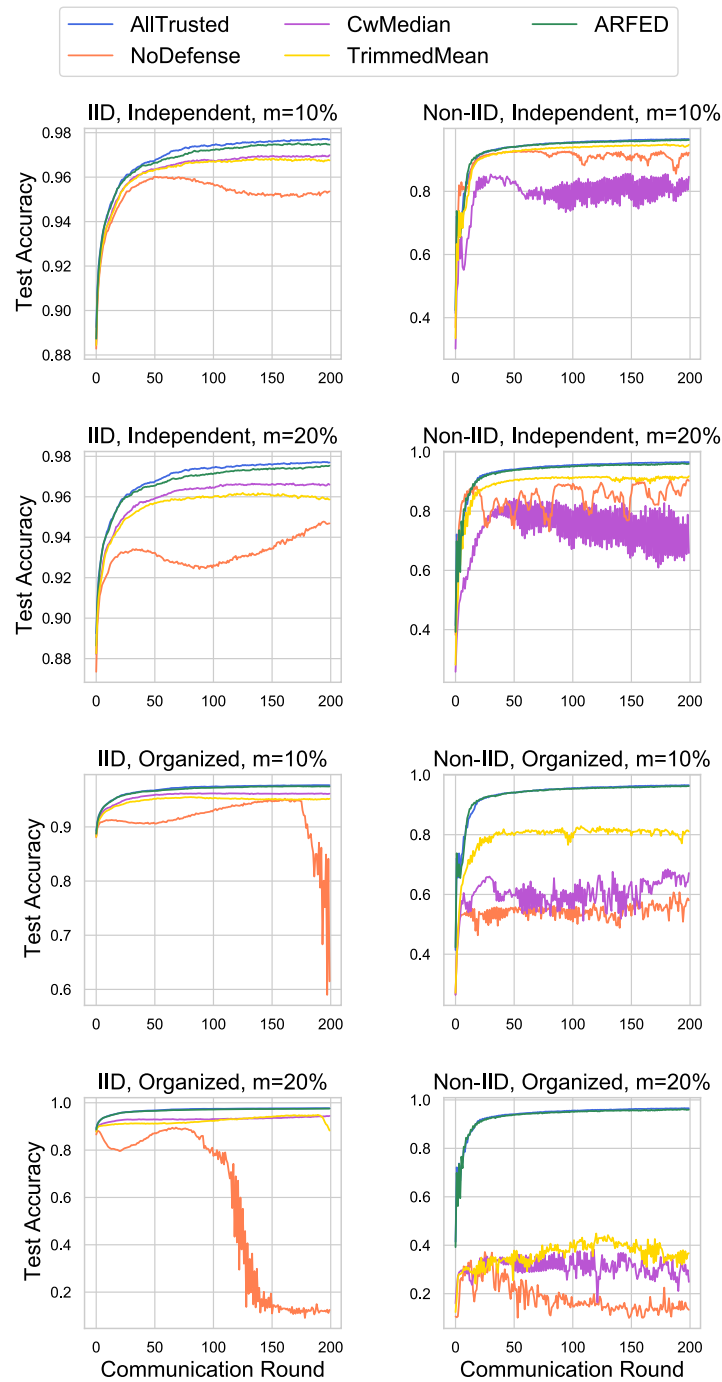


Fig. 5. Accuracy curves of different strategies for MNIST-2NN under partial knowledge attacks at different attacker ratios.

poisoning [19]. In data poisoning attacks, malicious participants manipulate or modify data, for instance, by adding noise to the training data or label flipping [23–26]. In model poisoning attacks, participants alter the models sent to the server in each iteration. Byzantine attack in which malicious participants send arbitrary updates is one of the prevalent model poisoning attacks [19, 23, 24, 27, 28]. Backdoor attacks aim to affect the global model adversarially on a particular sub-task, for example, by making the model classify “trucks” as “planes” by adding small visual artifacts to the training set [24, 25, 29–31].

In FL, the central server does not have access to the participants’ training data or control over the participants’ training process. Therefore, aggregation carried out by the server is the most appropriate step to defend against such security attacks.

Hence, many defense algorithms incorporated into aggregation rules are proposed to handle these attacks and prevent performance degradation caused by them [27, 29, 32–35]. Apart from the defense mechanisms incorporated into the aggregation process carried out by the central server, there are also decentralized solutions. With the distributed ledger technology provided by the blockchain, the need for a trusted central server can be removed [36]. Using blockchain technologies in FL can also provide robustness against adversarial attacks [37]. Although this is one of the promising approaches to defend against attacks, it is out of the scope of this paper.

Many studies have shown that the defense mechanisms proposed in the literature usually make assumptions that do not hold for practical FL settings [19, 30, 38, 39]. In particular, Non-IID

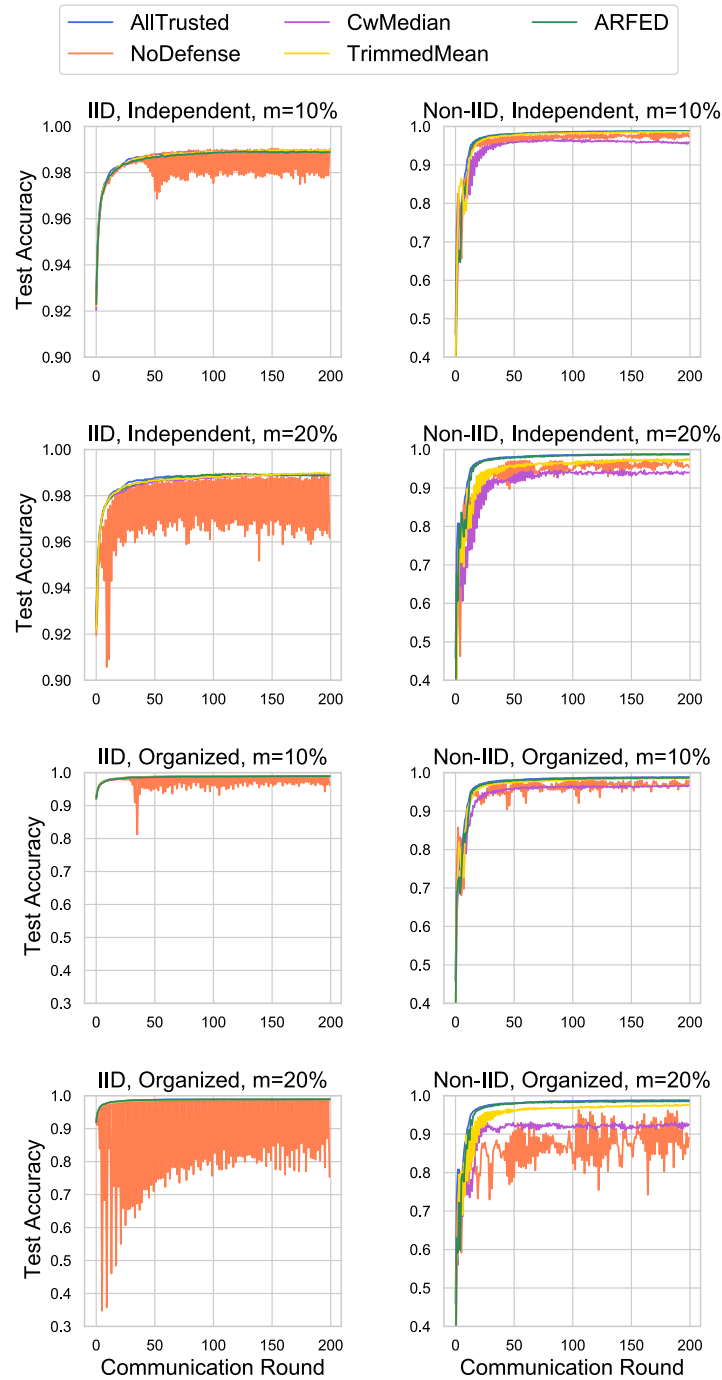


Fig. 6. Accuracy curves of different strategies for MNIST CNN under label flipping attacks at different attacker ratios.

datasets and organized (coordinated) attacks bring severe issues to the learning problem and invalidate assumptions of previous works. Moreover, defense strategies that require examining the local datasets and utilizing partial or complete knowledge of the training process (defense against backdoor attacks and approaches using data sanitization) are not appropriate in practical FL settings. Thus, analyzing and developing these approaches in realistic FL environments is an important study area.

Attack-robust FL has been a heavily studied topic in recent years. Yin et al. [35] introduced two different approaches instead of solely averaging gradients. The first method was the

coordinate-wise median and the second was the coordinate-wise trimmed mean that excludes the highest and smallest values with the given percentage. Blanchard et al. [27] proposed a Byzantine fault-tolerant SGD algorithm called Krum that combines the majority-based and square-distance methods. El Mhamdi et al. [28] introduced a method that combines Krum and trimmed mean, called Bulyan. These methods presume IID data and the ratio of malicious participants should be known in each communication round, which usually does not hold in FL settings.

There are also clustering or similarity metrics based methods that work under certain conditions and with certain assumptions.

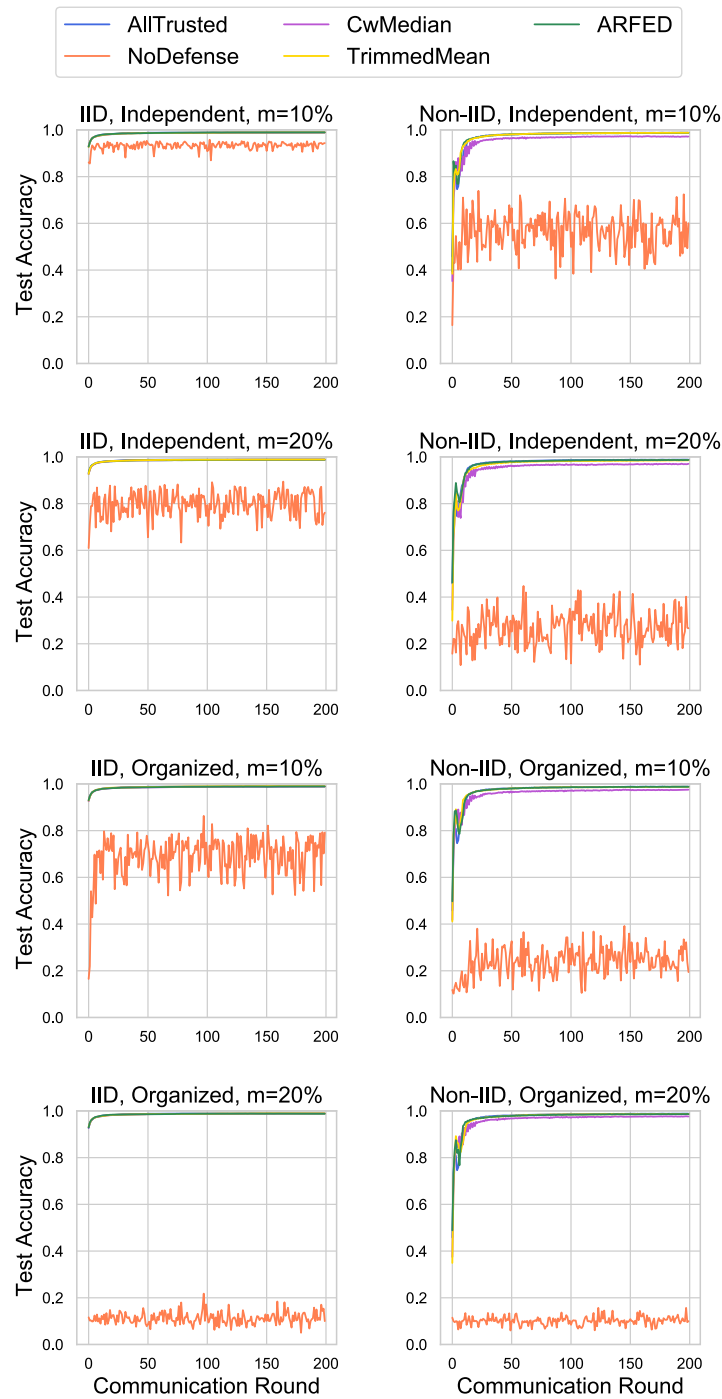


Fig. 7. Accuracy curves of different strategies for MNIST CNN under Byzantine attacks at different attacker ratios.

Fung et al. [25] assume that the trusted participants have a unique distribution and as a result, their gradient updates vary. Since malicious participants have a common goal, their gradient updates tend to be more similar. Based on this assumption, Fung et al. proposed the FoolsGold algorithm that identifies participants who make similar gradient updates with a method based on cosine similarity and reduces the learning rates of these participants. Sattler et al. [24] proposed a method that clusters the participants based on the pairwise cosine dissimilarities between their updates and considers the elements of the largest cluster as benign. Tolpegin et al. [18] presented a method based on identifying malicious participant clusters with a visualization that is obtained by

applying Principal Component Analysis to the parameters of the last layer of the participants' local models. Unlike Fung et al. [25], Sattler et al. [24] and Tolpegin et al. [18] worked with benign participants that have similar updates on IID data.

The data distribution and update similarities of participants are two essential factors that should be examined in detail. Most of the recent studies proposed methods for IID case [18,24,27,28,35]; however, Non-IID distribution of participants' data is one of the key properties of FL and it was emphasized that existing techniques for Byzantine tolerant distributed learning do not perform well when data of participants are Non-IID [9,29,30]. Although the proposed method in [25] addresses Non-IID data

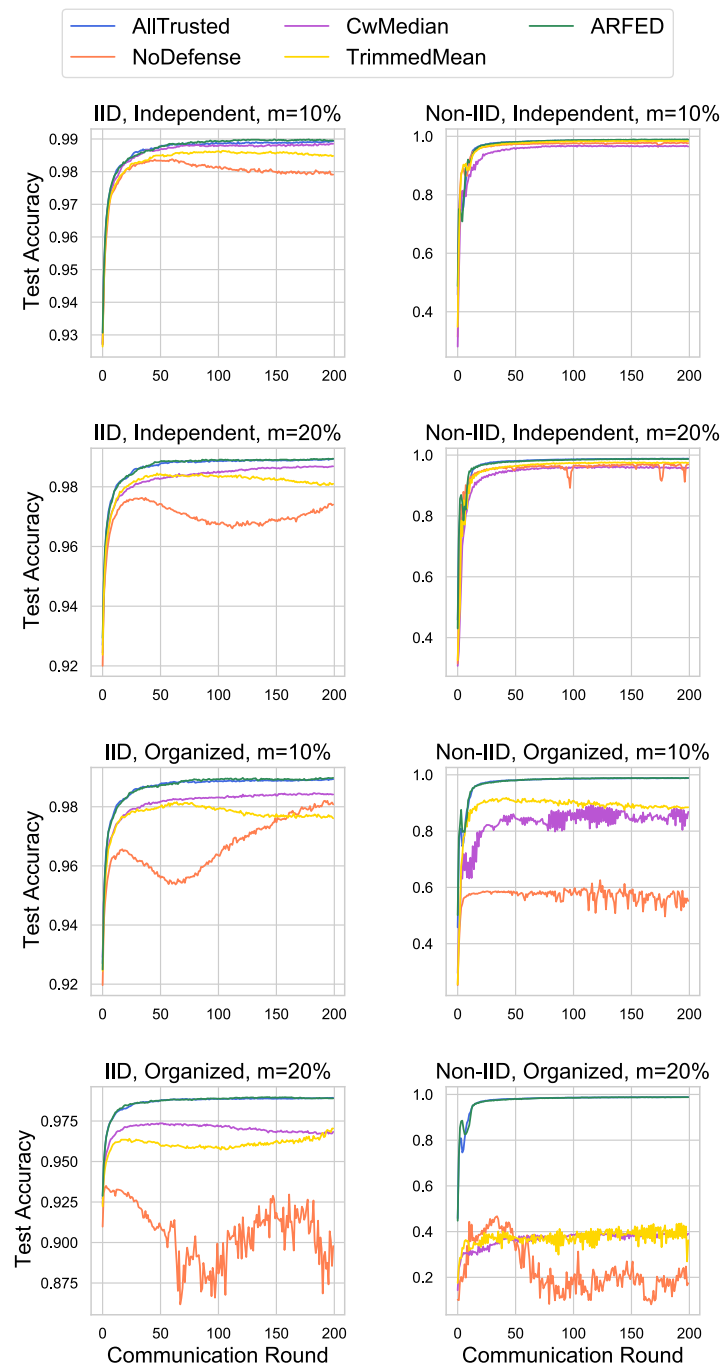


Fig. 8. Accuracy curves of different strategies for MNIST CNN under adaptive partial knowledge attacks at different attacker ratios.

distribution, it only covers a very specific case where trusted participants have unique updates and malicious participants have similar updates.

Although there have been notable new studies proposing aggregation methods for distributed learning that ensure the convergence of the global model, they sacrificed classification performance in exchange for convergence, resulting in ineffective strategies that are not useful for FL settings [23,27,28,35].

3. Attack-Resistant Federated Learning

The Attack-Resistant Federated Learning (ARFED) algorithm is based on the Federated Averaging (FedAvg) algorithm [9], which is one of the most widely used aggregation algorithms in FL [40].

In FedAvg, a server initializes a global parametric model such as a multi-layer neural network which the participants collectively train in several rounds. In each round, the server sends the current global model to the participants, which apply the gradient descent algorithm to train the current global model using their locally available data for several epochs and then send the resulting local models back to the server. The server aggregates the participants' locally trained models by calculating the weighted averages of corresponding parameters in the local models and updates the global model accordingly. The weights are determined according to the number of training examples in the participants.

When there is no attack and the global model converges, the latest local models received from the participants are unlikely to

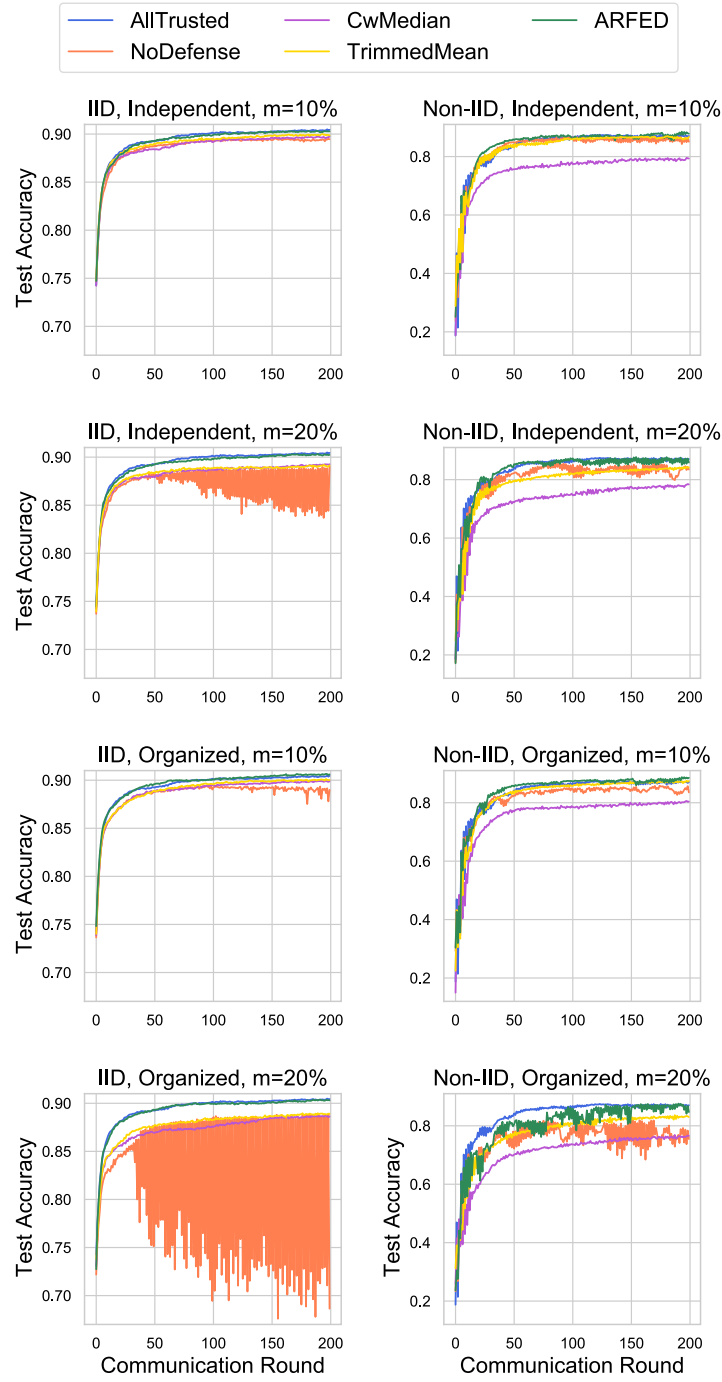


Fig. 9. Accuracy curves of different strategies for Fashion-MNIST under label flipping attacks at different attacker ratios.

be far from the global model. In contrast, if there are malicious participants, their models drift apart from the trusted ones and the global model [23]. Hence, outlier identification can be a promising approach to defend against model poisoning and data poisoning attacks by identifying and filtering out potentially malicious participants. To illustrate the situation and gain insight into the malicious participant detection problem, we carried out an experiment using the experimental setup introduced in Section 4. In this experiment, a set of trusted and malicious participants collaborate to train a three-layer neural network on the MNIST dataset. In order to quantify how far the local models are from the global model, we considered the differences between the parameters of the global model at the beginning of a round and the corresponding parameters of local models at the end of

that round. We used Euclidean distances between the global and the local model parameters to facilitate comparisons and outlier detection. Furthermore, to improve granularity, we calculated distances for each layer separately. To this end, each layer having K parameters (i.e., weights and biases) is represented by a point in K -dimensional space, and the Euclidean distances between the points corresponding to the global model's layer and the participants' models are computed for each layer. We considered the distribution of distances for each participant for each layer and used Inter Quartile Range (IQR) method to identify outliers. IQR is a measure of the spread of data and defined as the 25th percentile (i.e., Q_1 , the first quartile) and 75th percentile (i.e., Q_3 , the third quartile) of the data. That is, $IQR = Q_3 - Q_1$. According to the IQR technique, values less than $Q_1 - 1.5 \times IQR$ (i.e., the

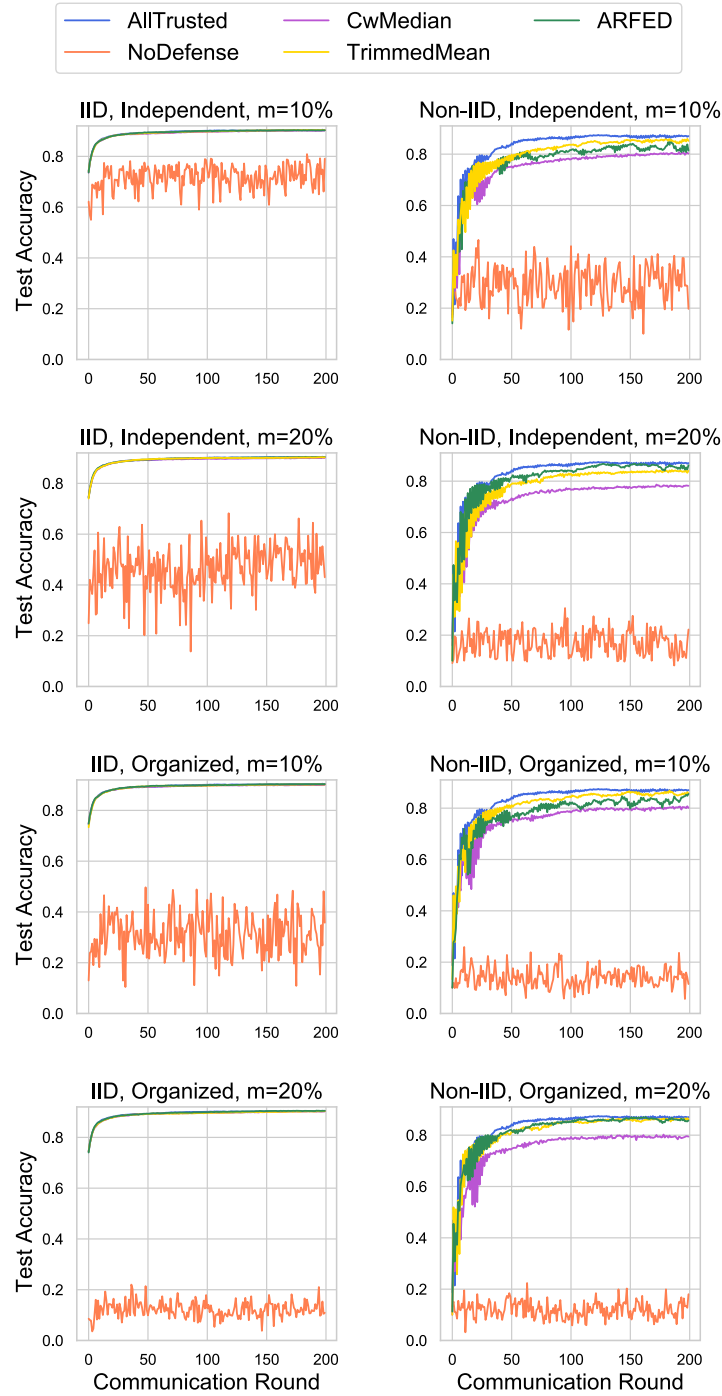


Fig. 10. Accuracy curves of different strategies for Fashion-MNIST under Byzantine attacks at different attacker ratios.

lower limit) or greater than $Q_3 + 1.5 \times IQR$ (i.e., the upper limit) are considered outliers. Fig. 1 shows the layer-wise distances of randomly selected trusted and malicious participants' models distances to the global model in successive federated learning rounds of the experiment. The shaded regions in the plots indicate $[Q_1 - 1.5 \times IQR, Q_3 + 1.5 \times IQR]$ range for non-outlier values through rounds. As seen from the plots, for this illustrative case, the selected trusted participant's distances to the global model almost always remain within the non-outlier range for all layers. On the other hand, the selected malicious participant's distances to the global models are almost always outliers. Inspired by this

experiment, we incorporated the IQR-based outlier identification technique to eliminate model updates from potentially malicious participants in the aggregation step of FedAvg.

The pseudocodes of the procedures carried by an ARFED server and ARFED participants are presented in Algorithm 1. Each participant executes the ParticipantUpdate procedure (Lines 01–05), which is invoked by the server, and the server executes the ServerUpdate procedure (Lines 06–29). The notation used in this algorithm is introduced in Table 1. For the sake of simplicity, we presumed that the algorithm is run for a predefined number of rounds, all participants are involved in all training rounds, and

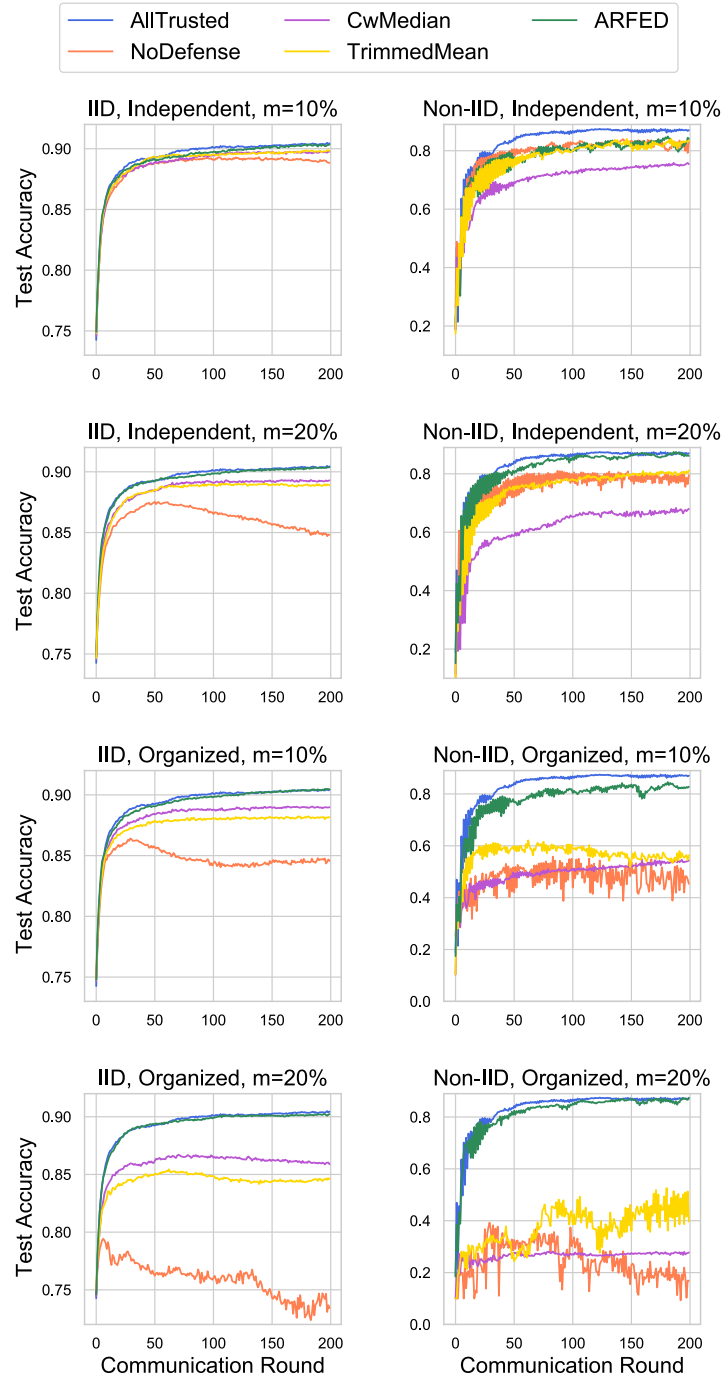


Fig. 11. Accuracy curves of different strategies for Fashion-MNIST under adaptive partial knowledge attacks at different attacker ratios.

participants employ batch gradient descent using a predefined learning rate in their local training process. However, this algorithm can easily be extended to realize other practices, such as repeating the process until convergence criteria are met, involving a subset of participants in each round, employing mini-batch gradient descent in local training, and learning rate scheduling.

ARFED is essentially an extension to FedAvg to train a parametric machine learning model such as a multi-layer neural network using the Gradient Descent algorithm. In ARFED, a server trains a randomly initialized global model (Line 07) in T rounds by collaborating with P participants (Lines 08–26). In each round t , participants receive the previous round's global model represented by parameters \mathbf{W}^{t-1} from the server (Lines 09–11). Then they apply the gradient descent algorithm to improve the model

using their locally available training data for E epochs according to a server-defined loss function $\mathcal{L}(x, y, \mathbf{W})$ (Lines 02–04). They finally return the resulting local model to the server (Line 05). The participant update in ARFED is essentially the same as the participant update of the vanilla FedAvg. Like FedAvg, the server initializes a global model (Line 07) and trains the model for T rounds (Lines 08–28) by involving participants and returns the resulting model (Line 29). The main difference between FedAvg and ARFED is in the aggregation process. In FedAvg, the local models (\mathbf{W}_p^t) received from the participants are aggregated to update the global model by computing weighted averages for each parameter (like Line 27, but including all participants). However, in ARFED, only the models received from participants deemed reliable (i.e., the participants in the reliable participant set r^t)

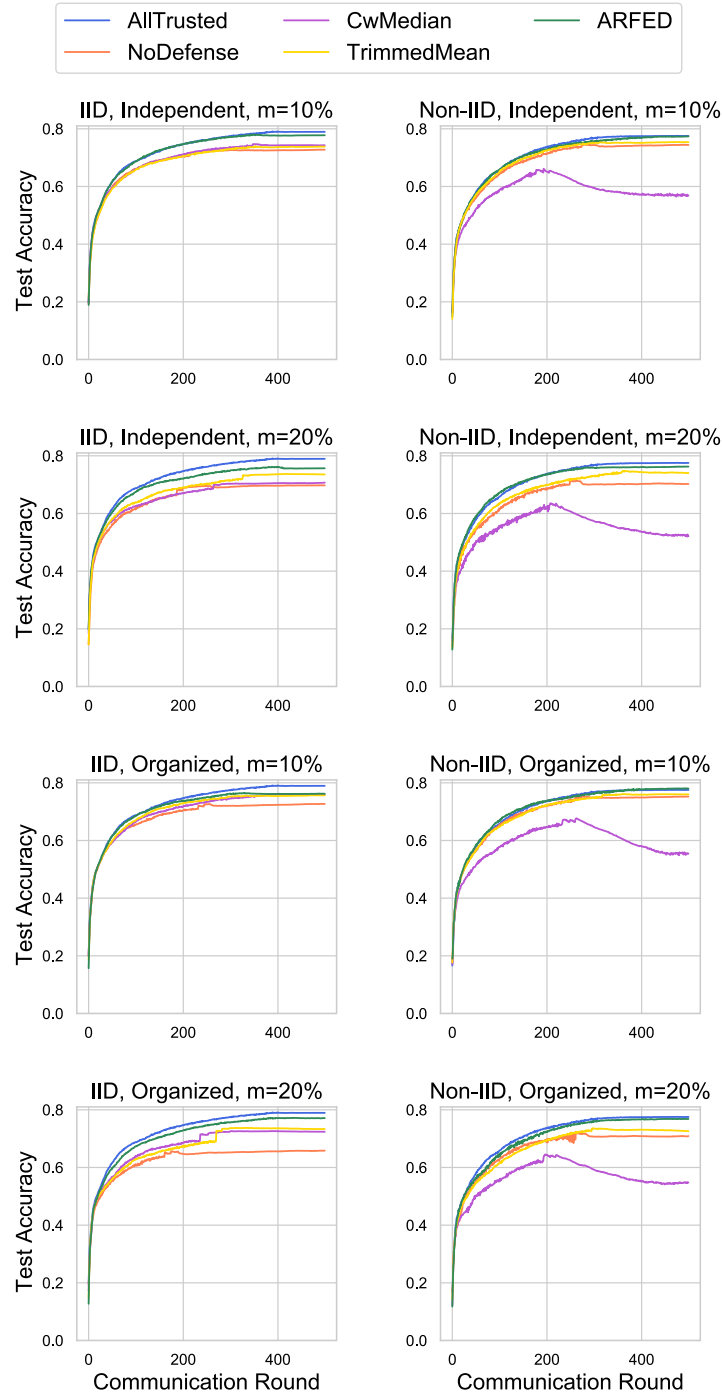


Fig. 12. Accuracy curves of different strategies for CIFAR10 under label flipping attacks at different attacker ratios.

in that round are included in the model aggregation step (Line 27). In the beginning, all participants are assumed to be reliable (Line 12), so r^t includes all participants. In order to identify the potentially malicious participants, for each layer l , a list of distances d_l^t between the received participant models and the global model is computed (Lines 15–18). Accordingly, a lower distance threshold $\min_d_l^t$ and an upper distance threshold $\max_d_l^t$ are determined by computing Q_1 , Q_2 , and IQR of distance values in d_l^t (Lines 19–20). Then, if a participant p is an outlier in a layer l (i.e., its distance $d_{p,l}^t$ is less than the lower distance threshold $\min_d_l^t$ or greater than the upper distance threshold $\max_d_l^t$ it is

considered malicious and removed from the reliable participant set r^t (Lines 21–25). Hence, if a participant is identified as potentially malicious according to at least one layer, it is identified as not reliable. After evaluating all layers, the federated averaging is applied to the local models of the participants deemed reliable (i.e., the participants in r^t) (Line 27).

Defense strategies such as trimmed mean and coordinate-wise median rely on including a model's parameters partially. Each parameter within the model is evaluated individually; some participants' updates are included, and the rest are discarded in the aggregation step for each parameter. Hence, a different group

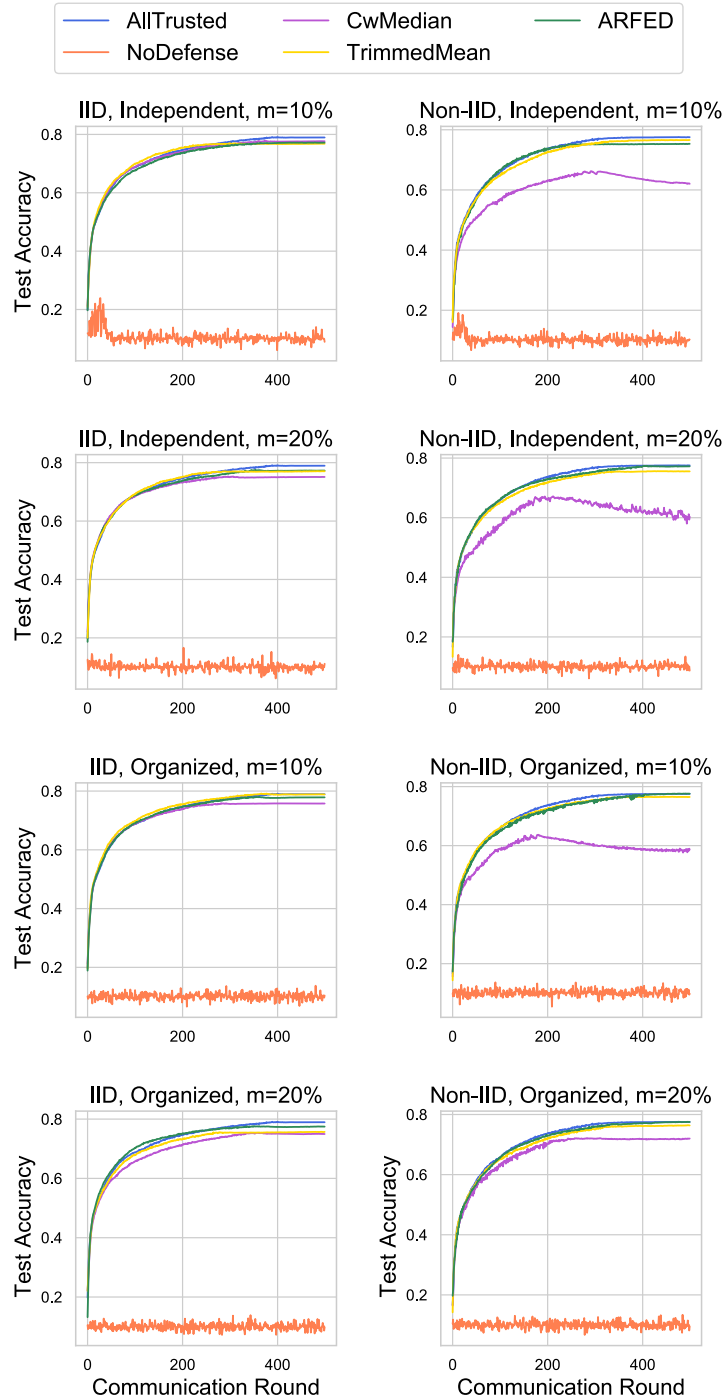


Fig. 13. Accuracy curves of different strategies for CIFAR10 under Byzantine attacks at different attacker ratios.

of participants can potentially contribute to each parameter. The primary motivation of the proposed all-or-nothing approach is that each participant is evaluated in a holistic approach. Parameters of a neural network are highly dependent on each other; therefore, independently evaluating each parameter may lead to misleading inferences. If any layer of a model update is an outlier, it is a sign of a malicious participant; therefore, it is not reasonable to include that participant in the aggregation step. In the proposed approach, for a participant's model to be included in the calculation of global model aggregation, each layer must fall within the safe interval calculated for that layer, i.e., a consensus should be ensured among all layers of the local model. Interestingly, experiments show that the ratio of unreliable participants

determined in the proposed approach is very close to actual malicious participant ratios (see Figs. 16 and 15 in Appendix).

Lemma 1. For an L -layer neural network with K parameters (weights and biases) in each layer collectively trained by P participants, the time complexity of reliable participant identification in ARFED is $\mathcal{O}(L \cdot P \cdot (K + \log P))$.

Proof. The server computes the differences and ℓ^2 -norms of K -dimensional vectors for each layer with K parameters, $\mathcal{O}(K)$. As there are L layers in the models received from P participants, the distance calculation is $\mathcal{O}(L \cdot P \cdot K)$. For each layer, the server sorts the distances of P participants to find Q_1 , Q_3 , and IQR, determining

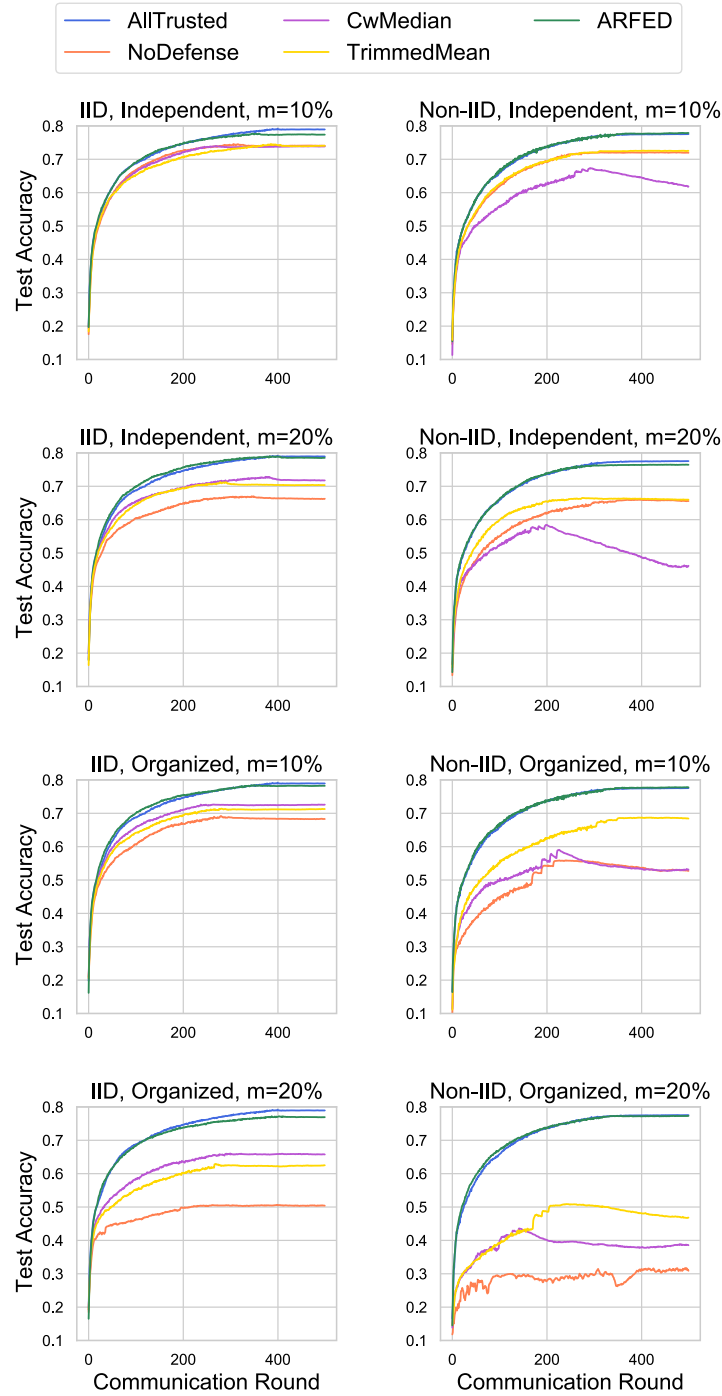


Fig. 14. Accuracy curves of different strategies for CIFAR10 under adaptive partial knowledge attacks at different attacker ratios.

lower and upper thresholds, which is $\mathcal{O}(P \cdot \log P)$. As there are L layers, the threshold determination is $\mathcal{O}(L \cdot P \cdot \log P)$. Finally, for each layer, the distances of P participants are checked if they are outliers or not, which is $\mathcal{O}(P)$. As there are L layers, the outlier detection is $\mathcal{O}(L \cdot P)$. As a result, the algorithm's time complexity can be found as $\mathcal{O}(L \cdot P \cdot K + L \cdot P \cdot \log P + L \cdot P) \rightarrow \mathcal{O}(L \cdot P \cdot (K + \log P + 1)) \rightarrow \mathcal{O}(L \cdot P \cdot (K + \log P))$.

The computational complexity is an essential factor in the practical use of an algorithm. Lemma 1 states that the time complexity of reliable participant identification (i.e., the extension made to FedAvg) in ARFED is $\mathcal{O}(L \cdot P \cdot (K + \log P))$. Hence it is much more efficient than that of Krum and its variant Bulyan, which are quadratic, $\mathcal{O}(K \cdot P^2)$. ARFED is slightly more efficient than

the coordinate-wise median and trimmed mean as they require sorting for all individual parameters (ARFED only makes sorting as many as the number of layers). \square

4. Experimental design

The dimensions of the experimental designs are the datasets, the data distribution of the participants, attack types, attacker types, and baseline method selection.

4.1. Datasets

We conducted experiments on MNIST [41], CIFAR10 [42], and Fashion-MNIST [43] datasets which are widely used by researchers

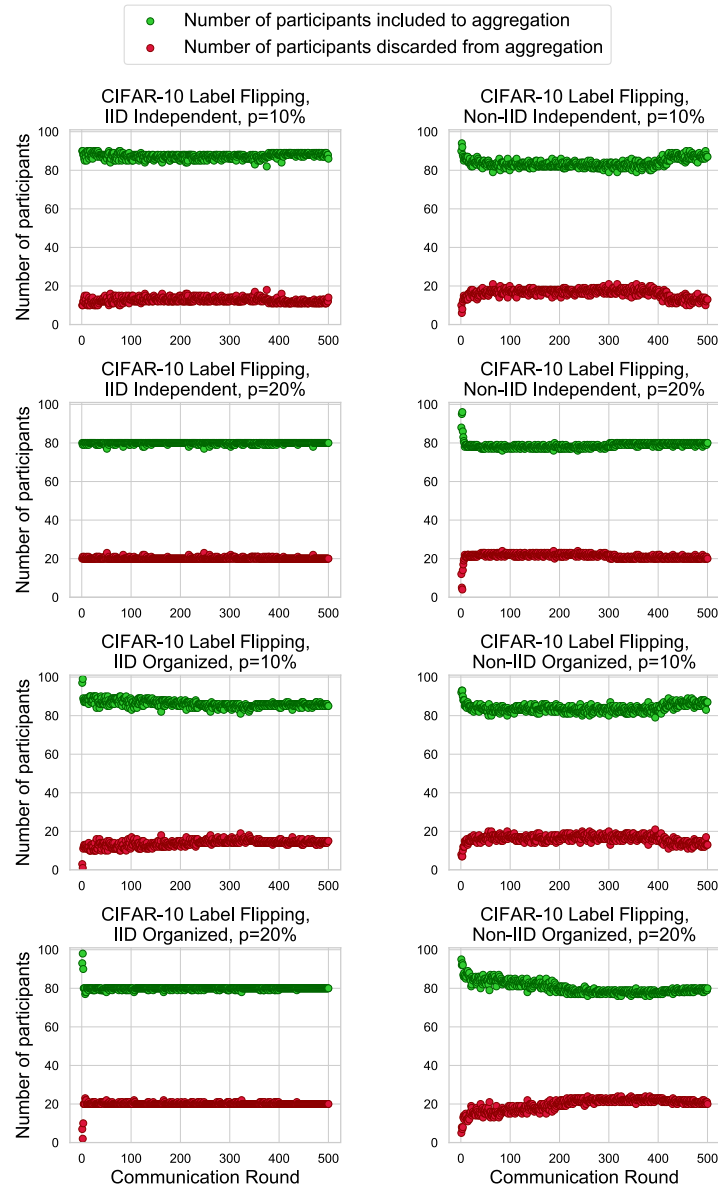


Fig. 15. Number of participants marked as reliable and outlier in CIFAR10 label flipping attacks.

to evaluate FL approaches [18,23,24,27–29,35]. MNIST dataset contains 28×28 grayscale images with 50,000 training images and 10,000 testing images. CIFAR10 dataset contains 32×32 color images with 50,000 training images and 10,000 testing images, and Fashion-MNIST contains 28×28 grayscale images with 60,000 training images and 10,000 testing images.

There are 100 participants in all experiments. There are two model architectures for MNIST, and only one architecture for CIFAR10 and Fashion-MNIST. The details of the model architectures, and the hyperparameters are given in Appendix A.1. For CIFAR10 experiments, data augmentation techniques such as horizontal flipping and random cropping, and training strategies like learning rate scheduling and gradient clipping were applied to enhance the model performance.

4.2. Data distributions

One of the dimensions of our experiments is the data distribution of the participants, which can be either IID or Non-IID.

For IID cases, training datasets are distributed to the participants randomly and uniformly, i.e., each participant has each class equally. On the other hand, in the Non-IID case, each participant has examples of only two randomly selected classes for MNIST and Fashion-MNIST and examples of only five randomly selected classes for CIFAR10.

4.3. Attack and attacker types

Another dimension is attack types. Three attack types are examined: label flipping attacks, Byzantine attacks, and adaptive partial knowledge attacks. In label flipping attacks, malicious nodes flip their ground truth labels with a target class label. In Byzantine attacks, malicious participants send random weight updates from a standard normal distribution with zero mean and unit standard deviation. In adaptive attacks, malicious participants use statistics of local models' parameter to manipulate sending weights.

Lastly, attacker types are investigated. Independent attackers are malicious participants incapable of coordinating with each

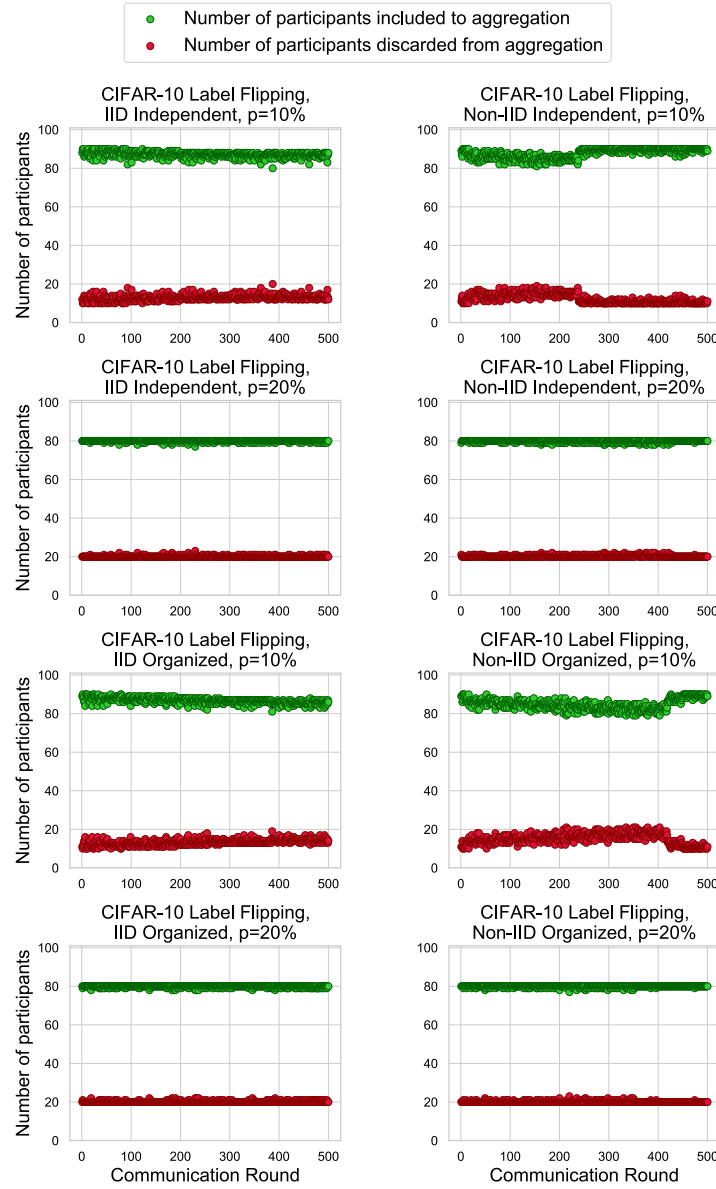


Fig. 16. Number of participants marked as reliable and outlier in CIFAR10 Byzantine attacks.

other, acting individually, and sending random updates to the server. Organized or coordinated attackers are malicious participants that carry out the attack in an organized or coordinated manner and send similar updates to the server. For example, in independent label flipping attacks, malicious participants flip their ground truth labels with an arbitrary target label, e.g., if there are two malicious participants with label 7 in their data sets, one flips 7 to 1, while the other flips to 4. On the other hand, in organized label flipping attacks, the malicious participants flip ground truth labels with consistent target labels, e.g., all malicious participants that have 7 in their datasets flip the label as 1.

In order to increase the success of the attacks and reduce the likelihood of malicious participants being caught, the replaced classes were chosen as semantically similar as possible. The replaced classes in the organized setting for each data set are presented in Table 2.

Similarly, malicious participants send different random weights for independent Byzantine attacks while they send the same random weights in the organized setting. The details of independent and organized adaptive partial knowledge attacks are given in Section 4.4

4.4. Adaptive partial knowledge attack

In the original partial knowledge attack in [19], the malicious participants train their local models with their local data. Then, for each parameter, mean, μ_w , and standard deviation, σ_w , are estimated among the malicious participants' parameters. Later, each malicious participant determines the update changing direction, s_w , for each parameter by looking at the global model they have received at the beginning of the FL round (if $w_m^{t+1} \geq w^t \rightarrow s_w = 1$, else $s_w = -1$).

If $s_w = -1$, each malicious participant replaces the parameter with a number uniformly sampled from the interval $[\mu_w + 3\sigma_w, \mu_w + 4\sigma_w]$. If $s_w = 1$, the malicious participant replaces the parameter with a number uniformly sampled from the interval $[\mu_w - 4\sigma_w, \mu_w - 3\sigma_w]$. We show the results of the original version of this attack under the independent experimental setting (See Tables 8 and 9).

In our experimental setting, the malicious participants send the same parameters to the server in the "Organized Byzantine" attacks. Based on this idea, we adopted the attack in [19] for the "Organized" version. For this time, to decide the direction

Algorithm 1 ARFED.

```

1: procedure PARTICIPANTUPDATE( $p, \mathbf{W}$ )
2:   for  $e = 1, 2, \dots, E$  do
3:      $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{1}{n_p} \sum_{i=1}^{n_p} \nabla_{\mathbf{W}} \mathcal{L}(x_{p,i}, y_{p,i}, \mathbf{W})$ 
4:   end for
5:   return  $\mathbf{W}$ 
6: procedure SERVERUPDATE
7:    $\mathbf{W}^0 \leftarrow$  initial weights(random)
8:   for  $t = 1, 2, \dots, T$  do
9:     for  $p = 1, 2, \dots, P$  do
10:       $\mathbf{W}_p^t \leftarrow$  PARTICIPANTUPDATE( $p, \mathbf{W}^{t-1}$ )
11:    end for
12:     $r^t \leftarrow \{p | p = 1, \dots, P\}$ 
13:    for  $l = 1, 2, \dots, L$  do
14:       $d_l^t = \{\}$ 
15:      for  $p = 1, 2, \dots, P$  do
16:         $d_{p,l}^t \leftarrow \|\text{flatten}(\mathbf{W}^{t-1}[l]) - \mathbf{W}_p^t[l]\|_2$ 
17:         $d_l^t \leftarrow d_l^t + [d_{p,l}^t]$   $\triangleright$  append to the list
18:      end for
19:       $\min\_d_l^t = Q_1(d_l^t) - 1.5 \times \text{IQR}(d_l^t)$ 
20:       $\max\_d_l^t = Q_3(d_l^t) + 1.5 \times \text{IQR}(d_l^t)$ 
21:      for each  $p$  in  $r^t$  do
22:        if  $d_{p,l}^t < \min\_d_l^t$  or  $d_{p,l}^t > \max\_d_l^t$  then
23:           $r^t \leftarrow r^t - \{p\}$   $\triangleright$  remove from the set
24:        end if
25:      end for
26:    end for
27:     $\mathbf{W}^t \leftarrow \frac{1}{\sum_{p \in r^t} n_p} \sum_{p \in r^t} n_p \times \mathbf{W}_p^t$ 
28:  end for
29:  return  $\mathbf{W}^t$ 

```

of change (s_w), the parameter of the global model is compared with the mean parameter, μ_w , for once instead of comparing separately for each participant's parameter. Then, the same (s_w) is used for each participant. If $s_w = -1$, the malicious participants replace the parameters with the same number sampled from the interval $[\mu_w + 3\sigma_w, \mu_w + 4\sigma_w]$. If $s_w = 1$, the malicious participants replace the parameter with the same number sampled from the interval $[\mu_w - 4\sigma_w, \mu_w - 3\sigma_w]$.

4.5. Baseline method selection

Fang et al. [19] have shown that trimmed mean and coordinate-wise median are more robust to attacks than Krum and its variant Bulyan. Unlike the coordinate-wise median, trimmed mean requires the knowledge of malicious participant ratio, which does not meet the nature of a realistic FL setting. Yet, both coordinate-wise median (will be referred to as CwMedian in the rest of the article) and trimmed mean (will be referred to as TrimmedMean in the rest of the article) are agnostic to update similarity of participants unlike [18,24,25], and they are outlier based methods as the proposed method ARFED; therefore, they are chosen as the baseline methods in our performance evaluations.

5. Experimental results & discussion

Attacks usually compromise convergence as well as the performance of the models and cause oscillations in test set accuracies; therefore, reporting only the score of the last communication round can be misleading because the peak point or lowest point of this oscillation may occur randomly. Thus, the minimum and the maximum accuracies achieved on the test set in the last ten

Table 1

Notation table.

| | |
|---------------------------------|--|
| P | Number of participants |
| p | Participant p , $p \in \{1, \dots, P\}$ |
| n_p | Number of training examples in participant p |
| $x_{p,i}$ | Features of i th training example in participant p |
| $y_{p,i}$ | Label of i th training example in participant p |
| T | Number of training rounds |
| t | Federated learning round t , $t \in \{1, \dots, T\}$ |
| L | Number of layers in the global model |
| l | Layer l in the global model, $l \in \{1, \dots, L\}$ |
| K_l | Number of parameters (weights and biases) in layer l |
| \mathbf{W}^t | Weights of the global model at round t |
| $\mathbf{W}^t[l]$ | Weights of the l th layer of the global model at round t |
| \mathbf{W}_p^t | Weights of the local model in participant p at round t |
| $\mathbf{W}_p^t[l]$ | Weights of the l th layer of the local model in participant p at round t |
| $\mathcal{L}(x, y, W)$ | Loss function for a training example x with label y on a model with weights W |
| $\nabla_W \mathcal{L}(x, y, W)$ | Gradient of the loss function with respect to weights W |
| E | Number of local iterations (epochs) in each round |
| e | Local training iteration e , $e \in \{1, \dots, E\}$ |
| η | Learning rate |
| $w = \text{flatten}(W)$ | flatten a tensor W to a vector w |
| $\ w\ _2$ | ℓ^2 -norm of a vector |
| $Q_q(d)$ | q th quartile of values in list d for $q \in \{1, 2, 3, 4\}$ |
| $\text{IQR}(d)$ | Inter Quartile Range of values in list d ; $\text{IQR}(d) = Q_3(d) - Q_1(d)$ |
| r^t | Set of participants marked as reliable at round t |
| $d_{p,l}^t$ | Distance between the global model's l th layer and participant p 's l th layer in round t |
| d_l^t | List of distances between the global model's l th layer and all participants' l th layers in round t |
| $\min_d_l^t$ | Lower distance threshold for layer l in round t to mark a participant reliable |
| $\max_d_l^t$ | Upper distance threshold for layer l in round t to mark a participant reliable |

Table 2

Replaced classes for organized label flipping attack.

| MNIST | | Fashion-MNIST | | CIFAR10 | |
|----------|----------|---------------|-------------|----------|----------|
| Original | Replaced | Original | Replaced | Original | Replaced |
| 0 | 9 | T-shirt/Top | Shirt | Plane | Bird |
| 1 | 7 | Trouser | Dress | Car | Truck |
| 2 | 5 | Pullover | Coat | Bird | Plane |
| 3 | 8 | Dress | Trouser | Cat | Dog |
| 4 | 6 | Coat | Pullover | Deer | Horse |
| 5 | 2 | Sandal | Sneaker | Dog | Cat |
| 6 | 4 | Shirt | T-shirt/Top | Frog | Ship |
| 7 | 1 | Sneaker | Ankle Boot | Horse | Deer |
| 8 | 3 | Bag | Sandal | Ship | Frog |
| 9 | 0 | Ankle Boot | Sneaker | Truck | Car |

rounds are reported in tables to show the severity of the oscillations created by attacks. In all experiments, **NoDefense** refers to the vanilla FedAvg, **CwMedian** refers to the coordinate-wise median and **TrimmedMean** refers to the trimmed mean.

The tables show the results of all datasets, while the figures show only the MNIST-2NN experiments in this section. The corresponding figures of other datasets and architectures can be found in the [Appendix A.3](#). Moreover, additional information about box plot factor comparison is presented in [Appendix A.4](#).

5.1. No attack

A robust defense strategy should not cause any noticeable performance loss when there is no attack in the FL system. [Fig. 2](#) and [Table 3](#) show the results of experiments when all participants are trusted (when there is no attack on any participant). Incorporating TrimmedMean and CwMedian strategies into FedAvg does not cause any performance degradation in the IID setting. Although the performance of the CwMedian strategy is slightly

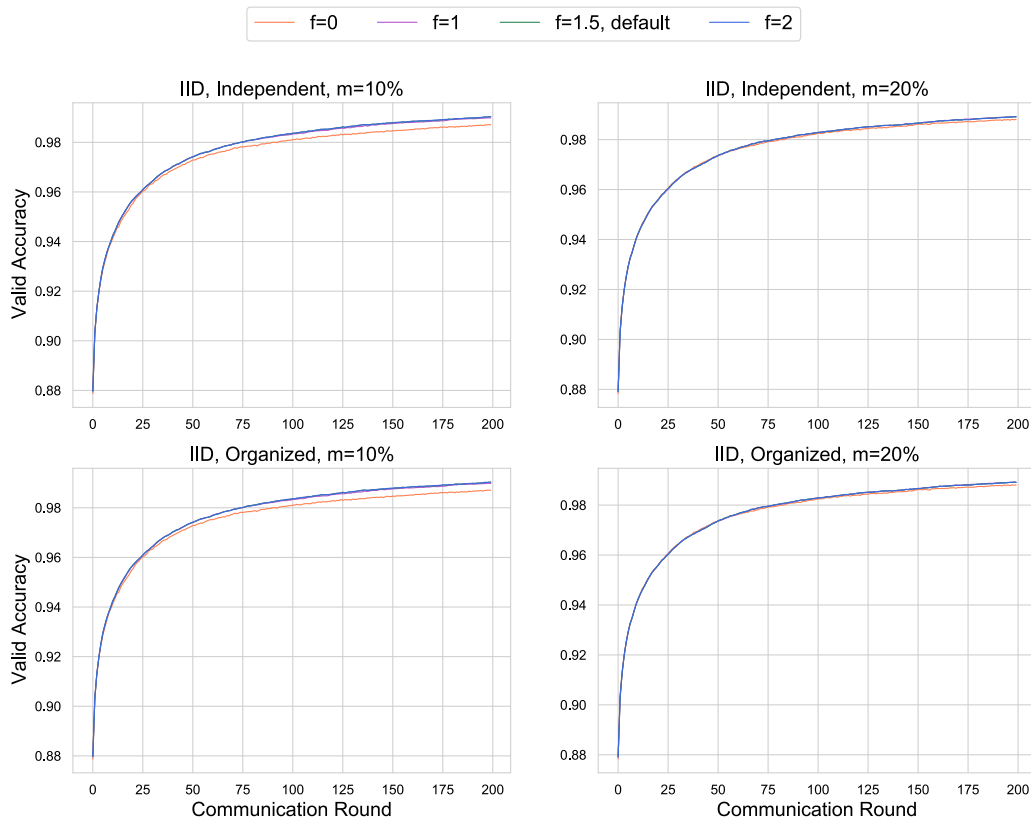


Fig. 17. Box plot factor comparison on accuracy for IID cases.

Table 3

Accuracies on test sets when all participants are trusted (i.e., $m = 0\%$). The worst results are bold.

| | MNIST-2NN | | MNIST-CNN | | Fashion-MNIST | | CIFAR10 | |
|--------------------------|-------------|-------------|-------------|-------------|---------------|-------------|-------------|-------------|
| IID | min | max | min | max | min | max | min | max |
| NoDefense | 97.7 | 97.7 | 98.9 | 98.9 | 90.4 | 90.5 | 78.9 | 79.0 |
| TrimmedMean ^a | 97.7 | 97.7 | 98.9 | 98.9 | 90.4 | 90.5 | 78.9 | 79.0 |
| CwMedian | 97.6 | 97.6 | 98.9 | 98.9 | 90.1 | 90.2 | 76.7 | 76.8 |
| ARFED | 97.6 | 97.6 | 98.9 | 98.9 | 90.4 | 90.5 | 78.4 | 78.4 |
| Non-IID | min | max | min | max | min | max | min | max |
| NoDefense | 96.4 | 96.6 | 98.8 | 98.8 | 86.8 | 87.4 | 77.5 | 77.6 |
| TrimmedMean ^a | 96.4 | 96.6 | 98.8 | 98.8 | 86.8 | 87.4 | 77.5 | 77.6 |
| CwMedian | 80.7 | 85.3 | 96.9 | 97.1 | 79.1 | 80.0 | 64.3 | 64.6 |
| ARFED | 96.2 | 96.4 | 98.7 | 98.8 | 82.4 | 84.3 | 77.8 | 77.9 |

^aThe same results as NoDefense.

worse in Fashion-MNIST and CIFAR10, it can be tolerable in an FL setting. On the other hand, when local datasets are Non-IID, the CwMedian strategy causes significant performance degradation, which points to the questionability of the method. When the malicious participant ratio is zero, in other words, when there is no attack in the system, no participants are discarded from the aggregation step; therefore, TrimmedMean gives the same result as NoDefense.

5.2. Label flipping attacks

Tables 4 and 5 show that as the ratio of malicious participants increases, vanilla FedAvg cannot avoid performance degradation, which requires that a defense mechanism should be incorporated. As indicated by Fig. 3 when malicious participants are organized, degradation becomes more severe and oscillation increases. The most performance degradation occurs when attacks are organized in Non-IID setting.

Table 4

Accuracies under label flipping attacks at different attacker ratios in IID settings. The best results are bold.

| | | Organized | | | | Independent | | | |
|---------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | | m=10% | | m=20% | | m=10% | | m=20% | |
| | | min | max | min | max | min | max | min | max |
| MNIST 2NN | NoDefense | 92.9 | 97.6 | 72.4 | 97.3 | 97.0 | 97.2 | 91.3 | 97.1 |
| | CwMedian | 97.4 | 97.4 | 97.2 | 97.3 | 97.4 | 97.5 | 97.0 | 97.1 |
| | TrimmedMean | 97.5 | 97.5 | 96.9 | 97.0 | 97.5 | 97.5 | 97.1 | 97.1 |
| | ARFED | 97.6 | 97.7 | 97.4 | 97.5 | 97.6 | 97.6 | 97.4 | 97.5 |
| MNIST CNN | NoDefense | 94.2 | 99.0 | 75.4 | 99.0 | 97.8 | 99.0 | 96.2 | 98.9 |
| | CwMedian | 98.9 | 98.9 | 98.8 | 98.8 | 98.9 | 98.9 | 98.9 | 98.9 |
| | TrimmedMean | 98.9 | 98.9 | 98.8 | 98.9 | 99.0 | 99.0 | 98.9 | 99.0 |
| | ARFED | 98.9 | 98.9 | 98.9 | 98.9 | 98.9 | 98.9 | 98.9 | 98.9 |
| Fashion MNIST | NoDefense | 87.8 | 89.3 | 68.6 | 88.9 | 89.2 | 89.5 | 83.7 | 89.0 |
| | CwMedian | 89.8 | 89.9 | 88.6 | 88.7 | 89.6 | 89.7 | 89.2 | 89.3 |
| | TrimmedMean | 90.0 | 90.1 | 88.8 | 88.9 | 89.9 | 90.0 | 89.0 | 89.2 |
| | ARFED | 90.5 | 90.7 | 90.3 | 90.4 | 90.2 | 90.3 | 90.2 | 90.3 |
| CIFAR10 | NoDefense | 72.7 | 72.7 | 65.8 | 65.9 | 72.8 | 72.8 | 69.7 | 69.8 |
| | CwMedian | 75.6 | 75.7 | 73.3 | 73.4 | 73.8 | 73.8 | 73.5 | 73.6 |
| | TrimmedMean | 75.6 | 75.7 | 73.3 | 73.4 | 73.8 | 73.8 | 73.5 | 73.6 |
| | ARFED | 76.2 | 76.2 | 77.0 | 77.2 | 77.7 | 77.8 | 75.6 | 75.7 |

For IID cases of MNIST-2NN, MNIST-CNN, and Fashion-MNIST, ARFED achieves a slightly higher accuracy score most of the time, but differences between ARFED, CwMedian, and trimmed mean are not significant. When comparing with all-trusted performance, all strategies can tolerate the negative effects of the label flipping attack. For IID cases of CIFAR10 experiments, ARFED achieves noticeably better performance than the others.

When the data of participants is Non-IID, CwMedian strategy performed worse than even the vanilla FedAvg. Although

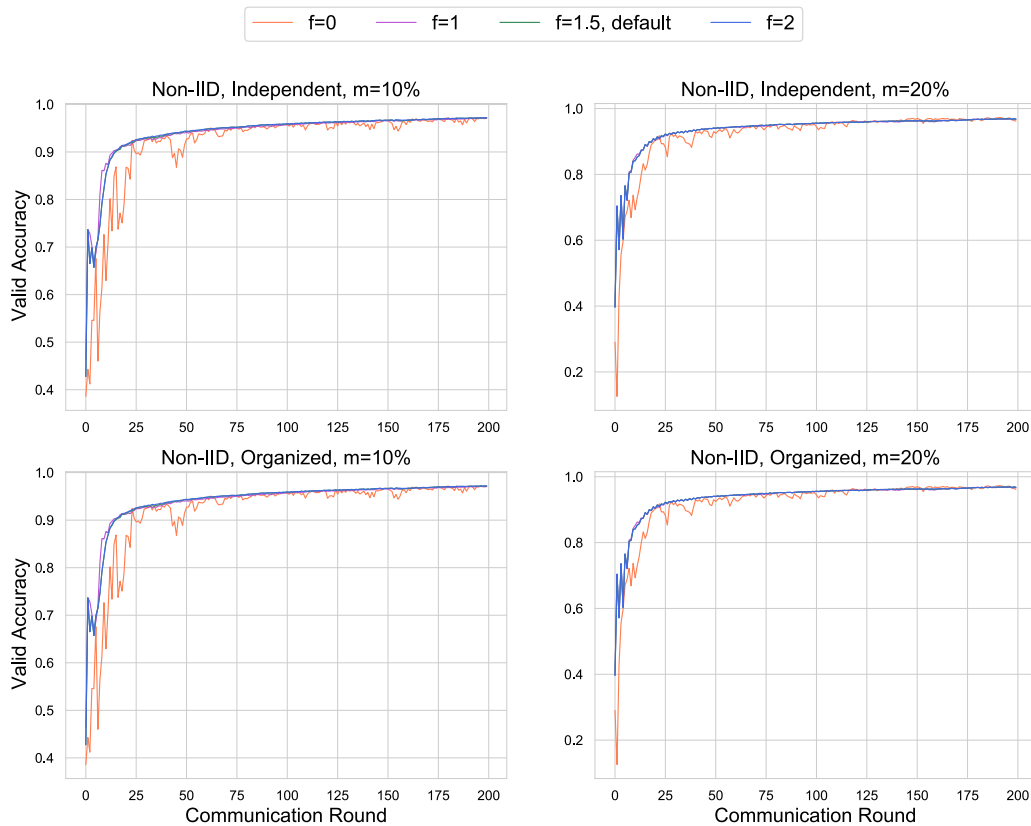


Fig. 18. Box plot factor comparison on accuracy for non-IID cases.

Table 5

Accuracies under label flipping attacks at different attacker ratios in Non-IID settings. The best results are bold.

| | | Organized | | | | Independent | | | |
|------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | | m=10% | | m=20% | | m=10% | | m=20% | |
| | | min | max | min | max | min | max | min | max |
| MNIST 2NN | NoDefense | 92.4 | 95.8 | 83.8 | 88.3 | 93.7 | 95.3 | 89.3 | 94.0 |
| | CwMedian | 75.1 | 83.8 | 67.7 | 75.0 | 80.8 | 83.3 | 67.8 | 75.9 |
| | TrimmedMean | 94.6 | 95.5 | 79.9 | 87.7 | 95.1 | 95.4 | 80.9 | 89.7 |
| | ARFED | 96.1 | 96.3 | 95.6 | 96.1 | 96.0 | 96.3 | 95.5 | 96.1 |
| MNIST CNN | NoDefense | 95.5 | 98.0 | 83.4 | 91.6 | 97.0 | 98.3 | 95.3 | 96.6 |
| | CwMedian | 96.4 | 96.7 | 91.4 | 93.1 | 95.4 | 95.9 | 93.5 | 94.2 |
| | TrimmedMean | 98.5 | 98.6 | 97.3 | 97.6 | 98.6 | 98.6 | 97.2 | 97.6 |
| | ARFED | 98.6 | 98.7 | 98.5 | 98.6 | 98.7 | 98.8 | 98.6 | 98.7 |
| Fashion MNIST | NoDefense | 83.0 | 85.7 | 73.7 | 79.2 | 84.3 | 86.7 | 81.5 | 84.6 |
| | CwMedian | 79.9 | 80.6 | 76.0 | 76.7 | 78.7 | 79.6 | 77.4 | 78.5 |
| | TrimmedMean | 86.7 | 87.5 | 82.7 | 83.4 | 85.8 | 86.7 | 83.8 | 84.3 |
| | ARFED | 87.7 | 88.8 | 84.2 | 87.6 | 87.5 | 88.5 | 85.1 | 87.6 |
| CIFAR10 | NoDefense | 75.2 | 75.3 | 70.8 | 70.9 | 74.4 | 74.4 | 70.2 | 70.3 |
| | CwMedian | 55.0 | 55.9 | 54.5 | 54.9 | 56.7 | 57.2 | 52.0 | 52.7 |
| | TrimmedMean | 76.0 | 76.0 | 72.6 | 72.7 | 75.4 | 75.4 | 74.0 | 74.1 |
| | ARFED | 78.0 | 78.1 | 76.8 | 76.9 | 77.3 | 77.4 | 76.2 | 76.3 |

trimmed mean achieves better performance than CwMedian and can remove the performance loss, it can be said that ARFED outperforms both of them and gives the highest accuracy scores among all these methods. In other words, ARFED successfully defends against malicious participants and gets an accuracy score very close to when all participants are trusted. The performance of accuracy curves for MNIST-2NN can be examined in Fig. 3 and for other datasets in Appendix A.3.

5.3. Byzantine attacks

In the literature, it has been shown that Byzantine attacks are more effective than data poisoning attacks [19,24,29]. Our experiments, which have shown that there is dramatic performance degradation under Byzantine attacks, are in line with the previous studies. In addition, the performance degradation caused by organized attackers is more severe than independent attackers.

When the data distributions of the participants are Non-IID, the performance scores get worse compared to IID cases and Non-IID attacks with organized attackers are the most harmful case for the performance. Again, as the number of malicious participants increases, the performance degradation increases, too.

For IID cases, all defense strategies are able to tolerate the negative effects of Byzantine attacks as if there has been no malicious participant (Table 6 and Fig. 4). For Non-IID cases, the CwMedian is able to prevent performance degradation up to a degree; however, it cannot eliminate the performance degradation caused by attacks as well as ARFED and trimmed mean. ARFED gets better scores than CwMedian for all data sets. As Fig. 4 and Table 7 shows CwMedian can catch the ARFED for only Non-IID experiments of MNIST-CNN.

5.4. Partial knowledge attack

In line with previous experiments on other attack types, Tables 8 and 9 show that Non-IID attacks are more severe than IID attacks. Moreover, as the ratio of malicious participants increases, the performance degrades more and when malicious participants are organized, the degradation becomes more severe for both IID and Non-IID cases.

For IID cases, all strategies can reverse the performance degradation caused by the partial knowledge attack but ARFED can

Table 6

Accuracies under Byzantine attacks at different attacker ratios in IID settings. The best results are bold.

| | | Organized | | | | Independent | | | |
|------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | | m=10% | | m=20% | | m=10% | | m=20% | |
| | | min max | | min max | | min max | | min max | |
| | | min | max | min | max | min | max | min | max |
| MNIST 2NN | NoDefense | 57.0 | 68.9 | 31.9 | 43.9 | 86.0 | 90.7 | 76.2 | 86.1 |
| | CwMedian | 97.2 | 97.2 | 97.3 | 97.4 | 97.2 | 97.3 | 97.3 | 97.3 |
| | TrimmedMean | 97.4 | 97.5 | 97.4 | 97.5 | 97.5 | 97.5 | 97.4 | 97.4 |
| | ARFED | 97.5 | 97.5 | 97.5 | 97.6 | 97.5 | 97.5 | 97.5 | 97.6 |
| MNIST CNN | NoDefense | 54.1 | 79.3 | 10.0 | 17.0 | 91.9 | 95.0 | 70.8 | 87.2 |
| | CwMedian | 98.9 | 98.9 | 98.9 | 98.9 | 98.9 | 98.9 | 98.8 | 98.8 |
| | TrimmedMean | 99.0 | 99.0 | 98.9 | 99.0 | 98.9 | 99.0 | 98.9 | 98.9 |
| | ARFED | 98.9 | 98.9 | 98.8 | 98.8 | 98.9 | 99.0 | 98.8 | 98.8 |
| Fashion MNIST | NoDefense | 15.3 | 48.2 | 9.80 | 21.1 | 65.0 | 79.2 | 36.8 | 60.2 |
| | CwMedian | 89.9 | 90.1 | 90.1 | 90.2 | 90.0 | 90.2 | 90.0 | 90.1 |
| | TrimmedMean | 90.2 | 90.3 | 90.2 | 90.3 | 90.4 | 90.4 | 90.2 | 90.3 |
| | ARFED | 90.2 | 90.3 | 90.4 | 90.5 | 90.2 | 90.3 | 90.3 | 90.4 |
| CIFAR10 | NoDefense | 8.7 | 12.2 | 8.3 | 11.3 | 8.9 | 13.2 | 8.9 | 11.1 |
| | CwMedian | 75.7 | 75.8 | 75.0 | 75.0 | 77.6 | 77.7 | 75.1 | 75.1 |
| | TrimmedMean | 78.8 | 78.9 | 75.6 | 75.7 | 76.8 | 76.9 | 77.0 | 77.1 |
| | ARFED | 77.8 | 77.9 | 77.4 | 77.5 | 77.1 | 77.2 | 77.3 | 77.3 |

Table 7

Accuracies under Byzantine attacks at different attacker ratios in Non-IID settings. The best results are bold.

| | | Organized | | | | Independent | | | |
|------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | | m=10% | | m=20% | | m=10% | | m=20% | |
| | | min max | | min max | | min max | | min max | |
| | | min | max | min | max | min | max | min | max |
| MNIST 2NN | NoDefense | 26.1 | 36.0 | 14.9 | 26.6 | 45.6 | 61.6 | 16.5 | 34.1 |
| | CwMedian | 85.6 | 89.5 | 90.3 | 92.7 | 83.4 | 88.8 | 89.3 | 90.9 |
| | TrimmedMean | 95.8 | 96.1 | 95.4 | 95.6 | 95.9 | 96.1 | 94.8 | 95.0 |
| | ARFED | 96.1 | 96.2 | 95.9 | 96.1 | 96.1 | 96.2 | 95.9 | 96.1 |
| MNIST CNN | NoDefense | 15.6 | 33.5 | 9.00 | 15.7 | 46.1 | 72.4 | 17.4 | 40.2 |
| | CwMedian | 97.4 | 97.5 | 97.6 | 97.7 | 97.0 | 97.2 | 96.8 | 97.1 |
| | TrimmedMean | 98.8 | 98.8 | 98.6 | 98.7 | 98.8 | 98.8 | 98.5 | 98.6 |
| | ARFED | 98.7 | 98.8 | 98.6 | 98.7 | 98.7 | 98.8 | 98.6 | 98.7 |
| Fashion MNIST | NoDefense | 5.70 | 23.6 | 8.60 | 18.0 | 19.7 | 39.6 | 10.0 | 22.2 |
| | CwMedian | 79.4 | 80.8 | 78.4 | 79.6 | 79.7 | 80.8 | 77.7 | 78.3 |
| | TrimmedMean | 84.9 | 86.2 | 85.9 | 86.6 | 84.2 | 86.0 | 83.2 | 84.3 |
| | ARFED | 82.3 | 85.6 | 85.3 | 86.5 | 80.8 | 83.8 | 84.5 | 86.3 |
| CIFAR10 | NoDefense | 8.9 | 11.4 | 8.2 | 11.5 | 7.9 | 10.3 | 8.7 | 13.4 |
| | CwMedian | 57.7 | 59.0 | 71.8 | 72.0 | 62.0 | 62.3 | 58.0 | 61.2 |
| | TrimmedMean | 76.4 | 76.5 | 76.3 | 76.4 | 76.6 | 76.6 | 75.5 | 75.6 |
| | ARFED | 77.6 | 77.6 | 77.5 | 77.6 | 75.2 | 75.3 | 77.2 | 77.2 |

achieve the highest score in all cases. The difference between ARFED and other defense strategies becomes more visible for Fashion-MNIST and CIFAR10.

When the data of participants are Non-IID and the attackers are independent, CwMedian achieves worse accuracy scores than the vanilla FedAvg while trimmed mean can tolerate the performance loss caused by the partial knowledge attack. Still, ARFED gets the highest scores among all of them. Moreover, When the data of participants are Non-IID and the attackers are organized, CwMedian can provide a slight improvement and TrimmedMean can get rid of the performance loss up to a point. ARFED successfully defends against the attack and achieves accuracy scores very close to when all participants are trusted. The performance of accuracy curves for MNIST-2NN can be examined in Fig. 5.

6. Conclusion

This study proposes ARFED, an assumption-free attack-resistant federated averaging algorithm based on outlier elimina-

Table 8

Accuracies under partial knowledge attacks at different attacker ratios in IID settings. The best results are bold.

| | | Organized | | | | Independent | | | |
|------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | | m=10% | | m=20% | | m=10% | | m=20% | |
| | | min max | | min max | | min max | | min max | |
| | | min | max | min | max | min | max | min | max |
| MNIST 2NN | NoDefense | 59.0 | 87.1 | 11.2 | 12.6 | 95.2 | 95.4 | 94.6 | 94.8 |
| | CwMedian | 96.1 | 96.2 | 94.1 | 94.4 | 96.9 | 97.0 | 96.6 | 96.6 |
| | TrimmedMean | 95.1 | 95.2 | 88.3 | 94.7 | 96.7 | 96.8 | 95.9 | 96.0 |
| | ARFED | 97.4 | 97.5 | 97.5 | 97.5 | 97.4 | 97.5 | 97.5 | 97.5 |
| MNIST CNN | NoDefense | 98.1 | 98.2 | 86.9 | 90.7 | 97.9 | 98.0 | 97.3 | 97.4 |
| | CwMedian | 98.4 | 98.4 | 96.7 | 96.8 | 98.8 | 98.9 | 98.7 | 98.7 |
| | TrimmedMean | 97.6 | 97.7 | 96.5 | 97.7 | 98.5 | 98.5 | 98.1 | 98.1 |
| | ARFED | 98.9 | 99.0 | 98.9 | 98.9 | 98.9 | 99.0 | 98.9 | 98.9 |
| Fashion MNIST | NoDefense | 84.4 | 84.8 | 73.1 | 74.7 | 88.8 | 89.1 | 84.7 | 85.0 |
| | CwMedian | 88.9 | 89.0 | 85.9 | 86.0 | 89.7 | 89.8 | 89.2 | 89.3 |
| | TrimmedMean | 88.0 | 88.2 | 84.4 | 84.7 | 89.7 | 89.9 | 88.8 | 89.0 |
| | ARFED | 90.4 | 90.5 | 90.1 | 90.2 | 90.2 | 90.3 | 90.3 | 90.4 |
| CIFAR10 | NoDefense | 68.3 | 68.4 | 50.4 | 50.5 | 74.1 | 74.1 | 66.2 | 66.3 |
| | CwMedian | 72.6 | 72.6 | 65.7 | 65.8 | 73.8 | 73.9 | 71.7 | 71.8 |
| | TrimmedMean | 71.2 | 71.3 | 62.4 | 62.5 | 74.0 | 74.1 | 70.3 | 70.4 |
| | ARFED | 78.2 | 78.3 | 76.9 | 77.0 | 77.4 | 77.4 | 78.5 | 78.5 |

Table 9

Accuracies under partial knowledge attacks at different attacker ratios in Non-IID settings. The best results are bold.

| | | Organized | | | | Independent | | | |
|------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | | m=10% | | m=20% | | m=10% | | m=20% | |
| | | min max | | min max | | min max | | min max | |
| | | min | max | min | max | min | max | min | max |
| MNIST 2NN | NoDefense | 50.9 | 60.7 | 13.2 | 16.9 | 87.5 | 92.3 | 85.2 | 90.7 |
| | CwMedian | 62.6 | 67.2 | 21.1 | 35.2 | 78.4 | 84.6 | 63.3 | 78.8 |
| | TrimmedMean | 77.1 | 81.5 | 32.0 | 38.0 | 94.1 | 94.7 | 91.0 | 91.8 |
| | ARFED | 96.1 | 96.3 | 95.9 | 96.2 | 96.1 | 96.3 | 95.9 | 96.2 |
| MNIST CNN | NoDefense | 51.8 | 58.8 | 16.5 | 24.7 | 97.6 | 97.9 | 91.2 | 97.1 |
| | CwMedian | 80.4 | 87.9 | 38.1 | 39.1 | 96.5 | 96.8 | 95.6 | 96.0 |
| | TrimmedMean | 88.0 | 88.7 | 27.0 | 43.6 | 98.3 | 98.4 | 97.3 | 97.6 |
| | ARFED | 98.8 | 98.9 | 98.7 | 98.8 | 98.8 | 99.0 | 98.6 | 98.7 |
| Fashion MNIST | NoDefense | 42.0 | 51.4 | 9.3 | 21.9 | 79.3 | 83.5 | 75.5 | 80.6 |
| | CwMedian | 53.6 | 54.6 | 26.8 | 28.0 | 74.8 | 75.8 | 66.4 | 68.0 |
| | TrimmedMean | 54.4 | 58.0 | 38.6 | 51.3 | 81.6 | 83.5 | 79.8 | 81.2 |
| | ARFED | 82.0 | 82.8 | 85.6 | 87.5 | 79.6 | 84.5 | 86.0 | 87.3 |
| CIFAR10 | NoDefense | 52.7 | 52.9 | 30.9 | 31.7 | 72.0 | 72.0 | 65.5 | 65.6 |
| | CwMedian | 53.0 | 53.3 | 38.6 | 38.8 | 61.7 | 62.1 | 45.8 | 46.3 |
| | TrimmedMean | 68.4 | 68.5 | 46.7 | 46.8 | 72.5 | 72.6 | 66.0 | 66.0 |
| | ARFED | 77.8 | 77.8 | 77.3 | 77.3 | 77.8 | 77.9 | 76.4 | 76.5 |

tion, and conducts comprehensive experiments in various FL settings. These experiments reveal that Byzantine attacks and partial knowledge attacks are dramatically more severe than label flipping attacks. Moreover, attacks in the Non-IID cases are more effective than IID cases and organized attackers can severely compromise the performance of the main model more compared to independent attackers.

Although CwMedian, TrimmedMean, and ARFED tolerate performance loss in the presence of attacks in IID cases, the CwMedian performs poorly in Non-IID cases; it may even perform worse than the vanilla FedAvg. For Non-IID cases, ARFED shows better performance recovery than CwMedian in all attack types. Moreover, ARFED outperforms TrimmedMean in label flipping attacks and partial knowledge attacks, but they get similar results in Byzantine attacks. The likely reason for this is that the parameter updates sent by malicious participants are extreme and lie in the distribution's tails for Byzantine attacks. In this way, TrimmedMean can detect poisoned parameter updates more

easily. On the other hand, in label flipping attacks and partial knowledge attacks, changes in parameter updates are more likely to be moderate; therefore, the TrimmedMean cannot provide the same performance. Moreover, it is also worth keeping in mind that TrimmedMean needs information about the malicious participant ratio in the system, while ARFED does not make such an assumption.

We put forward experimental evidence to show that ARFED removes performance loss even under organized attacks and in Non-IID cases. There are many attack-robust aggregation methods and mechanisms for FL in the literature, but they mainly focused on ensuring convergence under some assumptions such as data distribution, knowledge of malicious participant ratio, and update similarity of participants. Our work highlights the shortfall in current theoretical convergence guaranteed methods and presents a broader research goal to create aggregation mechanisms that work in harmony with Non-IID data, which is one of the key properties of FL.

Our method is mainly based on outlier elimination which may tolerate up to a certain number of malicious participants in the system. As the ratio of attackers increases in the FL setting, they will have a high impact on the distribution of distances. Distance distributions also depend on some other parameters, such as the severity of the Non-IID data and coordination of attackers; therefore, to what extent ARFED can handle malicious participants is beyond the scope of this study and reserved as future work.

Moreover, ARFED allows a participant to be included aggregation step of the FL round if the participant is reliable at all layers with *and*(\wedge) operation. The effect of losing participants on different layers will be examined for the more complex model architectures as future work.

Besides, MNIST, Fashion-MNIST, and CIFAR10 datasets were used in this study because they are widely used in federated learning and defense mechanism studies similar to the one we propose. However, the performance of our proposed method could be evaluated on different datasets in the future.

CRediT authorship contribution statement

Ece Isik-Polat: Conceptualization, Methodology, Software, Formal analysis, Writing – original draft. **Gorkem Polat:** Conceptualization, Writing – original draft. **Altan Kocyigit:** Conceptualization, Writing – review & editing, Supervision.

Declaration of competing 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.

Data availability

Please see the link: <https://github.com/eceisik/ARFED>.

Appendix

A.1. Model architectures

The model architectures used for each datasets are shown in Table 10 (MNIST-2NN [9]), Table 11 (MNIST-CNN [9]), Table 13 (CIFAR10 [18]) and Table 12 (Fashion-MNIST). All activation functions for all models are ReLU.

The FL setting parameters used for each dataset are shown in Table 14. Learning rate scheduling and clipping were not applied to MNIST and Fashion-MNIST, therefore related parameters are set as N/A (Not Applicable).

Table 10

Model architecture of MNIST-2NN.

| Layer | Size |
|-----------------|------------|
| Fully Connected | (784, 200) |
| Fully Connected | (200, 200) |
| Fully Connected | (200, 10) |

Table 11

Model architecture of MNIST-CNN.

| textbfLayer | Size |
|-----------------|-------------|
| Conv | 32@5 × 5 |
| Max Pooling | 2 × 2 |
| Conv | 64@5 × 5 |
| Max Pooling | 2 × 2 |
| Fully Connected | (1024, 512) |
| Fully Connected | (512, 10) |

Table 12

Model architecture of Fashion-MNIST.

| Layer | Size |
|-----------------|-----------------|
| Conv | 32@5 × 5, pad=2 |
| Max Pooling | 2 × 2 |
| Conv | 64@5 × 5, pad=2 |
| Max Pooling | 2 × 2 |
| Fully Connected | (3136, 500) |
| Fully Connected | (500, 10) |

Table 13

Model architecture of CIFAR10.

| Layer | Size |
|-----------------|------------------|
| Conv | 32@3 × 3, pad=1 |
| Conv | 32@3 × 3, pad=1 |
| Max Pooling | 2 × 2 |
| Conv | 64@3 × 3, pad=1 |
| Conv | 64@3 × 3, pad=1 |
| Max Pooling | 2 × 2 |
| Conv | 128@3 × 3, pad=1 |
| Conv | 128@3 × 3, pad=1 |
| Max Pooling | 2 × 2 |
| Fully Connected | (2048, 128) |
| Fully Connected | (128, 10) |

A.2. Machine configuration and used platforms

Tesla P100-PCI-E-16 GB, Tesla V100-SXM2-16 GB, and NVIDIA A100-SXM-80 GB were used for the experiments. According to the used data sets, the average running time of an experiment on the A100 machine is as follows: 1 h for MNIST-2NN, 1.4 h for MNIST-CNN, 3.3 h for FASHION-MNIST, and 7.6 h for CIFAR.

Due to the versatility of the experimental settings, using an off-the-shelf platform did not provide the necessary flexibility; therefore, we chose to code ourselves. All implemented methods and designed experiments can be seen via <https://github.com/eceisik/ARFED>. The required packages for the environment setup are also listed here.

A.3. Experiments

In order to evaluate the performance of ARFED, extensive experiments covering different scenarios such as whether the attackers are organized, different types of attacks, the effect of the data distribution, and malicious participant ratio have been carried out on different datasets with different model architectures.

The experimental results of

- the label flipping attacks are summarized for IID setting in Table 4 and for Non-IID setting in Table 5

Table 14
FL setting parameters used in experiments.

| Parameters | MNIST-2NN | MNIST-CNN | Fashion-MNIST | CIFAR10 |
|--|-----------|-----------|---------------|----------|
| Number of participant (n) | 100 | 100 | 100 | 100 |
| Communication round (t) | 200 | 200 | 200 | 500 |
| Number of label in each participant in IID setting | 2 | 2 | 2 | 5 |
| Number of label in each Participant in Non-IID setting | 10 | 10 | 10 | 10 |
| Batch size | 32 | 32 | 25 | 100 |
| Number of epoch | 10 | 10 | 10 | 10 |
| Momentum | 0.9 | 0.9 | 0.9 | 0.9 |
| Learning rate | 0.01 | 0.01 | 0.002 | 0.0015 |
| Minimum learning rate (min_lr) | N/A | N/A | N/A | 0.000010 |
| lr scheduler factor | N/A | N/A | N/A | 0.2 |
| Best threshold | N/A | N/A | N/A | 0.0001 |
| Clipping threshold | N/A | N/A | N/A | 10 |

- the Byzantine attacks for IID setting in Table 6 and for Non-IID setting in Table 7 and
- the partial knowledge attacks for IID setting in Table 8 and for Non-IID setting in Table 9

for each data set under Section 5. Due to the space limitation and to improve the readability of this study, the figures of the experiments of MNIST CNN, Fashion-MNIST, and CIFAR10 are presented here.

The figures reveal that the experiments run for different datasets and model architectures are in line with previous findings and are valid for all data sets. For example, the performance loss in Byzantine attacks is greatest, and partial knowledge attacks cause more performance degradation than label flipping attacks. When the data distributions of the participants are Non-IID, the performance degrades more compared to IID cases. As the ratio of malicious participants increases, performance degrades more. Organized attackers cause more performance degradation. The worst accuracy score is recorded when the attackers are organized and the data distribution of the participants is Non-IID.

ARFED could eliminate the harmful effects of all attack types for both IID and non-IID cases and achieve accuracy scores close to when all collaborators are trusted cases (no attack case). On the other hand, CwMedian is not able to handle the attacks in the Non-IID setting. CwMedian could tolerate the performance degradation in only IID cases. For non-IID cases, it could show only a slight improvement or worsen the performance degradation. ARFED generally performs better than TrimmedMean in label flipping attacks and partial knowledge attacks, but they get similar scores in Byzantine attacks. However, it is worth remembering that TrimmedMean requires information of the malicious collaborator ratio in the system, while ARFED does not make such an assumption.

A.3.1. MNIST CNN experiments

- Fig. 6 presents the results of the experiments carried out for label flipping attack
- Fig. 7 present the results of the experiments carried out for byzantine attack and
- Fig. 8 present the results of the experiments carried out for adaptive partial knowledge attack

on MNIST dataset with CNN model architecture.

A.3.2. Fashion-MNIST experiments

- Fig. 9 presents the results of the experiments carried out for label flipping attack
- Fig. 10 present the results of the experiments carried out for byzantine attack and
- Fig. 11 present the results of the experiments carried out for adaptive partial knowledge attack

on Fashion-MNIST.

A.3.3. CIFAR10 experiments

- Fig. 12 presents the results of the experiments carried out for label flipping attack
- Fig. 13 present the results of the experiments carried out for byzantine attack and
- Fig. 14 present the results of the experiments carried out for adaptive partial knowledge attack

on CIFAR10.

A.3.4. Number of reliable and outlier participants for CIFAR10 experiments

The working mechanism of ARFED decides which participant is eligible to be included in the aggregation step based on whether the parameters sent by the participant for each layer of the model architecture are in the safe interval. This all-or-nothing approach of ARFED may raise concerns that too many participants may be discarded from the aggregation step, and too much valuable information may be lost through the layers. The model architecture used for the CIFAR10 dataset has more layers than the architectures that are used for other datasets. For this reason, the risk of losing too many participants brought by the algorithm's $\text{and}(\wedge)$ operation is expected to be best observed in the CIFAR10 set.

Fig. 15 illustrates the number of participants marked as reliable and included aggregation step versus the number of participants marked as outliers and discarded from aggregation in label flipping attacks with different scenarios. Fig. 16 illustrates the number of participants marked as reliable and included aggregation step versus the number of participants marked as outliers and discarded from aggregation in Byzantine attacks with different scenarios. The number of discarded participants is in line with the malicious participant ratio.

A.4. Box plot factor comparison

Different factor values were tested to show the performance impact of how strict the algorithm is in labeling a participant as reliable. For this reason, different factor values were applied when determining the lower distance threshold ($\text{min_}d_l^i$) and upper distance threshold ($\text{max_}d_l^i$) in Lines 19–20 of the Algorithm 1.

Table 15 and Fig. 17 present the results of the partial knowledge attack scenarios when the participants' data are IID. Table 16 and Fig. 18 show the results of the partial knowledge attack scenarios when the participants' data are non-IID. These graphs show no significant difference between different factor values, i.e., 0, 1, 1.5, and 2, especially in the IID setting. However, in the non-IID setting, although $f=0$ achieves performance like others, the accuracy graph has oscillations that can signal a convergence problem. One possible reason might be that the data was non-IID, and the algorithm could not obtain a good enough sample

Table 15

Accuracy scores obtained on test set under partial knowledge attacks with different malicious participant ratios in the IID setting. The best results are bold.

| | Organized | | | | Independent | | | |
|--------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | m=10% | | m=20% | | m=10% | | m=20% | |
| | min | max | min | max | min | max | min | max |
| NoDefense | 60.4 | 87.6 | 11.3 | 12.8 | 96.4 | 96.5 | 94.8 | 95.0 |
| CwMedian | 97.4 | 97.5 | 94.7 | 94.8 | 98.2 | 98.2 | 97.5 | 97.5 |
| TrimmedMean | 96.4 | 96.5 | 89.0 | 94.9 | 98.3 | 98.3 | 97.1 | 97.1 |
| ARFED f1.5 | 99.0 | 99.0 | 98.9 | 98.9 | 99.0 | 99.0 | 98.9 | 98.9 |
| ARFED f0 | 98.7 | 98.7 | 98.8 | 98.8 | 98.7 | 98.7 | 98.8 | 98.8 |
| ARFED f1 | 98.9 | 99.0 | 98.9 | 98.9 | 98.9 | 99.0 | 98.9 | 98.9 |
| ARFED f2 | 99.0 | 99.0 | 98.9 | 98.9 | 99.0 | 99.0 | 98.9 | 98.9 |

Table 16

Accuracy scores obtained on test set under partial knowledge attacks with different malicious participant ratios in the Non-IID setting. The best results are bold.

| | Organized | | | | Independent | | | |
|--------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | m=10% | | m=20% | | m=10% | | m=20% | |
| | min | max | min | max | min | max | min | max |
| NoDefense | 51.1 | 60.4 | 13.2 | 16.7 | 88.5 | 92.9 | 85.6 | 91.0 |
| CwMedian | 62.8 | 67.0 | 21.8 | 35.6 | 78.4 | 84.8 | 63.6 | 79.6 |
| TrimmedMean | 77.2 | 81.9 | 31.6 | 38.4 | 94.8 | 95.4 | 91.5 | 92.7 |
| ARFED f1.5 | 97.1 | 97.2 | 96.8 | 96.9 | 97.1 | 97.2 | 96.8 | 96.9 |
| ARFED f0 | 96.4 | 97.2 | 96.2 | 97.3 | 96.4 | 97.2 | 96.2 | 97.3 |
| ARFED f1 | 97.0 | 97.1 | 96.6 | 96.9 | 97.0 | 97.1 | 96.6 | 96.9 |
| ARFED f2 | 97.1 | 97.2 | 96.8 | 96.9 | 97.1 | 97.2 | 96.8 | 96.9 |

space by eliminating too many participants from the main model aggregation step.

References

- [1] K. Osmundsen, J. Iden, B. Bygstad, Digital transformation: Drivers, success factors, and implications, in: MCIS, 2018, p. 37.
- [2] C. Naseeb, AI and ML-driving and exponentiating sustainable and quantifiable digital transformation, in: 2020 IEEE 44th Annual Computers, Software, and Applications Conference, COMPSAC, 2020, pp. 316–321, <http://dx.doi.org/10.1109/COMPSAC48688.2020.0-227>.
- [3] A. Sestino, M.I. Prete, L. Piper, G. Guido, Internet of things and big data as enablers for business digitalization strategies, *Technovation* 98 (2020) 102173, <http://dx.doi.org/10.1016/j.technovation.2020.102173>, URL <https://www.sciencedirect.com/science/article/pii/S0166497220300456>.
- [4] X.-W. Chen, X. Lin, Big data deep learning: Challenges and perspectives, *IEEE Access* 2 (2014) 514–525, <http://dx.doi.org/10.1109/ACCESS.2014.2325029>.
- [5] J. Wang, Y. Yang, T. Wang, R.S. Sherratt, J. Zhang, Big data service architecture: a survey, *J. Internet Technol.* 21 (2) (2020) 393–405.
- [6] N. Gruschka, V. Mavroeidis, K. Vishi, M. Jensen, Privacy issues and data protection in big data: A case study analysis under GDPR, in: 2018 IEEE International Conference on Big Data (Big Data), 2018, pp. 5027–5033, <http://dx.doi.org/10.1109/BigData.2018.8622621>.
- [7] B. Custers, A.M. Sears, F. Dechesne, I. Georgieva, T. Tani, S. Van der Hof, EU Personal Data Protection in Policy and Practice, Vol. 29, Springer, 2019.
- [8] P. Voigt, A. Von dem Bussche, The eu general data protection regulation (gdpr), in: A Practical Guide, first ed., Vol. 10, Springer International Publishing, Cham, 2017, 3152676.
- [9] B. McMahan, E. Moore, D. Ramage, S. Hampson, B.A. y Arcas, Communication-efficient learning of deep networks from decentralized data, in: Artificial Intelligence and Statistics, PMLR, 2017, pp. 1273–1282.
- [10] S. Singh, S. Rathore, O. Alfarraj, A. Tolba, B. Yoon, A framework for privacy-preservation of IoT healthcare data using federated learning and blockchain technology, *Future Gener. Comput. Syst.* 129 (2022) 380–388, <http://dx.doi.org/10.1016/j.future.2021.11.028>, URL <https://www.sciencedirect.com/science/article/pii/S0167739X21004726>.
- [11] Z. Zheng, Y. Zhou, Y. Sun, Z. Wang, B. Liu, K. Li, Applications of federated learning in smart cities: recent advances, taxonomy, and open challenges, *Connect. Sci.* 34 (1) (2022) 1–28.
- [12] T.R. Gadekallu, Q.-V. Pham, T. Huynh-The, S. Bhattacharya, P.K.R. Madrikunta, M. Liyanage, Federated learning for big data: A survey on opportunities, applications, and future directions, 2021, arXiv preprint [arXiv:2110.04160](https://arxiv.org/abs/2110.04160).
- [13] J. Konečný, H.B. McMahan, D. Ramage, P. Richtárik, Federated optimization: Distributed machine learning for on-device intelligence, 2016, arXiv preprint [arXiv:1610.02527](https://arxiv.org/abs/1610.02527).
- [14] F. Sattler, K.-R. Müller, W. Samek, Clustered federated learning: Model-agnostic multitask optimization under privacy constraints, *IEEE Trans. Neural Netw. Learn. Syst.* 32 (8) (2021) 3710–3722, <http://dx.doi.org/10.1109/TNNLS.2020.3015958>.
- [15] X. Ma, J. Zhu, Z. Lin, S. Chen, Y. Qin, A state-of-the-art survey on solving non-IID data in federated learning, *Future Gener. Comput. Syst.* 135 (2022) 244–258, <http://dx.doi.org/10.1016/j.future.2022.05.003>, URL <https://www.sciencedirect.com/science/article/pii/S0167739X22001686>.
- [16] T. Li, A.K. Sahu, A. Talwalkar, V. Smith, Federated learning: Challenges, methods, and future directions, *IEEE Signal Process. Mag.* 37 (3) (2020) 50–60, <http://dx.doi.org/10.1109/MSP.2020.2975749>.
- [17] V. Mothukuri, R.M. Parizi, S. Pouriyeh, Y. Huang, A. Dehghantanha, G. Srivastava, A survey on security and privacy of federated learning, *Future Gener. Comput. Syst.* 115 (2021) 619–640, <http://dx.doi.org/10.1016/j.future.2020.10.007>, URL <https://www.sciencedirect.com/science/article/pii/S0167739X20329848>.
- [18] V. Tolpegin, S. Truex, M.E. Gursoy, L. Liu, Data poisoning attacks against federated learning systems, in: European Symposium on Research in Computer Security, Springer, 2020, pp. 480–501.
- [19] M. Fang, X. Cao, J. Jia, N. Gong, Local model poisoning attacks to Byzantine-Robust federated learning, in: 29th USENIX Security Symposium (USENIX Security 20), USENIX Association, 2020, pp. 1605–1622, URL <https://www.usenix.org/conference/usenixsecurity20/presentation/fang>.
- [20] M. Duan, D. Liu, X. Chen, Y. Tan, J. Ren, L. Qiao, L. Liang, Astraea: Self-balancing federated learning for improving classification accuracy of mobile deep learning applications, in: 2019 IEEE 37th International Conference on Computer Design, ICCD, IEEE, 2019, pp. 246–254.
- [21] R. Pathak, M.J. Wainwright, FedSplit: an algorithmic framework for fast federated optimization, in: H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, H. Lin (Eds.), Advances in Neural Information Processing Systems, Vol. 33, Curran Associates, Inc., 2020, pp. 7057–7066, URL <https://proceedings.neurips.cc/paper/2020/file/4ebd440d99504722d80de606ea8507da-Paper.pdf>.
- [22] H. Yuan, T. Ma, Federated accelerated stochastic gradient descent, *Adv. Neural Inf. Process. Syst.* 33 (2020).
- [23] A.N. Bhagoji, S. Chakraborty, P. Mittal, S. Calo, Analyzing federated learning through an adversarial lens, in: International Conference on Machine Learning, PMLR, 2019, pp. 634–643.
- [24] F. Sattler, K.-R. Müller, T. Wiegand, W. Samek, 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.
- [25] C. Fung, C.J. Yoon, I. Beschastnikh, Mitigating sybils in federated learning poisoning, 2018, arXiv preprint [arXiv:1808.04866](https://arxiv.org/abs/1808.04866).
- [26] S. Shen, S. Tople, P. Saxena, Auror: Defending against poisoning attacks in collaborative deep learning systems, in: Proceedings of the 32nd Annual Conference on Computer Security Applications, 2016, pp. 508–519.
- [27] P. Blanchard, E.M. El Mhamdi, R. Guerraoui, J. Stainer, Machine learning with adversaries: Byzantine tolerant gradient descent, in: Proceedings of the 31st International Conference on Neural Information Processing Systems, 2017, pp. 118–128.
- [28] E.M. El Mhamdi, R. Guerraoui, S. Rouault, The hidden vulnerability of distributed learning in byzantium, in: Proceedings of the 35th International Conference on Machine Learning, in: Proceedings of Machine Learning Research, Vol. 80, PMLR, 2018, pp. 3521–3530.
- [29] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, V. Shmatikov, How to backdoor federated learning, in: International Conference on Artificial Intelligence and Statistics, PMLR, 2020, pp. 2938–2948.
- [30] P. Kairouz, H.B. McMahan, B. Avent, A. Bellet, M. Bennis, A.N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, et al., Advances and open problems in federated learning, 2019, arXiv preprint [arXiv:1912.04977](https://arxiv.org/abs/1912.04977).
- [31] H. Wang, K. Sreenivasan, S. Rajput, H. Vishwakarma, S. Agarwal, J.-y. Sohn, K. Lee, D. Papailiopoulos, Attack of the tails: Yes, you really can backdoor federated learning, in: H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, H. Lin (Eds.), Advances in Neural Information Processing Systems, Vol. 33, Curran Associates, Inc., 2020, pp. 16070–16084, URL <https://proceedings.neurips.cc/paper/2020/file/b8ffa41d4e492f0fad2f13e29e1762eb-Paper.pdf>.
- [32] L. Chen, H. Wang, Z. Charles, D. Papailiopoulos, Draco: Byzantine-resilient distributed training via redundant gradients, in: International Conference on Machine Learning, PMLR, 2018, pp. 903–912.
- [33] C. Xie, S. Koyejo, I. Gupta, Zeno: Distributed stochastic gradient descent with suspicion-based fault-tolerance, in: International Conference on Machine Learning, PMLR, 2019, pp. 6893–6901.
- [34] C. Xie, S. Koyejo, I. Gupta, Zeno+: Robust fully asynchronous SGD, in: International Conference on Machine Learning, PMLR, 2020, pp. 10495–10503.

- [35] D. Yin, Y. Chen, R. Kannan, P. Bartlett, Byzantine-robust distributed learning: Towards optimal statistical rates, in: International Conference on Machine Learning, PMLR, 2018, pp. 5650–5659.
- [36] J. Zhang, S. Zhong, T. Wang, H.-C. Chao, J. Wang, Blockchain-based systems and applications: a survey, J. Internet Technol. 21 (1) (2020) 1–14.
- [37] Y. Qu, S.R. Pokhrel, S. Garg, L. Gao, Y. Xiang, A blockchained federated learning framework for cognitive computing in industry 4.0 networks, IEEE Trans. Ind. Inform. 17 (4) (2021) 2964–2973, <http://dx.doi.org/10.1109/TII.2020.3007817>.
- [38] G. Baruch, M. Baruch, Y. Goldberg, A little is enough: Circumventing defenses for distributed learning, in: H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox, R. Garnett (Eds.), Advances in Neural Information Processing Systems, Vol. 32, Curran Associates, Inc., 2019, URL <https://proceedings.neurips.cc/paper/2019/file/ec1c59141046cd1866bbcbdfb6ae31d4-Paper.pdf>.
- [39] S. Rajput, H. Wang, Z. Charles, D. Papailiopoulos, DETOX: A redundancy-based framework for faster and more robust gradient aggregation, in: H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox, R. Garnett (Eds.), Advances in Neural Information Processing Systems, Vol. 32, Curran Associates, Inc., 2019, URL <https://proceedings.neurips.cc/paper/2019/file/415185ea244ea2b2bedeb0449b926802-Paper.pdf>.
- [40] X. Li, K. Huang, W. Yang, S. Wang, Z. Zhang, On the convergence of FedAvg on non-IID data, in: International Conference on Learning Representations, 2020, URL <https://openreview.net/forum?id=HjxNANVtDS>.
- [41] Y. LeCun, The MNIST database of handwritten digits, 1998, <http://yann.lecun.com/exdb/mnist/>.
- [42] A. Krizhevsky, G. Hinton, et al., Learning Multiple Layers of Features from Tiny Images, Citeseer, 2009.
- [43] H. Xiao, K. Rasul, R. Vollgraf, Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017, arXiv preprint [arXiv: 1708.07747](https://arxiv.org/abs/1708.07747).



Ece Isik-Polat is doing her Ph.D. in the Department of Information Systems at the Middle East Technical University (METU), and she is currently a research assistant in the same department. She received her M.Sc. degree from the same department, B.Sc. degree in the Department of Statistics, and minor degree in the Department of Sociology from the Middle East Technical University. She is particularly interested in Federated Learning, Data Science, and Machine Learning.

Gorkem Polat received his B.Sc. and M.Sc. in Electrical and Electronics Engineering at Middle East Technical University (METU). He is currently pursuing a Ph.D. degree in the Department of Medical Informatics at the same university. His main research interests are computer vision, particularly medical image processing and federated learning.

Altan Kocyigit received B.Sc., M.Sc., and Ph.D. degrees from METU Electrical and Electronics Engineering Department in 1993, 1997, and 2001, respectively. He has been with METU, Informatics Institute since 2002. His research interests include data science, computer networking, software engineering, and parallel/distributed processing.