# Byzantine Resilient Federated Multi-Task Representation Learning

Tuan Le and Shana Moothedath

*Abstract*—In this paper, we propose BR-MTRL, a Byzantine-resilient multi-task representation learning framework that handles faulty or malicious agents. Our approach leverages representation learning through a shared neural network model, where all clients share fixed layers, except for a client-specific final layer. This structure captures shared features among clients while enabling individual adaptation, making it a promising approach for leveraging client data and computational power in heterogeneous federated settings to learn personalized models. To learn the model, we employ an alternating gradient descent strategy: each client optimizes its local model, updates its final layer, and sends estimates of the shared representation to a central server for aggregation. To defend against Byzantine agents, we employ two robust aggregation methods for client-server communication, Geometric Median and Krum. Our method enables personalized learning while maintaining resilience in distributed settings. We implemented the proposed algorithm in a federated testbed built using Amazon Web Services (AWS) platform and compared its performance with various benchmark algorithms and their variations. Through experiments using real-world datasets, including CIFAR-10 and FEMNIST, we demonstrated the effectiveness and robustness of our approach and its transferability to new unseen clients with limited data, even in the presence of Byzantine adversaries.

*Index Terms*—Byzantine resilience, Multi-task representation learning, Federated learning, Alternating gradient descent

## I. Introduction

Machine learning has traditionally relied on centralized settings, where a single model is trained on large datasets stored in one central location. This approach requires gathering data from all clients into a central data center, which raises significant data privacy concerns. Federated Learning (FL) addresses these challenges by using a distributed framework that preserves data privacy while overcoming the problem of isolated data sources [1]. FL enables multiple clients, such as mobile devices or institutions, to collaboratively train a global model without sharing their raw data. A central server is responsible for coordinating the clients by aggregating the models received from each client and broadcasting the updated model back to them in each round of communication. McMahan *et al.* introduced the Federated Averaging (FedAvg) algorithm, which facilitates the training of a single shared global model designed to achieve optimal performance on average across all participating clients [1]. The approach offers several advantages, such as simplicity in implementation and a streamlined aggregation process, making it scalable for large client populations. Additionally, the global model promotes uniformity, ensuring consistent performance for clients with

T. Le is with the Department of Computer Science and S. Moothedath is with the Department of Electrical and Computer Engineering, Iowa State University, Ames, IA, USA. Email: {tuanle, mshana}@iastate.edu.

similar data distributions. However, FL approaches encounter three primary bottlenecks: (i) *Data heterogeneity:* In real-world scenarios, client tasks often exhibit significantly different underlying data distributions. In such environments, a global model optimized for average loss may underperform across diverse clients, as it fails to account for these variations in data [2]. (ii) *Generalization and adaptability:* Federated models may struggle to generalize to unseen clients that were not part of the training process, limiting their effectiveness in dynamic, real-world applications [3]. (iii) *Byzantine attacks:* FL is particularly vulnerable to Byzantine attacks, as it relies on distributed client updates; a faulty or malicious client can potentially corrupt the global model, compromising performance and reliability [4].

To solve data heterogeneity, many *personalized* FL algorithms have been developed. Yuyang Deng *et al.* proposed the Adaptive Personalized Federated Learning (APFL) algorithm to solve the problem of statistical heterogeneity in federated learning by applying a mix of local and global model optimization [5]. This method allows clients to train personalized models that better generalize on local distributions while still contributing to the learning of a robust global model. In addition, Alireza Fallah *et al.* proposed a method called Personalized Federated Learning via Model-Agnostic Meta-Learning (MAML) to solve the problem of personalized adaptation in federated learning by applying a meta-learning framework [6]. This method allows clients to quickly adapt a shared initial model to their local data distributions by performing a few steps of gradient descent, leading to more personalized models for each user.

To achieve generalizability and adaptability, Multi-Task Representation Learning (MTRL) has emerged as a promising methodology in machine learning, aiming to address multiple related tasks/clients simultaneously by leveraging shared representations [7]. The core principle of MTRL is based on the idea that, even though data in federated settings is often non-i.i.d. (non-independent and identically distributed) across clients, there exists a shared learned representation across multiple clients that effectively captures the common structural elements inherent to each client. This shared foundation enables the use of simple linear predictors in the final model layer, which are then fine-tuned to address the unique variations and requirements of each client. This approach significantly enhances the efficiency of model training processes, as it allows for the transfer and utilization of learned features across *different but related tasks*. Several works have studied MTRL to address different aspects of learning and model performance. Yu Xie *et al.* solved the problem of simultaneously optimizing node classification and link prediction tasks by applying a shared feature extraction module followed by client-specific

modules [8]. Yang *et al.* proposed a deep MTRL framework using tensor factorization, which enables automatic knowledge sharing across all layers of a deep neural network. This method applies to both homogeneous and heterogeneous settings and eliminates the need for manual specification of multi-task sharing strategies [9]. In addition, recent studies [3], [10], [11], [12] have focused on the theoretical understanding of MTRL, particularly when the shared model exhibits specific structures, such as a low-rank representation.

To address the challenge of Byzantine attacks in federated learning, Gouissem *et al.* developed a collaborative cross-check mechanism that enhances the system's resilience by allowing clients to collaboratively identify malicious nodes based on each other's contributions to model updates [13]. While effective, this approach introduces significant computational overhead due to the extensive cross-validation among clients, which can lead to inefficiencies. In contrast, Feng Lin *et al.* tackled Byzantine resilience using robust aggregation techniques, such as geometric median, Krum, and h-Krum, to address the matrix completion problem while ensuring resilience against Byzantine failures [14]. Recently, Li *et al.* studied Byzantine resilient FL for non-iid data by obtaining an inversed artificial gradient [4]. Zhang *et al.* investigated communication-efficient approaches for Byzantine-resilient FL and proposed a projected stochastic gradient descent (SGD) algorithm to enhance system robustness [15]. Their approach employed a Huber function-based robust aggregation with an adaptive threshold selecting strategy at the server to reduce the effects of Byzantine attacks.

Despite these advancements, a unified approach for Byzantine-resilient personalized FL with strong generalization remains unavailable. In this paper, we introduce the Byzantine-resilient Multi-Task Representation Learning (BR-MTRL) framework to jointly address three critical challenges in FL: data heterogeneity, generalization, and Byzantine resilience. By incorporating geometric median and Krum aggregation methods into personalized FL, we enhance the global model's resilience against Byzantine attacks. Our framework ensures that the models learned across distributed agents remain robust, reliable, and secure, even in the presence of adversarial influences, improving both performance and safety.

### A. Contributions

The contributions of this paper are threefold.

- We propose a Byzantine-resilient MTRL algorithm based on personalized federated learning. In our approach, all clients share a common representation while maintaining their local personalized heads. We introduce an alternating gradient descent approach that alternates between solving for the shared model and fine-tuning each client's personalized head. To ensure robustness, our method employs Geometric Median (GM) and Krum-based mechanisms for aggregating client updates at the server, allowing the clients to learn a Byzantine-resilient shared representation while preserving task-specific personalizations.
- We implemented our proposed approach in a truly federated environment by deploying it on a simulation testbed built on the Amazon Web Services (AWS) platform. We performed simulations on datasets, including CIFAR-10

and FEMNIST, and validated the effectiveness of BR-MTRL in mitigating the impact of Byzantine attacks while maintaining a robust learning performance.
- We conducted a meta-test to evaluate the transferability of the learned representation on new, unseen clients with limited data. The BR-MTRL approach was assessed by fixing the shared representation model and learning only the personalized heads of the new clients. We measured the impact of Byzantine clients on new clients' performance and demonstrated improved accuracy using BR-MTRL, ensuring the transferability of our approach.

### B. Related Work

Personalized Federated Learning (PFL) has received a lot of research interest in recent years. Liang *et al.* [16] proposed the Local Global Federated Averaging (LG-FedAvg) method to address data heterogeneity, which focuses on learning low-dimensional features from clients' data. Their approach emphasizes learning multiple local representations and a single global head. On the other hand, Arivazhagan *et al.* [17] proposed FedPer algorithm that explicitly divides the model architecture into base layers, trained across all devices via federated averaging and personalization layers, which are tailored locally on each device using its specific data. This structure not only promotes personalization by adapting the model to local data characteristics but also maintains a cohesive model performance across diverse environments by leveraging a shared base learning across all clients. In addition, Collins *et al.* [3] proposed the FedRep algorithm focusing on learning a global model shared across clients while allowing each client to maintain a personalized head model, similar to [17]. The algorithms in [3], [17] allow clients to collaborate on learning shared representations while maintaining personalized models tailored to their specific data distribution. [3], [17] are closely related to our work; however, they focus on PFL without considering the challenges posed by Byzantine attacks.

To counter Byzantine attacks, many PFL systems incorporate robust aggregation techniques at the server, such as the geometric median and its variations. For example, Li *et al.* introduced AutoGM, an auto-weighted geometric median aggregation rule designed to enhance robustness in the presence of Byzantine failures, specifically model and data poisoning attacks [18]. AutoGM improves the classic geometric median by automatically excluding extreme outliers and re-weighting the remaining parameter updates based on a skewness threshold, offering a more adaptive and flexible solution. Similarly, Wang *et al.* proposed WGM-dSAGA, a federated learning strategy that employs a weighted geometric median for robust aggregation against Byzantine attacks [19]. The strategy incorporates COPOD, a parameter-free outlier detection technique, to identify and assign weights to client updates based on their potential maliciousness. This method effectively mitigates the influence of Byzantine clients on the global model through the weighted geometric median. Meanwhile, He *et al.* tackled Byzantine attacks by introducing a novel projected stochastic block gradient descent method [20]. Unlike the GM-based approaches, their strategy employs a server aggregation technique that combines local models using

a Huber function-based descent step to manage the impact of malicious updates from Byzantine clients. While these studies have made significant strides in mitigating Byzantine attacks, they primarily focus on enhancing aggregation techniques on the server side. In contrast, our work introduces a Byzantine-resilient MTRL framework that learns and leverages a shared representation to enhance the quality of each client's model, while ensuring both transferability and resilience.

## II. PROBLEM SETTING

We address a multi-task supervised learning scenario with $n$ clients, where each $i$-th client possesses $M_i$ labeled samples $\{\mathbf{x}_i^j, y_i^j\}_{j=1}^{M_i}$, generated from a distribution $\mathcal{D}_i$. The goal is to learn a model $q_i : \mathbb{R}^d \to \mathcal{Y}$ maps input $\mathbf{x}_i \in \mathbb{R}^d$ to predicted labels $q_i(\mathbf{x}_i)$, such that $q_i(\mathbf{x}_i) \in \mathcal{Y}$ is as close as possible to the true label $y_i$. That is, to minimize the expected risk over $\mathcal{D}_i$ given by $f_i := \mathbb{E}_{(\mathbf{x}_i, y_i) \sim \mathcal{D}_i}[\ell(q(\mathbf{x}_i), y_i)]$, where $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ is a loss function that models the error between $q(\mathbf{x}_i)$ and $y_i$.

In FL paradigms, $M_i$'s are typically small, and if clients train a model using their local data separately, it may lead to poor generalization and performance due to limited data. To this end, agents collaborate by sharing model updates instead of raw data, leveraging the collective knowledge of all clients while preserving data privacy. Federated learning enables this cooperation through a central server, allowing clients to learn models using aggregated data while maintaining privacy. The standard FL approach thus learns a global model $q = q_1 = q_2 = \ldots = q_n$ by minimizing the average error $1/n \sum_{i=1}^n f_i(q)$ [1], [3]. Unlike the standard setting, we focus on personalized models and learning the shared representations.

We consider a model comprising a global representation, $\phi : \mathbb{R}^d \to \mathbb{R}^k$, and a local model for the $i$-th client, $h_i : \mathbb{R}^k \to \mathcal{Y}$. Specifically, the model for each client is defined as the composition of these two mappings, expressed as $q_i(\mathbf{x}_i) = (h_i \circ \phi)(\mathbf{x}_i)$. We assume that the dimensionality of the local model, $k \ll d$, meaning the number of parameters to be learned at the local client is much smaller than the overall model size. This captures the relationships among clients and enables model transferability in data-scarce scenarios, as only a small number of parameters are fine-tuned for new clients, while the global shared model $\phi$ remains fixed. The objective function of our learning process is thus formulated as

$$\min_{\phi \in \Phi} \frac{1}{n} \sum_{i=1}^n \min_{h_i \in \mathcal{H}} f_i(h_i \circ \phi), \qquad (1)$$

where $\Phi$ and $\mathcal{H}$ represent the classes of possible global and local representations, respectively.

### A. System Model

Our scheme comprises $n$ clients and a central server, with a subset of clients being malicious. The system model of the proposed MTRL scheme is presented in Figure 1. In each iteration round $t$ the proposed scheme follows these steps: (1) Each client fine-tunes and updates its local head $h_i^t$ using local data. (2) Clients update their estimates of the shared representation $\phi_i^t$. (3) Clients upload their updates to the server. (4) The server applies geometric median-based robust
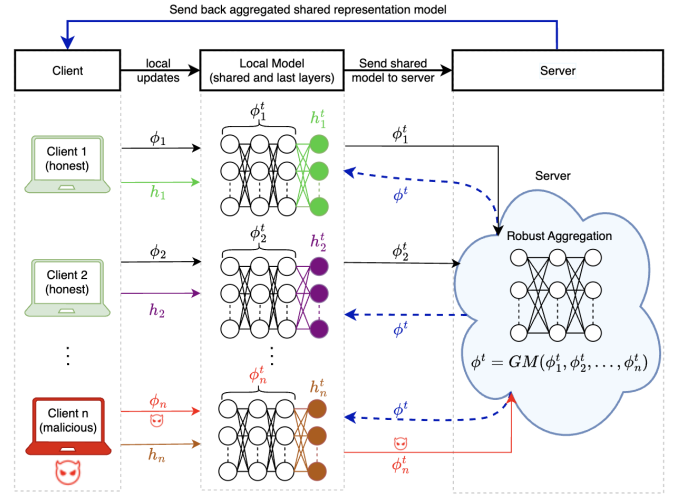


Fig. 1: The system model of proposed MTRL.

aggregation to compute the new shared representation $\phi^t$. (5) The updated representation is sent back to all clients.

**Malicious clients.** We consider two types of clients: Byzantine and honest. Byzantine clients manipulate the model sent to the central server by providing malicious model updates, aiming to corrupt the final aggregation. In contrast, honest clients contribute genuine updates, ensuring reliable and accurate training. We consider a white-box attack scenario, where the adversary has the capability to poison the global model.

**Robust server.** To minimize the impact of Byzantine clients, we apply the Geometric Median (GM) and Krum to the weights received from clients, which is then transmitted to the server. These techniques offer a more resilient aggregation by effectively reducing the influence of outlier updates. By using the median instead of the mean, GM reduces outlier distortion, ensuring the global model reflects honest client updates. Krum enhances resilience by selecting the client update that is most similar to its nearest neighbors, thereby mitigating the impact of malicious updates on the global model.

### B. Design Goals

In order to solve the MTRL problem under the above threats, our design should achieve the following goals.

**Byzantine robustness.** The proposed scheme learns an unbiased, personalized model, which minimizes the impact of malicious model updates.

**Representation learning and accuracy.** The proposed scheme must learn the shared model and the local client head and needs to ensure that the range of accuracy is acceptable.

**Transferability.** The proposed scheme should generalize to new, unseen clients by leveraging the learned shared model while training only the local client head, enabling efficient learning even with limited data.

## III. PROPOSED BYZANTINE RESILIENT MTRL ALGORITHM

We propose an alternating gradient descent (GD) approach, the pseudocode of which is presented in Algorithm 1. The goal of our algorithm is to jointly learn the global model $\phi$ and to learn the local client heads $h_i$, for $i \in [n]$, locally to minimize Eq. (1). The algorithm begins with the central server

initializing the shared model as $\phi^0$ randomly. Then, each round of the algorithm involves two main steps: (i) Client update via GD to update $h_i$ after fixing $\phi$ and (ii) Server update via GD at the client side to update $\phi_i$ by fixing $h_i$, followed by robust aggregation at the server. We elaborate each of the steps below.

**Client Update.** In each round $t \in [T]$, we select a fraction $\alpha \in (0, 1]$ of clients to participate in the client update. During this update, client $i$ performs $\tau_h$ local gradient steps to optimize its head model based on the current global representation $\phi^{t-1}$, which is shared by the server. Specifically, for each step $s = 1, \ldots, \tau_h$, client $i$ updates its head model as follows:

$$h_i^{t,s} = \text{GD}(f_i(h_i^{t,s-1}, \phi^{t-1}), h_i^{t,s-1}, \eta).$$

$\text{GD}(f, h, \eta)$ represents a generic gradient descent update of the variable $h$ using the gradient of the function $f$ with respect to $h$, and the step size $\eta$. Different gradient descent methods can be used, such as SGD or SGD with momentum. Each client makes many such local updates, i.e., $\tau_h$ is large, to optimize its local head based on the most recent global representation $\phi^{t-1}$ received from the server, leveraging its own local data. Next, the client performs $\tau_\phi$ local updates to refine its estimate of shared representation $\phi_i$ after fixing the local head estimate $h_i^{t,\tau_h}$ as, for $s = 1, \ldots, \tau_\phi$,

$$\phi_i^{t,s} = \text{GD}(f_i(h_i^{t,\tau_h}, \phi_i^{t,s-1}), \phi_i^{t,s-1}, \eta).$$

**Server Update.** Once the local updates for the head and shared representation are completed, the client participates in the server update by sending its locally updated representation $\phi_i^{t,\tau_\phi}$ to the server. Prior to sending the updates, Byzantine clients manipulate the local representations by adding noise,

$$\phi_i^{t,\tau_\phi} = \phi_i^{t,\tau_\phi} + \sigma \cdot \mathcal{N}(0, I),$$

where $\sigma$ is a scalar and $\mathcal{N}(0, I)$ denotes standard Gaussian noise. All clients $i \in [n]$ sends their local estimate of the representation, $\phi_i^{t,\tau_\phi}$, to the central server.

**Robust Aggregation.** To ensure robust aggregation, the server calculates Geometric Median for each layer of the shared model by finding a matrix $P$ by solving the following optimization

$$\text{GM}(\phi_1^{t,\tau_\phi}, \cdots, \phi_n^{t,\tau_\phi}) = \arg \min_{P \in \mathbb{R}^{d \times k}} \sum_{i=1}^{n} \|\phi_i^{t,\tau_\phi} - P\|_F,$$

where $\|\cdot\|_F$ denote Frobenius norm. The server sets the GM as the estimate of the representation $\phi^t$ and sends it to all clients. The three steps are repeated for $T$ iterations.

Further, we perform a Krum approach [14] detailed below.

$$\text{Krum}(\phi_1^{t,\tau_\phi}, \cdots, \phi_n^{t,\tau_\phi}) = \phi_{k*}^{t,\tau_\phi}, \text{ where}$$

$$k^* = \arg \min_{k \in [n]} \sum_{k \to k'} \|\phi_k^{t,\tau_\phi} - \phi_{k'}^{t,\tau_\phi}\|_F.$$

For $k \neq k'$, $k \to k'$ is the set of $n - f - 2$ nearest neighbors of $\phi_k^{t,\tau_\phi}$, where $f$ is the number of Byzantine clients. Thus $\phi_{k*}^{t,\tau_\phi}$ is the model with minimal summed distance to its neighbors.

**Remark III.1.** Krum requires to know the number of Byzantine clients in advance, and its computational complexity scales quadratically with the number of clients, resulting in significantly higher computational requirements as compared to GM.

---

**Algorithm 1:** Byzantine-Resilient Multi-Task Representation Learning (BR-MTRL) Algorithm

---

**Parameters:** Participation rate $\alpha$, step size $\eta$, number of local updates for the head $\tau_h$ and for the global representation $\tau_\phi$, number of communication rounds $T$.
**Initialize:** $\phi^0, h_1^{0,\tau_h}, \ldots, h_n^{0,\tau_h}$

1: **for** $t = 1, 2, \ldots, T$ **do**
2:     Server receives a batch of clients $I^t$ of size $rn$
3:     Server sends current representation $\phi^{t-1}$ to these clients
4:     **for** each client $i \in I^t$ **do**
5:         Client $i$ initializes $h_i^{t,0}$, $h_i^{t,0} \leftarrow h_i^{t-1,\tau_h}$
6:         **for** $s = 1$ to $\tau_h$ **do**   ▷ Client $i$ makes $\tau_h$ head updates
7:             $h_i^{t,s} \leftarrow \text{GD}(f_i(h_i^{t,s-1}, \phi^{t-1}), h_i^{t,s-1}, \eta)$
8:         **end for**
9:         Client $i$ initializes $\phi_i^{t,0}$ $\phi_i^{t,0} \leftarrow \phi^{t-1}$
10:        **for** $s = 1$ to $\tau_\phi$ **do**   ▷ Client $i$ makes $\tau_\phi$ updates to $\phi$
11:           $\phi_i^{t,s} \leftarrow \text{GD}(f_i(h_i^{t,\tau_h}, \phi_i^{t,s-1}), \phi_i^{t,s-1}, \eta)$
12:        **end for**
13:        **if** Client $i$ is malicious **then**
14:           $\phi_i^{t,\tau_\phi} = \phi_i^{t,\tau_\phi} + \sigma \cdot \mathcal{N}(0, I)$
15:        **end if**
16:        Client $i$ sends updated representation $\phi_i^{t,\tau_\phi}$ to server
17:     **end for**
18:     **for** each client $i \notin I^t$ **do**
19:        Set $h_i^{t,\tau_h} \leftarrow h_i^{t-1,\tau_h}$
20:     **end for**
21:     Server computes the new representation as $\phi^t \leftarrow \text{GM}(\phi_1^{t,\tau_\phi}, \cdots, \phi_n^{t,\tau_\phi})$ or $\text{Krum}(\phi_1^{t,\tau_\phi}, \cdots, \phi_n^{t,\tau_\phi})$
22: **end for**

---

## IV. EXPERIMENTAL ANALYSIS

### A. Experimental Setup

**Datasets.** We evaluated our algorithm on two real-world datasets, CIFAR-10 and Federated Extended MNIST (FEMNIST). CIFAR-10 consists of 60,000 $32 \times 32$ color images in 10 classes, with 50,000 training images and 10,000 test images. FEMNIST consists of 145,600 $28 \times 28$ pixel images of handwritten characters, with 124,800 training images and 20,800 test images. It incorporates both uppercase and lowercase letters, and the data is partitioned based on the individual writers of characters, offering a broader classification task.

**Model.** We use the AlexNet architecture, which is composed of five convolutional blocks followed by a linear layer for the shared representation and an additional linear layer for the personalized client heads. Each convolutional block includes a convolutional layer with a stride and padding of 1, a ReLU activation function, and a max-pooling layer with a kernel size of 2. We optimize the model using stochastic gradient descent (SGD) with a learning rate $\eta = 0.01$ and momentum $\beta = 0.9$.

**Hyperparameters.** We initialized all models with random parameters. For CIFAR-10, training was conducted for $T = 100$ communication rounds, while FEMNIST models were trained for $T = 200$ rounds. In each round, 20% of clients ($\alpha = 0.2$) were selected to perform local updates and receive the aggregated model from the server. Each selected client

executed 10 local epochs of SGD with momentum for head training and one epoch for updating the global representation.

Data was distributed in a non-iid manner, with each client assigned different classes. We considered 100 and 1000 clients for CIFAR-10 and 150 clients for FEMNIST. Among them, 20 out of 100, 100 out of 1000, and 50 out of 150 clients were designated as malicious. By varying the proportion of Byzantine clients, we analyzed their impact on learned model. **Benchmark Algorithms.** We evaluated our algorithm with benchmark approaches, including FedAvg [1], FedPer [17], and FedRep [3]. FedPer and FedRep are similar approaches that learn a shared representation with personalized heads, but FedPer updates both simultaneously, maintaining an equal number of updates per local round. In contrast, FedAvg learns a single global model without personalization. As a best-case benchmark, we considered these methods without malicious clients. Additionally, we analyzed two variants: (i) introducing malicious clients (Byzantine FedPer/FedRep/FedAvg), (ii) incorporating geometric median-based aggregation (FedAvg+GM, FedPer+GM), and (iii) incorporating Krum-based aggregation (FedAvg+Krum, FedPer+Krum). For meta-test, we consider an additional baseline, Naive, which trains new clients independently without leveraging the MTRL framework.

**Performance metrics.** We evaluate the performance of our algorithm by calculating the average classification accuracy across all clients in the presence of malicious clients.

### B. Amazon Web Services (AWS)-Based Simulation Testbed

We deployed our proposed approach on AWS using two experimental setups. In the first setup, called *sequential AWS*, a single EC2 (Elastic Compute Cloud) instance is used for the server, while another EC2 instance hosts all the clients. In the second setup, *federated AWS*, each client is allocated a separate EC2 instance. In both setups, during each communication round, the server retrieves the shared model from S3 (Simple Storage Service) cloud storage and sends it along with the selected clients' IDs to the clients. In the sequential AWS setup, clients update their respective heads sequentially and send the updated model back to the server. In the federated AWS setup, the server sends the shared model and client IDs to the respective EC2 instances, where each client performs local learning, updates its head, and sends the updated model back to the server in parallel. Thus, sequential AWS mimics a standard simulation of an FL algorithm, while federated AWS represents a genuinely federated simulation platform. We compared the communication and computation times for both setups and evaluated their performance.

We considered a synchronous broadcast model, where the server waits until it receives the updated models from all selected clients before aggregation. After updating its head, the client stores the new model in its S3 bucket. We deployed the system on an AWS EC2 t2.large instance (2 vCPUs, 8 GiB RAM, 64 GB EBS storage). Communication between the server and clients is handled via Flask, which enables HTTP requests for model updates. The server sends requests to clients, and Flask processes like receiving the shared model estimates and sending back updated models. This setup ensures efficient communication and aggregation of client models. We present

| Process | Federated AWS | Sequential AWS |
|---|---|---|
| Sending from server to clients | 0.31 | 0.13 |
| Training clients' model | 15.41 | 92.79 |
| Sending from clients to server | 4.68 | 1.20 |
| Average time cost per round | 26.50 | 121.75 |

TABLE I: Comparison of average runtime (in milliseconds) for federated AWS and sequential AWS implementations.

| | CIFAR-10 | | | FEMNIST |
|---|---|---|---|---|
| (# Clients n, # Classes per Client S, # Byzantine Clients) | (100,2,20) | (100,5,20) | (1000,2,100) | (150,3,50) |
| FedAvg [1] | 40.08 | 62.62 | 36.66 | 50.42 |
| FedPer[17] | 90.81 | 82.29 | 79.23 | 76.82 |
| FedRep [3] | 87.70 | 80.17 | 80.79 | 76.63 |
| Byzantine FedAvg | 13.44 | 16.45 | 22.33 | 17.80 |
| Byzantine FedPer | 52.72 | 31.11 | 57.26 | 35.21 |
| Byzantine FedRep | 52.08 | 32.69 | 58.17 | 35.25 |
| FedAvg+GM | 36.38 | 41.77 | 31.84 | 38.43 |
| FedAvg+Krum | 31.52 | 42.25 | 30.08 | 38.20 |
| FedPer+GM | 72.37 | 50.08 | 63.03 | 52.47 |
| FedPer+Krum | 71.90 | 61.21 | 68.17 | 53.12 |
| BR-MTRL+GM | 72.94 | 58.55 | 75.92 | 55.28 |
| BR-MTRL + Krum | 80.25 | 61.80 | 72.76 | 54.07 |

TABLE II: Average test accuracy across all clients with some clients under byzantine attacks for different algorithms.

the running time comparison in Table I. Training time in Table I represents the total time for all clients to train their models. The results show that federated AWS achieves significantly faster communication rounds than sequential AWS due to parallel client training, which greatly reduces overall training time. Although federated AWS incurs slightly higher communication overhead between the server and clients, the overall efficiency gains make FL an effective approach.

### C. Experiment Results

**Test accuracy.** As shown in Figures 2a and 2b and Table II, the FedRep and FedPer algorithms consistently achieved higher accuracy than FedAvg across various experiments, irrespective of the number of clients or classes per client, demonstrating the effectiveness of personalized FL. However, the performance of these three baselines (without any malicious clients) significantly declined under Byzantine attacks. This decline validates that (personalized) FL frameworks susceptible to Byzantine attacks will lead to unreliable models with poor performance. Increasing the number of classes per client from two to five resulted in a substantial accuracy drop of $21.61\%$. This is because Byzantine clients can negatively impact the global model's performance across a greater number of classes. Furthermore, reducing the proportion of Byzantine clients relative to the total number of clients improves performance, as seen in the comparison between experiments with $(100, 2, 20)$ and $(1000, 2, 100)$, i.e., $20\%, 10\%$ malicious clients, respectively, where the latter exhibits greater resilience.

Our experimental results further demonstrate that the proposed approaches (BR-MTRL+GM and BR-MTRL+Krum) achieves a higher accuracy than both the GM and Krum-adapted versions of FedPer. This is because FedPer updates the server model at the same frequency as client updates, making it more susceptible to Byzantine clients, which can adversely affect learning accuracy. In contrast, our approach updates

(a) CIFAR-10    (b) FEMNIST    (c) Meta test on CIFAR-10    (d) Meta test on FEMNIST
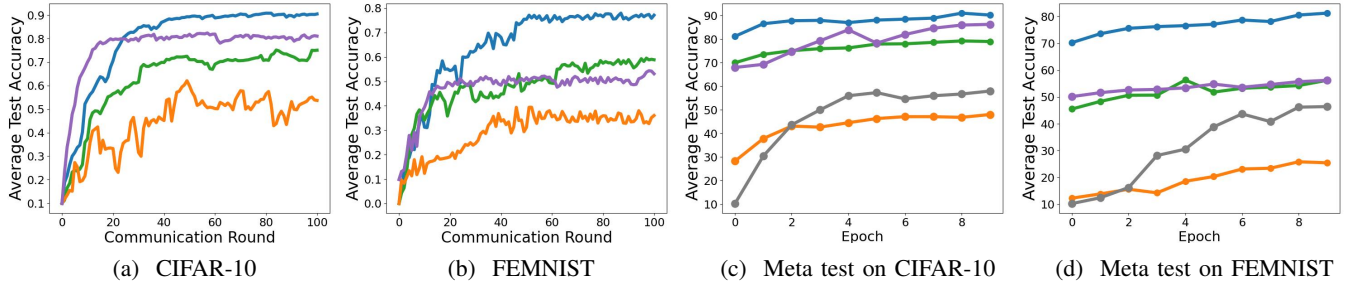
Fig. 2: ■ FedRep, ■ Byzantine FedRep, ■ BR-MTRL+GM (proposed), ■ BR-MTRL+Krum (proposed) ■ Naive. **Average test accuracy across all clients:** Figure 2a shows the results for CIFAR-10 with 100 clients with 20 malicious clients. Figure 2b presents the results for FEMNIST with 150 clients with 20 malicious clients. We compared our proposed BR-MTRL with the FedRep algorithm under benign conditions (i.e., no malicious client) and FedRep subjected to Byzantine attacks (Byzantine FedRep). **Meta test for new clients:** Figures 2c and 2d present the average test accuracy across 10 new clients (CIFAR-10 and FEMNIST) after fine-tuning the local heads for 10 epochs using the shared model. We first learned the shared model using source clients and then used the shared model to optimize the local head of new clients.

the shared model only once after updating the client heads, leading to improved overall performance. All experimental results consistently showcase the improved performance of proposed approach. Krum outperforms GM in experiments with 100 clients; however, has a lower accuracy when the number of clients increases. Krum requires to know the number of malicious clients and computational complexity is higher.

**Results for Meta-test.** We tested the transferability of the proposed approach. We first learn the shared model $\phi$ using the source clients and then utilize this model to train the local heads of new, previously unseen clients. For each new client, we fine-tune the client head for 10 epochs, using the transferred shared model. For CIFAR-10, we randomly selected 10 clients as new clients and trained the shared model using the remaining 90 clients, including 20 Byzantine clients. Similarly, for FEMNIST, we selected 10 new clients and trained the shared model with the remaining 140 clients, of which 10 were Byzantine. We evaluated the average accuracy for four approaches. To highlight the benefits of MTRL for transfer learning, we considered a Naive setting, where each client trains their model independently, resulting in poor accuracy. Figures 2c and 2d illustrate the results: FedRep, without Byzantine attacks, demonstrates effective transferability of the shared model to new clients. In contrast, Byzantine MTRL, with malicious clients, exhibits significantly lower performance, highlighting the vulnerability of FL systems. Finally, BR-MTRL, our proposed approach, improves learning accuracy by employing GM and Krum aggregation, validating its effectiveness.

## V. CONCLUSION

In this paper, we proposed a Byzantine-resilient federated MTRL framework to learn robust shared representations across *related yet different* tasks. Our approach employs an alternating gradient descent strategy, where the shared model is updated while keeping the local head fixed, and vice versa. To ensure Byzantine resilience, we incorporated geometric median and Krum-based aggregation at the server, mitigating the impact of malicious clients sending corrupted updates. We implemented our federated learning approach on AWS, simulating a real-world federated setting. Extensive experiments on CIFAR-10 and FEMNIST validate the effectiveness of our approach and its transferability to new clients with limited data.

## REFERENCES

[1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*, 2017, pp. 1273–1282.

[2] Y. Jiang, J. Konečnỳ, K. Rush, and S. Kannan, "Improving federated learning personalization via model agnostic meta learning," *arXiv:1909.12488*.

[3] L. Collins, H. Hassani, A. Mokhtari, and S. Shakkottai, "Exploiting shared representations for personalized federated learning," in *International conference on machine learning*. PMLR, 2021, pp. 2089–2099.

[4] M. Li, X. Tang, T. Zhang, G. Liu, and Y. Weng, "Byzantine-robust federated learning on non-iid data via inversing artificial gradients," in *IEEE Global Communications Conference*, 2024, pp. 355–360.

[5] Y. Deng, M. M. Kamani, and M. Mahdavi, "Adaptive personalized federated learning," *arXiv:2003.13461*, 2020.

[6] A. Fallah, A. Mokhtari, and A. Ozdaglar, "Personalized federated learning: A meta-learning approach," *arXiv:2002.07948*, 2020.

[7] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.

[8] Y. Xie, P. Jin, M. Gong, C. Zhang, and B. Yu, "Multi-task network representation learning," *Frontiers in Neuroscience*, vol. 14, p. 1, 2020.

[9] Y. Yang and T. Hospedales, "Deep multi-task representation learning: A tensor factorisation approach," *arXiv:1605.06391*, 2016.

[10] S. S. Du, W. Hu, S. M. Kakade, J. D. Lee, and Q. Lei, "Few-shot learning via learning the representation, provably," in *International Conference on Learning Representations (ICLR)*, 2021.

[11] J. Lin, S. Moothedath, and N. Vaswani, "Fast and sample efficient multi-task representation learning in stochastic contextual bandits," in *International Conference on Machine Learning*, 2024.

[12] J. Lin and S. Moothedath, "Fast and sample-efficient relevance-based multi-task representation learning," *IEEE Control Systems Letters*, 2024.

[13] A. Gouissem, K. Abualsaud, E. Yaacoub, T. Khattab, and M. Guizani, "Collaborative byzantine resilient federated learning," *IEEE Internet of Things Journal*, vol. 10, no. 18, pp. 15 887–15 899, 2023.

[14] F. Lin, Q. Ling, and Z. Xiong, "Byzantine-resilient distributed large-scale matrix completion," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 8167–8171.

[15] J. Zhang, X. He, Y. Huang, and Q. Ling, "Byzantine-robust and communication-efficient personalized federated learning," *IEEE Transactions on Signal Processing*, 2024.

[16] P. P. Liang, T. Liu, L. Ziyin, N. B. Allen, R. P. Auerbach, D. Brent, R. Salakhutdinov, and L.-P. Morency, "Think locally, act globally: Federated learning with local and global representations," *arXiv:2001.01523*.

[17] M. G. Arivazhagan, V. Aggarwal, A. K. Singh, and S. Choudhary, "Federated learning with personalization layers," *arXiv:1912.00818*, 2019.

[18] S. Li, E. Ngai, and T. Voigt, "Byzantine-robust aggregation in federated learning empowered industrial iot," *IEEE Transactions on Industrial Informatics*, vol. 19, no. 2, pp. 1165–1175, 2021.

[19] X. Wang, H. Zhang, A. Bilal, H. Long, and X. Liu, "Wgm-dsaga: Federated learning strategies with byzantine robustness based on weighted geometric median," *Electronics*, vol. 12, no. 5, p. 1190, 2023.

[20] X. He, J. Zhang, and Q. Ling, "Byzantine-robust and communication-efficient personalized federated learning," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2023, pp. 1–5.