

# LoRA-FL: A Low-Rank Adversarial Attack for Compromising Group Fairness in Federated Learning

Anonymous authors  
Paper under double-blind review

## Abstract

Federated Learning (FL) enables collaborative model training without requiring participants to share raw data, and is increasingly deployed in regulated domains such as healthcare, finance, and large-scale personalization. FL offers privacy and governance benefits, it can also obscure fairness risks: heterogeneity in client data distributions may lead to models that systematically disadvantage minority groups. Ensuring fairness in such settings is not only an ethical concern but also a regulatory requirement under frameworks such as GDPR and anti-discrimination law. Existing adversarial manipulations in FL, such as noise injection or scaling attacks, typically degrade predictive performance or are mitigated by robust aggregation rules (e.g., KRUM or FLAME), limiting their practical relevance. In this work, we introduce LoRA-FL, a stealthy fairness attack that leverages low-rank adapters to inject group-level bias while preserving accuracy. By constraining adversarial updates to a compact subspace that aligns with benign client variation, LoRA-FL evades both standard and robust aggregators, even under heterogeneous (non-IID) data distributions. We provide empirical results, across widely used fairness benchmarks, including tabular datasets such as Adult and Bank. With LoRA-FL as few as 10–20% adversarial clients can increase violations of demographic parity and equalized odds by over 40%, while maintaining comparable predictive performance.

## 1 Introduction

Federated Learning (FL) (McMahan et al., 2017) enables multiple data holders to collaboratively train deep learning models without sharing raw data, thereby respecting data ownership and supporting regulatory compliance (Parliament & of the European Union, 2016). FL has been widely adopted in high-stakes domains such as healthcare (Sheller et al., 2020; Xu et al., 2021), finance (Yang et al., 2019), and personalized language modeling (Hard et al., 2018), where centralized data collection is either infeasible or undesirable. As FL systems increasingly influence decisions about individuals and communities, concerns about fairness have become central to their responsible deployment (Barocas & Selbst, 2016; Barocas et al., 2023).

**Real-world Instances (Salazar et al., 2024).** In federated credit scoring or risk assessment systems, banks and financial institutions collaboratively train models while retaining control over proprietary customer data. Such institutions may have economic incentives to preserve lending practices that disadvantage certain demographic groups, so long as overall predictive performance remains strong. Similarly, large digital platforms participating in federated training for ranking, recommendation, or personalization systems may prioritize revenue or engagement objectives that conflict with fairness constraints, resulting in models that encode systematic representational biases.

From a standard learning viewpoint, success in FL is typically measured by the convergence of the global loss and overall predictive accuracy. In the idealized IID setting—where each client’s local data are drawn independently from a common distribution—such metrics often suffice to characterize model performance. In contrast, in realistic deployments, FL operates under non-IID data distributions, where the joint distribution of features and labels varies significantly across clients. In decision-making systems such as credit approval, however, model behavior across demographic groups, such as gender, age, or race, constitutes a

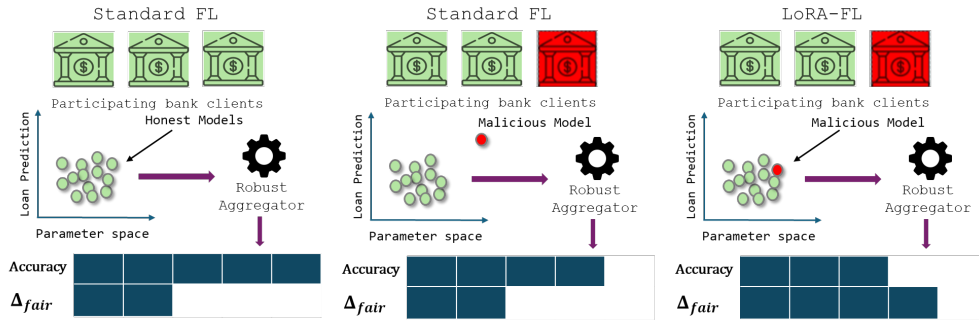


Figure 1: LoRA-FL constrains adversarial updates to a low-rank subspace, making malicious updates appear similar to benign ones and allowing robust aggregators to accept them. Consequently, overall accuracy remains stable, while group-wise bias increases.

first-order concern (Chouldechova, 2017; Hardt et al., 2016a). A converged FL model may achieve strong aggregate accuracy while approving loans for one group at substantially higher rates than another, even among applicants with comparable financial profiles (Angwin et al., 2016; Fabris et al., 2025). In practice, such fairness degradation is rarely treated as anomalous; instead, it is often viewed as a natural byproduct of maximizing accuracy under client heterogeneity, particularly in non-IID FL settings. This tension between performance and bias has motivated a formal approach in quantifying and enforcing fairness constraints.

**Fair Federated Learning.** A substantial body of work has formalized group fairness through criteria such as Demographic Parity (Chouldechova, 2017), Equalized Odds (Hardt et al., 2016a), and Equal Opportunity (Hardt et al., 2016a), and studied their enforcement and evaluation in machine learning systems (Bilal Zafar et al., 2015; Agarwal et al., 2018; Madras et al., 2018; Friedler et al., 2019; Pessach & Shmueli, 2022; Rabonato & Berton, 2025). In FL, however, achieving such fairness guarantees is fundamentally more challenging. Even small and individually insignificant disparities at the client level can accumulate and amplify into significant group-level unfairness in the global model (Chang & Shokri; Wang et al.; Salazar et al., 2024). The fairness challenge becomes substantially more complex when clients are *malicious*. In FL, participants retain control over their local training procedures and may intentionally deviate from the intended optimization process by manipulating local updates, causing adversarial attacks (Bhagoji et al., 2019; Shejwalkar & Houmansadr, 2021; Xia et al., 2023). Such attacks can further worsen disparities across sensitive groups (Solans et al., 2020).

**Adversarial Attacks in FL.** Prior work on adversarial behavior in FL has largely examined model poisoning strategies in which malicious participants manipulate local updates to influence the global model either by degrading overall predictive accuracy or by inducing targeted misclassification that benefits specific objectives (Cao et al., 2020; So et al., 2020). Such attacks are particularly relevant in settings where participants have incentives to shape model outcomes while remaining indistinguishable from benign clients. To defend against malicious participants, a substantial literature has proposed Byzantine-resilient aggregation rules that suppress updates deviating from the majority (Blanchard et al., 2017a; Yin et al., 2018; Nguyen et al., 2022).

**Robust FL blind to fairness manipulation.** The proposed robust aggregators, both in IID and non-IID, are designed to detect *accuracy-disruptive* behavior. Classical defenses such as KRUM (Blanchard et al., 2017a) and trimmed-mean (TM) (Yin et al., 2018) operate under IID assumptions by selecting or averaging updates that are closest in parameter space, while more recent mechanisms such as FLAME (Nguyen et al., 2022) explicitly account for non-IID heterogeneity by clustering or reweighting updates before aggregation. Updates that remain statistically consistent with benign optimization, yet systematically influence fairness-relevant directions in parameter space, do not violate the criteria enforced by these aggregators. In fact, under non-IID assumptions, the natural variance of benign updates further enlarges the set of directions that such defenses must tolerate, creating room for fairness manipulation that preserves both convergence and accuracy.

**Our key insight.** The above observations suggest a structural gap between how robustness is enforced in FL and how fairness violations manifest. Robust aggregators are designed to detect deviations that disrupt accuracy, implicitly treating all directions in parameter space as equally relevant. In contrast, fairness-relevant behavior often depends on a restricted subspace tied to features correlated with sensitive attributes, whose influence on global performance may be subtle. We hypothesize that confining adversarial perturbations to this subspace allows fairness manipulations to remain hidden and preserves accuracy.

**Our Approach and Contributions.** Motivated by this hypothesis, we propose LoRA-FL that leverages low-rank parameterizations to restrict adversarial updates to a compact subspace (Figure 1). Using LoRA-FL clients jointly optimize for task performance and fairness degradation with minimal additional computational overhead. As we demonstrate empirically, this design allows adversarial behavior to blend seamlessly with benign updates, even under robust aggregation and non-IID data heterogeneity. Our contributions are as follows:

1. We introduce LoRA-FL, a fairness-attack strategy that augments standard FL training by leveraging low-rank adapters to inject fairness-oriented bias into federated models without compromising predictive performance (Algorithm 1). Operating within a low-dimensional subspace, LoRA-FL allows adversarial clients to craft updates that remain statistically consistent with benign ones (Section 4.2).
2. We demonstrate the effectiveness of how LoRA-FL reliably degrades Demographic Parity, Equalized Odds, and Equal Opportunity on standard benchmarks (Adult Dua & Graff (2017), Bank Moro et al. (2014), and Dutch liobaite et al. (2011)). With only 10–20% adversarial clients, LoRA-FL reduces fairness metrics (DP, EO, and EOpp) by over 40%, while preserving high accuracy under both IID (Table 1) and non-IID (Table 2), robust aggregators such as KRUM (Blanchard et al., 2017a) and FLAME (Nguyen et al., 2022) (Section 5). An ablation study further confirms LoRA-FL’s robustness to varying numbers of agents. Additionally, increasing the adapter rank makes adversarial updates more detectable, highlighting the importance of low-rank constraints.
3. We analyze why low-rank adversarial updates evade detection by studying their alignment with dominant gradient subspaces and aggregation geometry (Section 6). A detailed interpretability study shows how low-rank adapters are effective. Lower-rank perturbations closely align with benign update distributions, thereby evading detection. Moreover, these adapters disrupt internal neuron-level representations, thereby systematically skewing predictions across demographic groups. Finally, we show that even higher-rank adapters concentrate changes along a few principal directions – explaining why low-rank updates suffice to achieve the attack’s effect while remaining covert.

Together, our results show that fairness manipulation can be embedded directly into the optimization dynamics of FL, remaining invisible to accuracy-centric robust aggregators. This exposes a fundamental gap between robustness and fairness in decentralized optimization.

## 2 Related Work

Data poisoning attacks corrupt local training datasets by injecting false or misleading examples, thereby influencing the learned model (Biggio et al., 2012; Rubinstein et al., 2009; Li et al., 2016; Fang et al., 2020; Dai & Li, 2023). Prior work has shown that algorithmic fairness can be compromised through data poisoning in centralized learning. Solans et al. (2020) propose a gradient-based poisoning attack that induces disparities across demographic groups, degrading fairness metrics such as Demographic Parity and Equalized Odds.

In contrast, local model poisoning is a much more severe threat in FL. Byzantine-robust stochastic gradient descent and aggregation methods include distance-based and coordinate-wise defenses (Blanchard et al., 2017a; Yin et al., 2018; Nguyen et al., 2022). Following the robustness definition of Li et al. (2021), these approaches aim to preserve convergence and predictive performance in the presence of adversarial clients. Studies have demonstrated that targeted manipulation of training data can exacerbate bias while maintaining high predictive accuracy (Van et al., 2022; Mehrabi et al., 2021). These results establish that fairness can be manipulated independently of accuracy, but they focus on instance-level attacks in centralized settings. Unlike these, we propose an attack for the FL setting that introduces a substantially different threat model.

In FL, adversaries can control entire clients, influence local optimization procedures, and repeatedly affect the global model through aggregated updates. To mitigate fairness concerns in this setting, prior work has proposed fairness-aware aggregation and optimization mechanisms, such as FairFed (Ezzeldin et al., 2023) and GIFAIR-FL (Yue et al., 2023), which explicitly incorporate fairness objectives into federated training. Unlike our work, these approaches typically assume benign participants and do not consider adversarial manipulation of local updates.

Only a limited number of works study adversarial attacks on fairness in FL. PFAttack (Gao et al., 2024) demonstrates that fairness-aware aggregators can be subverted without degrading accuracy, but the attack is tailored to demographic parity and requires careful empirical tuning to remain stealthy. He et al. (2025) propose robust aggregation while maintaining *performance fairness*. This notion of fairness differs fundamentally from group fairness definitions we consider, which focus on disparities across individuals or demographic groups. EAB-FL (Meerza & Liu, 2024) proposes a model poisoning strategy that increases bias while preserving utility, but is ineffective against robust aggregation rules. DiTTO (Li et al., 2021) explores targeted fairness manipulation through client-level optimization, yet remains sensitive to aggregation mechanisms and data heterogeneity. In contrast, our proposed LoRA-FL attack is model-agnostic and effective against robust aggregators as it preserves accuracy.

### 3 Background

We consider a standard classification setting in which each data point is a tuple  $(x, y, a)$  drawn from an unknown distribution  $\mathcal{D}$  over  $\mathcal{X} \times \mathcal{Y} \times \mathcal{A}$ . The instance space is  $\mathcal{X} \subseteq \mathbb{R}^d$ , where  $x \in \mathcal{X}$  denotes a  $d$ -dimensional feature vector. The label space is  $\mathcal{Y} = [C]$  for a  $C$ -class classification problem, and  $\mathcal{A}$  denotes the space of sensitive attributes. Each  $a \in \mathcal{A}$  encodes a sensitive group membership (e.g., gender, age, caste), and is observed together with the input-label pair.

#### 3.1 Federated Learning (FL)

In a typical Federated Learning (FL) setup, we consider (i) a set of agents  $[K] = \{1, \dots, K\}$ ; and (ii) a central aggregator. Each agent  $k$  holds a local dataset  $\mathcal{D}^{(k)} = \{(x_i, y_i, a_i)\}_{i=1}^{n_k}$  formed by sampling from  $\mathcal{D}$  according to a client-specific distribution. To model statistical heterogeneity across agents, we assume that class proportions at each agent are drawn from a Dirichlet distribution parameterized by a concentration parameter  $\rho > 0$ . Larger values of  $\rho$  correspond to more homogeneous data distributions across agents, with the limit  $\rho \rightarrow \infty$  recovering the i.i.d. setting. Conversely, smaller values of  $\rho$  induce increasing non-i.i.d. behavior, resulting in highly skewed local datasets dominated by a subset of classes or sensitive groups. At the outset, the aggregator initializes global model parameters, denoted by  $\Theta_0$ . In each training round  $t$ , every agent  $k$  updates its local model parameters  $\theta_{k,t}$  using its dataset  $\mathcal{D}^{(k)}$  and sends the updated parameters to the aggregator for aggregation.

At each communication round  $t$ , the central aggregator computes the global model parameters  $\Theta_t$  by aggregating local updates  $\{\theta_{k,t}\}_{k=1}^K$  according to an aggregation rule **Agg**:

$$\Theta_t = \text{Agg}(\{\theta_{k,t}\}_{k=1}^K, \{w_k\}_{k=1}^K),$$

where  $\theta_{k,t}$  denotes the local model parameters of agent  $k$ , and  $w_k$  is the aggregation weight assigned to agent  $k$ . A widely used instantiation is the *Federated Averaging* (**FedAvg**) algorithm McMahan et al. (2017), where the weights are proportional to the number of data points held by each agent. The global model is updated as:

$$\Theta_t^{\text{FedAvg}} = \sum_{k \in \mathcal{S}_t} \frac{|\mathcal{X}^{(k)}|}{\sum_{j \in \mathcal{S}_t} |\mathcal{X}^{(j)}|} \cdot \theta_{k,t} \quad (1)$$

Here,  $\mathcal{S}_t \subseteq [K]$  represents the (random) set of participating agents in round  $t$ , and  $|\mathcal{X}^{(k)}|$  denotes the size of agent  $k$ 's dataset. The weights,  $w_k = \frac{|\mathcal{X}^{(k)}|}{\sum_{j \in \mathcal{S}_t} |\mathcal{X}^{(j)}|}$  for each agent  $k$ . This weighting ensures that agents with larger datasets have a proportionally greater influence on the global model. The process repeats over multiple rounds until convergence, resulting in a final global model  $\Theta^*$  at round  $T$ .

### 3.1.1 Robust Aggregators

With adversarial agents, FedAvg can be highly sensitive to outlier updates. This sensitivity to outliers has motivated the development of robust aggregation rules such as KRUM, TM and FLAME, which aim to limit the influence of anomalous or adversarial.

*m*-KRUM Blanchard et al. (2017b). *m*-KRUM selects updates that are most consistent with the majority. Given  $K$  agents in any FL round, each providing an update  $\theta_k$  for  $k = 1, 2, \dots, K$ , and assuming up to  $\tilde{q} = \lfloor qK \rfloor$  adversarial agents, *m*-KRUM proceeds as follows. First, it computes pairwise distances  $d_{i,j} = \|\theta_i - \theta_j\|_2$  for all  $i \neq j$ . For each agent  $i$ , it then selects the set  $\mathcal{N}_i$  of the  $K - \tilde{q} - 2$  updates closest to  $\theta_i$  and computes the outlier score  $s_i = \sum_{j \in \mathcal{N}_i} d_{i,j}^2$ . Next, it chooses the  $m$  updates with the smallest scores – ensuring  $m \geq K - \tilde{q} -$  and aggregates them via a weighted average proportional to each agents sample size. By excluding highscore (outlying) updates, *m*-KRUM effectively filters adversarial contributions, yielding a more robust global model  $\Theta$ . The formal aggregation is summarized in Appendix B.1

Trimmed-Mean (TM) (Yin et al., 2018). *f*-Trimmed-Mean assumes up to  $f$  of these may be adversarial updates and the  $f$  largest and  $f$  smallest values among  $\{\theta_k\}_{k=1}^K$  are removed. Specifically, for each parameter  $\theta$ , the aggregator collects all corresponding agent values, sorts them element-wise, discards the extreme  $2f$  values, and takes the mean of the remaining  $K - 2f$  entries to compute the aggregated parameter  $\Theta^{\text{TM}}$ . This process is repeated for all parameters in the model. The formal algorithm is available in Appendix B.2.

FLAME (Nguyen et al., 2022). FLAME computes pairwise cosine similarities between updates and clusters them to identify a majority benign group. Updates outside this cluster are discarded or down-weighted, and the benign updates are averaged with adaptive Gaussian noise calibrated to their variance, providing robustness under adversarial behavior and non-IID data. The formal procedure is summarized in Appendix B.3.

## 3.2 Group Fairness

Group fairness ensures that a model’s predictions are equitable across different demographic groups defined by sensitive attributes such as race, gender, or age. We consider a parameterized classifier  $f_\Theta : \mathcal{X} \rightarrow \mathcal{Y}$ , where  $\theta$  are the (learned) model parameters and  $\mathcal{X}$  is the input space. For each input sample  $x \in \mathcal{X}$ , the predicted label is given by  $\hat{y} = f_\Theta(x)$ . We denote by  $\hat{Y}$  the set of predicted labels for a dataset, i.e.,  $\hat{Y} = \{f_\Theta(x_i)\}_{i=1}^n$  for samples  $\{x_i\}_{i=1}^n \subseteq \mathcal{X}$ . Formally, we focus on the following notions of fairness.

**Demographic Parity (DP) Dwork et al. (2012).** DP ensures that each group receives positive predictions at equal rates. In our running example of loan approval, DP looks only at the overall rate of approvals: it requires that  $f_\Theta$  approves individuals at equal rates across groups  $a$  &  $b$ , regardless of actual qualification.

**Definition 1** (Demographic Parity (DP) Dwork et al. (2012)). *A classifier  $f_\Theta$  satisfies DP if the probability of a positive prediction is the same across all groups, regardless of the actual outcomes. Formally, for all groups  $a, b \in \mathcal{A}$ :*

$$\Pr(\hat{Y} = 1 \mid \mathcal{A} = a) = \Pr(\hat{Y} = 1 \mid \mathcal{A} = b)$$

Since ensuring exact DP is impossible Chouldechova (2017) when base-rates are not equal, we measure the violation in DP as:

$$\Delta_{DP} := |\Pr(\hat{Y} = 1 \mid \mathcal{A} = a) - \Pr(\hat{Y} = 1 \mid \mathcal{A} = b)| \quad (2)$$

**Equalized Odds (EO) Hardt et al. (2016b).** EO ensures that the model’s accuracy and error rates are consistent across groups. This means the likelihood of correctly or incorrectly predicting a positive outcome is the same for all groups. In our loan approval example, EO ensures that qualified and unqualified individuals are treated similarly across groups  $a$  and  $b$ , i.e., both true positive and false positive rates align.

**Definition 2** (Equalized Odds (EO) Hardt et al. (2016b)). *A classifier  $f_\Theta$  satisfies EO if all groups have equal true positive rates (TPR) and false positive rates (FPR). Formally, for all groups  $a, b \in \mathcal{A}$ :*

$$\begin{aligned}\Pr(\hat{Y} = 1 \mid \mathcal{Y} = 1, \mathcal{A} = a) &= \Pr(\hat{Y} = 1 \mid \mathcal{Y} = 1, \mathcal{A} = b) \\ \Pr(\hat{Y} = 1 \mid \mathcal{Y} = 0, \mathcal{A} = a) &= \Pr(\hat{Y} = 1 \mid \mathcal{Y} = 0, \mathcal{A} = b)\end{aligned}$$

We define the violation in EO as:

$$\Delta_{EO} := \max\{\Delta_{\text{TPR}}, \Delta_{\text{FPR}}\}, \text{ where} \quad (3)$$

$$\begin{aligned}\Delta_{\text{TPR}} &:= |\Pr(\hat{Y} = 1 \mid \mathcal{Y} = 1, \mathcal{A} = a) - \Pr(\hat{Y} = 1 \mid \mathcal{Y} = 1, \mathcal{A} = b)| \\ \Delta_{\text{FPR}} &:= |\Pr(\hat{Y} = 1 \mid \mathcal{Y} = 0, \mathcal{A} = a) - \Pr(\hat{Y} = 1 \mid \mathcal{Y} = 0, \mathcal{A} = b)|\end{aligned}$$

**Equal Opportunity (EOpp) Hardt et al. (2016b).** EOpp focuses on ensuring that qualified individuals (i.e., those with  $\mathcal{Y} = 1$ ) have an equal chance of being correctly identified by the model, regardless of their group membership. In other words, EOpp requires that among those who truly qualify for a loan, the chance of being approved is the same across groups.

**Definition 3** (Equal Opportunity (EOpp) Hardt et al. (2016b)). *A classifier  $f_\theta$  satisfies EOpp if it has equal true positive rates (TPR) across all groups. Formally, for all groups  $a, b \in \mathcal{A}$ :*

$$\Pr(\hat{Y} = 1 \mid Y = 1, \mathcal{A} = a) = \Pr(\hat{Y} = 1 \mid Y = 1, \mathcal{A} = b)$$

The violation in EOpp is straightforward from Definition 3, and implies that EOpp is a weaker fairness notion than EO.

$$\Delta_{EOpp} := \Delta_{\text{TPR}} \quad (4)$$

**Accuracy Parity (AP).** AP ensures that the classifier achieves equal predictive accuracy across sensitive groups. Unlike EO and EOpp, which condition on specific outcomes, AP enforces fairness in terms of overall correctness.

**Definition 4** (Accuracy Parity (AP)). *A classifier  $f_\theta$  satisfies Accuracy Parity (AP) if it has equal accuracy across all groups. Formally, for all groups  $a, b \in \mathcal{A}$ :*

$$\Pr(\hat{Y} = Y \mid \mathcal{A} = a) = \Pr(\hat{Y} = Y \mid \mathcal{A} = b)$$

We define the violation in Accuracy Parity as:

$$\Delta_{\text{AP}} := \left| \Pr(\hat{Y} = Y \mid \mathcal{A} = a) - \Pr(\hat{Y} = Y \mid \mathcal{A} = b) \right| \quad (5)$$

## 4 Methodology

In this section, we formalize the overall optimization objective in the FL setting for classification. We then describe the optimization problem solved by a strategic adversary aiming to amplify group-level bias in the resulting global model. Specifically, the adversary seeks to perform a model poisoning attack that circumvents state-of-the-art robust aggregation mechanisms.

### 4.1 FL Optimization

Recall that our goal is to train a global classifier  $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  that optimizes for global accuracy across a set of  $\mathcal{K}$  agents each having a private dataset  $\mathcal{D}^{(k)} = \left\{ (x_i^{(k)}, y_i^{(k)}, a_i^{(k)}) \right\}_{i=1}^{n_k}$ . Each agent  $k \in \mathcal{K}$  defines a local objective  $F_k(\theta)$  given by,

$$F_k(\theta) = \mathbb{E}_{x^{(k)}, y^{(k)} \sim \mathcal{D}^{(k)}} \left[ \ell_{CE}(\theta; x^{(k)}, y^{(k)}) \right] \quad (6)$$

where  $\ell_{CE}(\cdot)$  is the standard cross-entropy loss. The global objective is to minimize the **weighted aggregation** of local losses:  $\min_{\Theta} F(\Theta) = \sum_{k=1}^K w_k F_k(\Theta)$ . Here,  $\Theta$  denotes the model parameters to be optimized,  $w_k = \frac{n_k}{n}$  is the weight for agent  $k$ , with  $n_k = |\mathcal{D}_k|$  being the size of agent  $k$ 's dataset, and  $n = \sum_{k=1}^K n_k$  representing the total number of data points across all agents.

#### 4.1.1 FL Optimization: Adversarial agent

When a subset of agents is malicious, they aim to compromise the fairness of the global model by launching model poisoning attacks. Let  $K_A \subset K$  denote the set of *adversarial agents* among the  $K$  total agents. While honest agents minimize the standard local objective defined in Equation 6, adversarial agents aim to increase the demographic bias of the aggregated model while maintaining acceptable predictive accuracy. Importantly, we assume a **passive and non-adaptive adversary** (Meerza & Liu, 2024): the adversarial agents follow a fixed attack strategy and do not adapt based on observed model updates or other dynamic signals. Each adversarial agent  $k \in K_A$  solves the following:

$$\begin{array}{ll} \max_{\theta} \ell_F(\theta; \mathcal{D}^{(k)}) & \triangleright \text{Maximize Bias} \\ \text{s.t. } \mathbb{E}[\ell_{CE}(\theta, x^{(k)}, y^{(k)})] \leq \epsilon & \triangleright \text{Maintain Accuracy} \end{array}$$

Here,  $\ell_F(\cdot)$  is a differentiable surrogate objective designed to increase group-level fairness violations (e.g., for Demographic Parity or Equal Opportunity), and  $\ell_{CE}$  is the standard cross-entropy loss. The threshold  $\epsilon$  defines the maximum allowable performance degradation – controlling the *stealthiness* of the attack. Without this constraint, the adversary's update could be easily flagged by robust aggregators due to poor accuracy.

**Surrogate for Equalized Odds (Padala & Gujar, 2020).** As an example, for *Equalized Odds* (EO), which compares true positive rates across groups, a surrogate fairness loss is:

$$\ell_{EO} = \left| \frac{\sum_i (1-p_i) a_i y_i}{\sum_i a_i y_i} - \frac{\sum_i (1-p_i) (1-a_i) y_i}{\sum_i (1-a_i) y_i} \right| + \left| \frac{\sum_i p_i a_i (1-y_i)}{\sum_i a_i (1-y_i)} - \frac{\sum_i p_i (1-a_i) (1-y_i)}{\sum_i (1-a_i) (1-y_i)} \right|$$

Where  $p_i = f_{\theta}(x_i)$  denotes the predicted logits (or probability scores),  $a_i \in \{0, 1\}$  indicates binary group membership (e.g., gender), and  $y_i \in \{0, 1\}$  is the ground-truth label. Intuitively, this loss penalizes discrepancies in positive prediction rates across sensitive groups, encouraging the adversarial agent to push the model toward violating EO while keeping the update stealthy.

#### 4.1.2 Naïve Model Poisoning Attack

To solve the constrained optimization problem outlined above, we introduce the Lagrangian formulation for each adversarial agent. The Lagrangian multiplier  $\lambda \in \mathbb{R}_{\geq 0}$  will enforce the constraint on  $\mathbb{E}[\ell_{CE}(\cdot)] \leq \epsilon$  while allowing the adversarial agent to optimize for fairness (or bias) maximization. For all  $k \in K_A$ ,

$$F_k(\theta) = -\ell_F(\theta; \mathcal{D}^{(k)}) + \lambda \cdot \left( \mathbb{E} \left[ \ell_{CE}(\theta, x^{(k)}, y^{(k)}) \right] - \epsilon \right) \quad (7)$$

The empirical version of the above objective for  $\epsilon = 0$  is given by,

$$F_k(\theta) = -\ell_F(\theta; \mathcal{D}^{(k)}) + \lambda \cdot \frac{1}{n_k} \sum_{i=1}^{n_k} \ell_{CE}(\theta; x_i^{(k)}, y_i^{(k)}) \quad (8)$$

Equation 8 optimizes for maximizing fairness violation and minimizing accuracy simultaneously. The optimal parameter updates deviate significantly from an honest agent that solves Equation 6. Thus, the scores computed in  $m$ -KRUM (Algorithm B.1 (Line 4)) would easily help detect the adversarial agent. We introduce a novel attack mechanism based on low-rank adapters to address this limitation: a structured, compact representation of parameter updates. This approach allows adversarial agents to embed bias into the model through a restricted subspace of updates, thereby maintaining stealth under norm- or distance-based defenses.

**Algorithm 1** LoRA-FL

---

**Require:** Global model parameters  $\Theta_0$ , number of rounds  $T$ , local epochs  $E$ , adversarial local epochs  $E_{\mathcal{A}}$ , number of agents  $m$ , agent optimizer  $\text{OPT}$ , adversarial optimizers  $\text{OPT}_{\text{REG}}, \text{OPT}_{\text{F}}$ , scaling factor  $\alpha \in (0, 1]$ , aggregator function  $\text{Agg}$

**Ensure:** Aggregated global model  $\Theta_T$

```

1: for each round  $t = 1, 2, \dots, T$  do
2:   Server samples a subset of agents  $\mathcal{S}_t \subseteq \{1, \dots, K\}$ 
3:   for each agent  $i \in \mathcal{S}_t$  in parallel do
4:      $\theta_{i,t} \leftarrow \Theta_t$ 
5:     /* Agent  $i$  updates  $\theta_{i,t}$  locally using optimizer  $\text{OPT}$  */
6:     for each local epoch  $e = 1$  to  $E$  do
7:        $\theta_{i,t} \leftarrow \text{OPT}(\theta_{i,t}, \nabla \ell_{\text{CE}}(\theta_{i,t}; \mathcal{D}_i))$ 
8:     end for
9:     if  $i$  is Adversarial then
10:      Initialize adapter parameters  $A_{i,t}, B_{i,t}$ 
11:      for each adversarial local epoch  $e = 1$  to  $E_{\mathcal{A}}$  do
12:        for each batch in adversarial data do
13:          /* Phase 1: Train adapters for Accuracy */
14:           $(A_{i,t}, B_{i,t}) \leftarrow \text{OPT}_{\text{REG}}(A_{i,t}, B_{i,t}, \nabla \ell_{\text{REG}}(\cdot))$ 
15:          /* Phase 2: Train Adapters to Compromise Fairness */
16:           $(A_{i,t}, B_{i,t}) \leftarrow \text{OPT}_{\text{F}}(A_{i,t}, B_{i,t}, -\nabla \ell_{\text{F}}(\cdot))$ 
17:        end for
18:      end for
19:       $\theta_{i,t} \leftarrow \Theta_t + \alpha \cdot A_{i,t} B_{i,t}^\top$ 
20:      Agent sends updated model  $\theta_{i,t}$  to server
21:    end for
22:    Server aggregates agent updates:
23:       $\Theta_{t+1} \leftarrow \text{Agg}(\{\theta_{i,t}\}_{i \in [n]})$ 
24:  end for
25: return Final global model  $\Theta_T$ 

```

---

**4.2 LoRA-FL: Achieving Unfairness through Low-rank Adapters in FL**

Notice that the update to agent  $k$ 's local model after training at round  $t$  is the decomposition  $\theta_{k,t} = \Theta_{t-1} + \Delta\theta$ , where  $\Delta\theta$  represents the parameter change during the local update. Intuitively, an adversary's objective (from the discussion in Section 4.1.2) is that the update  $\Delta\theta$  encodes information that compromises fairness, while remaining close (in the parameter space) to  $\Theta$  to avoid detection by robust aggregators.

In our attack, namely LoRA-FL, an adversary achieves its objective by training low-rank matrices (aka *adapters*) that replace  $\Delta\theta$ , ensuring the desired behavior. Algorithm 1 presents the formal attack. For our discussion, consider the local model  $\theta \in \mathbb{R}^{d \times k}$  with adapters  $A \in \mathbb{R}^{d \times r}$  and  $B \in \mathbb{R}^{k \times r}$  as low-rank matrices such that  $r \ll \min(d, k)$ .

*Local Training.* Each agent  $k \in [K]$  (whether adversarial or benign) performs standard local updates as in conventional FL (Line 6, Algorithm 1). Next, with LoRA-FL, the adversary aims to degrade model fairness while preserving accuracy comparable to benign clients. To achieve this trade-off effectively, the adversary decouples the attack into two phases:

**Phase 1: Train Adapters for Accuracy.** The first phase focuses on stealth part of the attack, by improving the adapters' accuracy, using a regularizer that constrains the adapters to be close to  $\Delta\theta$ . Formally, for  $\theta_{k,t}$  as the current model (Line 12, Algorithm 1),

$$\ell_{\text{REG}}(A_{k,t}, B_{k,t}; \cdot) := \|A_{k,t} \cdot B_{k,t}^\top - (\Theta_{t-1} - \tilde{\theta}_{k,t})\|_2 \quad (9)$$

The optimizer  $\text{OPT}_{\text{REG}}$  minimizes  $\ell_{\text{REG}}$  during Phase 1, effectively driving the low-rank update  $A_{k,t} \cdot B_{k,t}^\top$  to be close to  $\Delta\theta$  – providing performance gains and avoiding detection.

Table 1: Comparison of different aggregators on accuracy and fairness metrics across the three tabular datasets for the IID setting. Here,  $r = 4$  for Adult and Dutch, and  $r = 2$  for Bank. **Acc** denotes accuracy, **Adv** denotes adversarial clients, and results are reported as  $\text{mean}_{\text{std}}$  over four independent runs. **Bold** values indicate highest accuracy and highest fairness violation, highlighting the strongest performance and most severe fairness degradation, respectively.

Agg	% Adv	Dataset											
		Adult				Bank				Dutch			
		Acc ( $\uparrow$ )	$\Delta_{EO}$ ( $\downarrow$ )	$\Delta_{EOpp}$ ( $\downarrow$ )	$\Delta_{DP}$ ( $\downarrow$ )	Acc ( $\uparrow$ )	$\Delta_{EO}$ ( $\downarrow$ )	$\Delta_{EOpp}$ ( $\downarrow$ )	$\Delta_{DP}$ ( $\downarrow$ )	Acc ( $\uparrow$ )	$\Delta_{EO}$ ( $\downarrow$ )	$\Delta_{EOpp}$ ( $\downarrow$ )	$\Delta_{DP}$ ( $\downarrow$ )
FedAvg	–	85.00 <sub>0.10</sub>	0.118 <sub>0.007</sub>	0.096 <sub>0.016</sub>	0.185 <sub>0.006</sub>	91.68 <sub>0.09</sub>	0.151 <sub>0.020</sub>	0.111 <sub>0.026</sub>	<b>0.211</b> <sub>0.043</sub>	<b>82.24</b> <sub>0.29</sub>	0.061 <sub>0.010</sub>	0.051 <sub>0.012</sub>	<b>0.188</b> <sub>0.009</sub>
KRUM	–	84.76 <sub>0.08</sub>	<b>0.129</b> <sub>0.019</sub>	<b>0.106</b> <sub>0.014</sub>	<b>0.195</b> <sub>0.011</sub>	91.49 <sub>0.08</sub>	<b>0.156</b> <sub>0.026</sub>	<b>0.122</b> <sub>0.018</sub>	0.203 <sub>0.039</sub>	82.26 <sub>0.20</sub>	0.061 <sub>0.011</sub>	0.052 <sub>0.015</sub>	0.187 <sub>0.011</sub>
TM	–	<b>85.05</b> <sub>0.08</sub>	0.101 <sub>0.012</sub>	0.073 <sub>0.011</sub>	0.185 <sub>0.015</sub>	<b>91.70</b> <sub>0.02</sub>	0.149 <sub>0.020</sub>	0.117 <sub>0.027</sub>	0.199 <sub>0.036</sub>	80.07 <sub>0.75</sub>	<b>0.073</b> <sub>0.009</sub>	<b>0.054</b> <sub>0.017</sub>	0.179 <sub>0.016</sub>
FedAvg	10%	<b>84.73</b> <sub>0.19</sub>	0.152 <sub>0.016</sub>	0.131 <sub>0.018</sub>	0.218 <sub>0.008</sub>	<b>91.68</b> <sub>0.11</sub>	0.172 <sub>0.0167</sub>	0.140 <sub>0.0189</sub>	0.238 <sub>0.0430</sub>	<b>80.50</b> <sub>0.20</sub>	0.072 <sub>0.002</sub>	0.029 <sub>0.006</sub>	0.207 <sub>0.009</sub>
FedAvg	20%	84.39 <sub>0.14</sub>	0.187 <sub>0.030</sub>	0.171 <sub>0.035</sub>	0.238 <sub>0.009</sub>	91.38 <sub>0.16</sub>	0.199 <sub>0.028</sub>	0.135 <sub>0.020</sub>	0.269 <sub>0.040</sub>	78.18 <sub>0.55</sub>	0.125 <sub>0.015</sub>	0.063 <sub>0.009</sub>	0.254 <sub>0.006</sub>
FedAvg	30%	83.51 <sub>0.25</sub>	0.266 <sub>0.030</sub>	0.261 <sub>0.032</sub>	0.285 <sub>0.008</sub>	90.83 <sub>0.24</sub>	0.262 <sub>0.0150</sub>	0.166 <sub>0.0314</sub>	0.338 <sub>0.0377</sub>	74.36 <sub>0.92</sub>	0.189 <sub>0.014</sub>	0.125 <sub>0.020</sub>	0.290 <sub>0.011</sub>
FedAvg	40%	81.95 <sub>0.43</sub>	<b>0.374</b> <sub>0.052</sub>	<b>0.371</b> <sub>0.055</sub>	<b>0.344</b> <sub>0.016</sub>	89.86 <sub>0.30</sub>	<b>0.312</b> <sub>0.015</sub>	<b>0.174</b> <sub>0.029</sub>	<b>0.384</b> <sub>0.016</sub>	70.08 <sub>0.48</sub>	<b>0.244</b> <sub>0.010</sub>	<b>0.216</b> <sub>0.021</sub>	<b>0.326</b> <sub>0.009</sub>
KRUM	10%	<b>84.56</b> <sub>0.28</sub>	0.157 <sub>0.017</sub>	0.135 <sub>0.019</sub>	0.213 <sub>0.016</sub>	<b>91.39</b> <sub>0.21</sub>	0.174 <sub>0.028</sub>	0.136 <sub>0.025</sub>	0.222 <sub>0.044</sub>	<b>82.11</b> <sub>0.22</sub>	0.065 <sub>0.011</sub>	0.053 <sub>0.010</sub>	0.188 <sub>0.009</sub>
KRUM	20%	84.14 <sub>0.56</sub>	0.181 <sub>0.068</sub>	0.157 <sub>0.070</sub>	0.234 <sub>0.043</sub>	91.21 <sub>0.36</sub>	0.196 <sub>0.055</sub>	0.150 <sub>0.037</sub>	0.258 <sub>0.068</sub>	78.85 <sub>0.41</sub>	0.091 <sub>0.008</sub>	0.029 <sub>0.007</sub>	0.220 <sub>0.006</sub>
KRUM	30%	82.87 <sub>1.25</sub>	<b>0.372</b> <sub>0.119</sub>	<b>0.367</b> <sub>0.125</sub>	<b>0.310</b> <sub>0.052</sub>	90.97 <sub>0.15</sub>	0.225 <sub>0.047</sub>	0.157 <sub>0.039</sub>	0.287 <sub>0.052</sub>	77.84 <sub>0.13</sub>	0.115 <sub>0.007</sub>	0.055 <sub>0.015</sub>	0.242 <sub>0.010</sub>
KRUM	40%	81.76 <sub>2.12</sub>	0.318 <sub>0.130</sub>	0.304 <sub>0.143</sub>	0.304 <sub>0.068</sub>	87.92 <sub>0.70</sub>	<b>0.396</b> <sub>0.036</sub>	<b>0.211</b> <sub>0.019</sub>	<b>0.463</b> <sub>0.046</sub>	76.85 <sub>0.29</sub>	<b>0.122</b> <sub>0.008</sub>	<b>0.076</b> <sub>0.007</sub>	<b>0.250</b> <sub>0.006</sub>
TM	10%	<b>84.73</b> <sub>0.26</sub>	0.108 <sub>0.010</sub>	0.091 <sub>0.013</sub>	0.170 <sub>0.019</sub>	<b>91.60</b> <sub>0.13</sub>	0.165 <sub>0.027</sub>	0.139 <sub>0.025</sub>	0.215 <sub>0.065</sub>	79.12 <sub>1.39</sub>	0.081 <sub>0.028</sub>	0.027 <sub>0.007</sub>	0.207 <sub>0.025</sub>
TM	20%	84.33 <sub>0.55</sub>	0.197 <sub>0.033</sub>	0.195 <sub>0.035</sub>	0.183 <sub>0.028</sub>	91.37 <sub>0.13</sub>	0.197 <sub>0.023</sub>	0.142 <sub>0.019</sub>	0.255 <sub>0.042</sub>	<b>79.79</b> <sub>0.37</sub>	0.087 <sub>0.008</sub>	0.025 <sub>0.006</sub>	0.220 <sub>0.009</sub>
TM	30%	84.24 <sub>0.33</sub>	0.201 <sub>0.052</sub>	0.192 <sub>0.060</sub>	0.224 <sub>0.011</sub>	91.16 <sub>0.13</sub>	0.216 <sub>0.006</sub>	0.159 <sub>0.024</sub>	0.278 <sub>0.031</sub>	79.43 <sub>0.32</sub>	0.087 <sub>0.011</sub>	0.033 <sub>0.004</sub>	0.227 <sub>0.009</sub>
TM	40%	82.82 <sub>1.05</sub>	<b>0.281</b> <sub>0.076</sub>	<b>0.280</b> <sub>0.077</sub>	<b>0.270</b> <sub>0.053</sub>	90.21 <sub>0.35</sub>	<b>0.299</b> <sub>0.030</sub>	<b>0.173</b> <sub>0.017</sub>	<b>0.358</b> <sub>0.053</sub>	79.18 <sub>0.64</sub>	<b>0.100</b> <sub>0.013</sub>	<b>0.036</b> <sub>0.013</sub>	<b>0.235</b> <sub>0.010</sub>

**Phase 2: Train Adapters to Compromise Fairness.** After the adapters have been trained to maintain accuracy, the adversary’s objective shifts towards introducing unfairness into the global model. The adversary aims to minimize the fairness loss  $\ell_{\text{UF}}$  concerning the adapter parameters. Formally, for  $\tilde{\theta}_{k,t}$  as the current model (Line 13, Algorithm 1), this loss function is designed to maximize the bias in the model’s predictions. Specifically,

$$\ell_{\text{UF}}(A_{i,t}, B_{i,t}; \cdot) := -\ell_F \quad \text{s.t.} \quad F \in \{\text{EO}, \text{EOpp}, \text{DP}, \text{AP}\} \quad (10)$$

where  $\ell_F$  represents a surrogate fairness loss for the chosen fairness metric. The optimizer  $\text{OPT}_F$  minimizes  $\ell_{\text{UF}}$  during this phase, effectively guiding the adapter parameters to introduce unfairness.

**Communication & Parameter Complexity.** In LoRA-FL, an adversary incurs no additional communication cost compared to the standard FL setup. From Algorithm 1 (Line 16), the adversary fuses the adapters into the base model and communicates the resulting model, identical to that of honest agents, to the Aggregator. For a base model with dimensions  $d \times k$  and adapter parameters  $d \times r$  and  $r \times k$ , the total number of parameters is  $\mathcal{O}(d \times k + r \times (d + k))$ . Since the adapters are low-rank, i.e.,  $r \ll \min(d, k)$ , the parameter increase is slight and *independent* in terms of  $d$  and  $k$ . Hence, for both honest and adversarial agents, the effective count is  $\mathcal{O}(d \times k)$ , with minimal overhead from the adapters.

## 5 Experiments & Results

**Datasets.** We evaluate our approach on four datasets spanning both tabular and image domains. Three are tabular binary classification datasets: **Adult** (Dua & Graff, 2017), which predicts whether an individual’s income exceeds \$50K using demographic attributes with **sex** and **race** treated as sensitive attributes; **Bank** (Moro et al., 2014), which predicts subscription to a term deposit based on demographic and contact-related features, where **age** (individuals aged 25-60 as the privileged group) is the sensitive attribute; **Dutch Census** (Iiobaite et al., 2011), which predicts occupation and uses **gender** as the binary sensitive attribute. These tabular datasets contain approximately 40K–60K samples each. We also consider image-based classification using **UTKFace** (Zhang et al., 2017), a dataset of approximately 20K face images. The task is *multi-class* ethnicity classification and **gender** is treated as the sensitive attribute. Together, these datasets enable evaluation of fairness and robustness across diverse data modalities and classification settings.

**IID and Non-IID Client Data.** We evaluate our approach under both IID and non-IID FL settings. In the IID setting, each client’s local dataset is formed by uniformly sampling from the global dataset, so that all clients share the same distribution over features and labels. In the non-IID setting, clients possess partitioned data using a Dirichlet distribution (Ng et al., 2011) over *class labels* with concentration parameter  $\rho \in \{0.25, 0.5, 0.75\}$ , where smaller values of  $\rho$  induce stronger label and feature skew across clients.

## 5.1 Model Architectures and Adapter Placement

### 5.1.1 Adapters in Dense Layers (Tabular Datasets)

In the first FL setup, corresponding to tabular datasets, each honest client employs a two-layer multilayer perceptron (MLP) with hidden layer sizes of 64 and 32, together with dataset-specific input and output heads. Adversarial agents augment this architecture by inserting low-rank adapters into the hidden layers of the MLP, with adapter rank set to  $r = 4$ . All agents use ReLU activations (Nair & Hinton, 2010). Optimization is performed using AdamW (Loshchilov & Hutter, 2019) without momentum for both honest and adversarial updates, denoted by OPT, OPT<sub>F</sub>, and OPT<sub>REG</sub>. Additional training details and hyperparameter settings are provided in Appendix A.1.

### 5.1.2 Adapters in CNN and Dense Layers (Image Datasets)

In the second FL setup, corresponding to image datasets, each honest agent uses an untrained ResNet-18 backbone (He et al., 2016) followed by a two-layer MLP classification head identical to the tabular setting. This setting is substantially more challenging than the tabular case due to the depth and hierarchical structure of convolutional representations, which distribute task- and fairness-relevant features across multiple layers. Adversarial agents insert low-rank adapters into both the convolutional layers of the ResNet backbone and the subsequent dense layers. This design choice is non-trivial: restricting adapters to only the final classification head is often insufficient to influence higher-level semantic features that correlate with sensitive attributes. By fusing adapters across both convolutional and dense layers, adversarial agents can steer fairness-relevant representations throughout the network while maintaining accuracy. The adapter design follows prior work on parameter-efficient fine-tuning of convolutional networks, such as LoRA-C (Ding et al., 2024). All experiments are implemented in PyTorch (Paszke et al., 2019) and conducted on an NVIDIA L40S GPU with 48 GB memory. Further hyperparameter details are provided in Appendix A.1.

### 5.1.3 Performance & Fairness Measures

We benchmark LoRA-FL’s efficacy based on accuracy as our performance measure, and  $\Delta_{DP}$ ,  $\Delta_{EOpp}$ , and  $\Delta_{EO}$  as fairness measures for binary classifications and  $\Delta_{AP}$  as a fairness metric for multiclass classification. Here, we note that improvements in fairness typically come at the cost of accuracy (Bilal Zafar et al., 2015; Madras et al., 2018). The trade-off between performance and fairness measures is a *crucial* factor<sup>1</sup> in evaluating the overall performance.

### 5.1.4 Other Details

We report **average** values across all agents using 25% of their local data for testing. We benchmark performance and fairness against FedAvg, KRUM, and TM for IID distributions and include FLAME for non-IID distributions. Adversarial agent percentages varies across {10%, 20%, 30%, 40% }.

For KRUM and TM, we take  $m = 60\%$  of all clients for all adversary settings. This corresponds to the *worst-case* permissible choice across our experiments and ensures the aggregation is defined even for the highest adversary ratio. Additionally, we do not evaluate LoRA-FL against fair aggregators like FairFed (Ezzeldin et al., 2023) because adversarial clients can easily manipulate these systems by reporting inaccurate fairness scores, rendering them ineffective.

<sup>1</sup>**Example:** A model always predicting a single class may appear fair under Demographic Parity (DP), but its poor accuracy makes it ineffective.

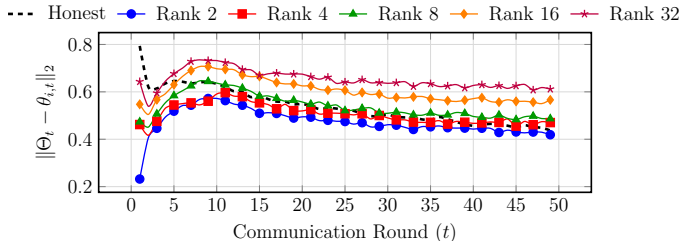


Figure 2: Effect of low-rank constraints on model divergence. As the rank increases, the  $\ell_2$  distance  $\|\Theta_t - \theta_{i,t}\|_2$  grows, indicating that adversarial updates diverge more from the global model.

## 5.2 Results

Table 1 and Table 2 present results across tabular datasets under both IID and non-IID data distributions. The numbers show the effectiveness of LoRA-FL across aggregation rules, datasets, and non-IID parameter  $\rho$ . To evaluate cross-domain generality of LoRA-FL, we perform image-based multi-class classification on UTKFace (Table 3) and compare against existing attacks such as EAB-FL (Meerza & Liu, 2024) (Table 4). Details follow.

**Effectiveness of LoRA-FL across aggregators.** Across both IID (Table 1) and non-IID settings (Table 2), LoRA-FL consistently achieves a strong trade-off between accuracy preservation and fairness degradation, even in the presence of robust aggregation rules. Under IID data, LoRA-FL induces substantial increases in fairness gaps with minimal loss in accuracy across FedAvg, KRUM, and TM. Notably, this behavior persists across non-IID data: for all values of  $\rho$ , both FLAME and KRUM exhibit fairness degradation patterns comparable to those of FedAvg. Overall, the results from these tables demonstrate that LoRA-FL generalizes across aggregation rules and data regimes.

In Table 1 for Adult with TM, under a 30% adversary setup, accuracy declines by 0.81% (from 85.05% to 84.24%), yet the EO gap nearly doubles – from 0.101 to 0.201 (a 99.0% surge) – and the DP gap rises from 0.185 to 0.224 (21.1%). Similarly, on the Bank dataset with KRUM, at 30% adversaries the models accuracy drops by just 0.52% (from 91.49% to 90.97%), while the EO gap increases from 0.156 to 0.225 (44.2%) and the DP gap from 0.203 to 0.287 (41.4%). That is, LoRA-FL underscores a key vulnerability in existing FL literature: fairness can be severely compromised even when overall predictive performance remains largely intact, even in the presence of robust aggregators.

**Effect of data heterogeneity ( $\rho$ ).** Table 2 shows that data heterogeneity substantially amplifies the impact of LoRA-FL. As  $\rho$  decreases (i.e., client data becomes more non-IID), fairness violations increase sharply across all aggregators and datasets, often accompanied by higher variance. In contrast, when  $\rho$  is large, fairness degradation remains pronounced but is relatively stable. When compared against the IID baseline under the same rank in Table 1, the results indicate that non-IID data not only worsens fairness outcomes but also obscures adversarial behavior, allowing LoRA-FL to remain stealthy while inducing larger disparities.

**Dataset-dependent behavior.** Comparing across datasets in both tables reveals consistent but dataset-specific vulnerability patterns. The **Adult** and **Bank** datasets exhibit the most severe fairness degradation under LoRA-FL, with EO and DP gaps often increasing by 40–100% at moderate adversarial fractions while accuracy drops remain below 1–2%. These effects are further magnified under non-IID settings, particularly for smaller values of  $\rho$ . In contrast, the **Dutch** dataset shows smaller absolute fairness gaps in both IID and non-IID regimes, though the monotonic increase in violations with adversarial participation persists. This suggests that while LoRA-FL is broadly effective, datasets with richer demographic structure and stronger correlations between sensitive attributes and labels are especially susceptible to fairness-targeted manipulation.

**LoRA-FL generalizes across Domain.** As discussed earlier, we evaluate LoRA-FL on UTKFace dataset by considering multi-class image classification using a ResNet-based convolutional architecture. Table 3

Table 2: Comparison of different aggregators on accuracy and fairness metrics across the three tabular datasets for the **non-IID** setting. Here,  $\rho \in \{0.5, 0.75\}$  and  $r = 2$  for all the datasets. Also, **Acc**: Accuracy, **Adv**: Adversary, and we report  $\text{mean}_{\text{std}}$  across **four** independent runs. **Bold** values indicate highest accuracy and highest fairness violation, highlighting the strongest performance and most severe fairness degradation, respectively.

	Method	% Adv	$\rho = 0.5$				$\rho = 0.75$			
			Acc ( $\uparrow$ )	$\Delta_{EO}$ ( $\downarrow$ )	$\Delta_{EOpp}$ ( $\downarrow$ )	$\Delta_{DP}$ ( $\downarrow$ )	Acc ( $\uparrow$ )	$\Delta_{EO}$ ( $\downarrow$ )	$\Delta_{EOpp}$ ( $\downarrow$ )	$\Delta_{DP}$ ( $\downarrow$ )
Adult	FedAvg	-	<b>84.80</b> <sub>0.58</sub>	0.161 <sub>0.024</sub>	0.150 <sub>0.026</sub>	0.191 <sub>0.018</sub>	<b>84.72</b> <sub>0.78</sub>	0.162 <sub>0.016</sub>	0.139 <sub>0.014</sub>	0.193 <sub>0.011</sub>
		10%	84.79 <sub>0.59</sub>	0.168 <sub>0.033</sub>	0.159 <sub>0.038</sub>	0.195 <sub>0.012</sub>	84.55 <sub>0.72</sub>	0.159 <sub>0.013</sub>	0.137 <sub>0.013</sub>	0.195 <sub>0.007</sub>
		20%	84.55 <sub>0.51</sub>	0.168 <sub>0.026</sub>	0.156 <sub>0.034</sub>	0.205 <sub>0.015</sub>	83.87 <sub>1.48</sub>	0.180 <sub>0.034</sub>	0.146 <sub>0.025</sub>	0.221 <sub>0.020</sub>
		30%	84.58 <sub>0.65</sub>	0.176 <sub>0.025</sub>	0.165 <sub>0.033</sub>	0.212 <sub>0.007</sub>	83.22 <sub>2.13</sub>	0.201 <sub>0.068</sub>	0.173 <sub>0.058</sub>	0.245 <sub>0.048</sub>
		40%	82.91 <sub>1.37</sub>	<b>0.240</b> <sub>0.043</sub>	<b>0.219</b> <sub>0.044</sub>	<b>0.265</b> <sub>0.034</sub>	82.88 <sub>0.73</sub>	<b>0.204</b> <sub>0.026</sub>	<b>0.175</b> <sub>0.025</sub>	<b>0.253</b> <sub>0.039</sub>
	FLAME	-	<b>84.81</b> <sub>0.59</sub>	0.162 <sub>0.027</sub>	0.151 <sub>0.032</sub>	0.192 <sub>0.019</sub>	<b>84.76</b> <sub>0.73</sub>	0.164 <sub>0.017</sub>	0.140 <sub>0.014</sub>	0.194 <sub>0.011</sub>
		10%	84.71 <sub>0.64</sub>	0.158 <sub>0.023</sub>	0.144 <sub>0.027</sub>	0.195 <sub>0.014</sub>	84.58 <sub>0.79</sub>	0.161 <sub>0.014</sub>	0.135 <sub>0.012</sub>	0.195 <sub>0.007</sub>
		20%	84.48 <sub>0.55</sub>	0.163 <sub>0.028</sub>	0.149 <sub>0.036</sub>	0.204 <sub>0.014</sub>	83.84 <sub>1.48</sub>	0.180 <sub>0.035</sub>	0.147 <sub>0.025</sub>	0.221 <sub>0.020</sub>
		30%	84.55 <sub>0.59</sub>	0.175 <sub>0.023</sub>	0.164 <sub>0.031</sub>	0.211 <sub>0.007</sub>	83.14 <sub>2.19</sub>	0.203 <sub>0.070</sub>	0.173 <sub>0.060</sub>	0.245 <sub>0.052</sub>
		40%	82.89 <sub>1.36</sub>	<b>0.245</b> <sub>0.045</sub>	<b>0.223</b> <sub>0.048</sub>	<b>0.266</b> <sub>0.032</sub>	82.86 <sub>0.62</sub>	<b>0.204</b> <sub>0.027</sub>	<b>0.174</b> <sub>0.024</sub>	<b>0.257</b> <sub>0.043</sub>
	KRUM	-	<b>80.34</b> <sub>1.05</sub>	0.264 <sub>0.044</sub>	0.212 <sub>0.043</sub>	0.308 <sub>0.035</sub>	<b>81.79</b> <sub>1.03</sub>	0.253 <sub>0.042</sub>	0.206 <sub>0.045</sub>	0.273 <sub>0.022</sub>
		10%	80.81 <sub>0.90</sub>	0.279 <sub>0.066</sub>	0.242 <sub>0.075</sub>	0.304 <sub>0.021</sub>	81.59 <sub>1.98</sub>	0.259 <sub>0.080</sub>	0.213 <sub>0.086</sub>	0.283 <sub>0.046</sub>
20%		80.46 <sub>1.26</sub>	0.288 <sub>0.056</sub>	0.243 <sub>0.058</sub>	0.312 <sub>0.038</sub>	80.71 <sub>1.51</sub>	0.265 <sub>0.019</sub>	0.207 <sub>0.029</sub>	0.301 <sub>0.051</sub>	
30%		79.67 <sub>0.91</sub>	0.335 <sub>0.053</sub>	0.294 <sub>0.061</sub>	0.349 <sub>0.033</sub>	79.50 <sub>1.84</sub>	0.313 <sub>0.033</sub>	0.254 <sub>0.038</sub>	0.343 <sub>0.055</sub>	
40%		75.56 <sub>6.91</sub>	<b>0.412</b> <sub>0.172</sub>	<b>0.387</b> <sub>0.162</sub>	<b>0.363</b> <sub>0.107</sub>	78.30 <sub>3.42</sub>	<b>0.435</b> <sub>0.213</sub>	<b>0.423</b> <sub>0.221</sub>	<b>0.383</b> <sub>0.127</sub>	
TM	-	<b>84.00</b> <sub>0.29</sub>	0.167 <sub>0.063</sub>	0.148 <sub>0.067</sub>	0.224 <sub>0.025</sub>	83.73 <sub>1.15</sub>	0.177 <sub>0.053</sub>	0.151 <sub>0.058</sub>	0.217 <sub>0.020</sub>	
	10%	83.76 <sub>0.36</sub>	0.291 <sub>0.114</sub>	0.280 <sub>0.117</sub>	0.264 <sub>0.040</sub>	<b>84.20</b> <sub>1.06</sub>	0.171 <sub>0.081</sub>	0.142 <sub>0.095</sub>	0.214 <sub>0.040</sub>	
	20%	83.49 <sub>0.42</sub>	0.270 <sub>0.100</sub>	0.254 <sub>0.118</sub>	0.257 <sub>0.036</sub>	83.77 <sub>0.88</sub>	0.224 <sub>0.030</sub>	0.196 <sub>0.046</sub>	0.232 <sub>0.032</sub>	
	30%	82.09 <sub>1.48</sub>	0.405 <sub>0.187</sub>	0.400 <sub>0.190</sub>	0.326 <sub>0.077</sub>	82.79 <sub>1.58</sub>	0.351 <sub>0.154</sub>	0.342 <sub>0.159</sub>	0.279 <sub>0.057</sub>	
	40%	78.30 <sub>4.01</sub>	<b>0.607</b> <sub>0.256</sub>	<b>0.599</b> <sub>0.270</sub>	<b>0.427</b> <sub>0.137</sub>	81.32 <sub>2.53</sub>	<b>0.527</b> <sub>0.200</sub>	<b>0.520</b> <sub>0.209</sub>	<b>0.322</b> <sub>0.099</sub>	
Bank	FedAvg	-	<b>91.31</b> <sub>0.41</sub>	0.272 <sub>0.055</sub>	0.212 <sub>0.063</sub>	0.247 <sub>0.021</sub>	90.89 <sub>0.21</sub>	0.203 <sub>0.038</sub>	0.143 <sub>0.038</sub>	0.240 <sub>0.019</sub>
		10%	91.24 <sub>0.41</sub>	0.280 <sub>0.052</sub>	0.209 <sub>0.063</sub>	0.267 <sub>0.025</sub>	<b>90.94</b> <sub>0.24</sub>	0.208 <sub>0.038</sub>	0.146 <sub>0.043</sub>	0.244 <sub>0.017</sub>
		20%	90.97 <sub>0.32</sub>	0.292 <sub>0.057</sub>	0.205 <sub>0.060</sub>	0.295 <sub>0.042</sub>	90.01 <sub>1.75</sub>	0.247 <sub>0.090</sub>	0.145 <sub>0.042</sub>	0.296 <sub>0.080</sub>
		30%	90.96 <sub>0.23</sub>	0.310 <sub>0.059</sub>	<b>0.219</b> <sub>0.076</sub>	0.315 <sub>0.030</sub>	89.14 <sub>3.09</sub>	0.293 <sub>0.084</sub>	0.186 <sub>0.035</sub>	0.335 <sub>0.071</sub>
		40%	90.01 <sub>0.84</sub>	<b>0.350</b> <sub>0.085</sub>	0.216 <sub>0.070</sub>	<b>0.387</b> <sub>0.067</sub>	85.00 <sub>10.20</sub>	<b>0.345</b> <sub>0.083</sub>	<b>0.192</b> <sub>0.048</sub>	<b>0.365</b> <sub>0.046</sub>
	FLAME	-	<b>91.30</b> <sub>0.47</sub>	0.275 <sub>0.054</sub>	0.214 <sub>0.063</sub>	0.244 <sub>0.027</sub>	90.90 <sub>0.22</sub>	0.199 <sub>0.038</sub>	0.143 <sub>0.043</sub>	0.238 <sub>0.018</sub>
		10%	91.25 <sub>0.41</sub>	0.279 <sub>0.051</sub>	0.210 <sub>0.061</sub>	0.267 <sub>0.027</sub>	<b>90.96</b> <sub>0.22</sub>	0.207 <sub>0.036</sub>	0.147 <sub>0.042</sub>	0.239 <sub>0.016</sub>
		20%	90.99 <sub>0.35</sub>	0.287 <sub>0.057</sub>	0.203 <sub>0.064</sub>	0.298 <sub>0.032</sub>	90.00 <sub>1.80</sub>	0.252 <sub>0.094</sub>	0.155 <sub>0.037</sub>	0.302 <sub>0.086</sub>
		30%	90.90 <sub>0.21</sub>	0.305 <sub>0.060</sub>	0.214 <sub>0.073</sub>	0.315 <sub>0.029</sub>	89.13 <sub>3.08</sub>	0.291 <sub>0.082</sub>	0.182 <sub>0.033</sub>	0.334 <sub>0.071</sub>
		40%	90.00 <sub>0.77</sub>	<b>0.351</b> <sub>0.083</sub>	<b>0.215</b> <sub>0.070</sub>	<b>0.392</b> <sub>0.064</sub>	85.42 <sub>9.25</sub>	<b>0.338</b> <sub>0.075</sub>	<b>0.193</b> <sub>0.050</sub>	<b>0.360</b> <sub>0.046</sub>
	KRUM	-	<b>89.42</b> <sub>0.72</sub>	0.404 <sub>0.081</sub>	0.306 <sub>0.092</sub>	0.425 <sub>0.073</sub>	<b>89.39</b> <sub>0.44</sub>	0.356 <sub>0.042</sub>	0.221 <sub>0.029</sub>	0.392 <sub>0.077</sub>
		10%	87.71 <sub>1.23</sub>	0.451 <sub>0.087</sub>	0.244 <sub>0.078</sub>	0.485 <sub>0.061</sub>	88.33 <sub>1.22</sub>	0.369 <sub>0.073</sub>	0.224 <sub>0.051</sub>	0.401 <sub>0.053</sub>
20%		87.23 <sub>0.84</sub>	0.484 <sub>0.063</sub>	0.278 <sub>0.044</sub>	0.502 <sub>0.055</sub>	87.81 <sub>1.43</sub>	0.393 <sub>0.061</sub>	0.221 <sub>0.056</sub>	0.430 <sub>0.053</sub>	
30%		85.54 <sub>2.41</sub>	0.494 <sub>0.014</sub>	0.246 <sub>0.064</sub>	0.479 <sub>0.063</sub>	86.45 <sub>1.53</sub>	0.415 <sub>0.079</sub>	0.248 <sub>0.067</sub>	0.448 <sub>0.078</sub>	
40%		84.97 <sub>2.12</sub>	<b>0.516</b> <sub>0.053</sub>	<b>0.309</b> <sub>0.064</sub>	<b>0.512</b> <sub>0.022</sub>	80.88 <sub>6.45</sub>	<b>0.439</b> <sub>0.021</sub>	<b>0.263</b> <sub>0.061</sub>	<b>0.474</b> <sub>0.034</sub>	
TM	-	<b>90.89</b> <sub>0.74</sub>	0.288 <sub>0.084</sub>	0.193 <sub>0.086</sub>	0.285 <sub>0.070</sub>	90.90 <sub>0.79</sub>	0.222 <sub>0.045</sub>	0.175 <sub>0.034</sub>	0.259 <sub>0.035</sub>	
	10%	90.75 <sub>0.54</sub>	0.301 <sub>0.038</sub>	0.234 <sub>0.079</sub>	0.297 <sub>0.079</sub>	<b>91.35</b> <sub>0.31</sub>	0.204 <sub>0.054</sub>	0.152 <sub>0.038</sub>	0.246 <sub>0.042</sub>	
	20%	90.13 <sub>1.07</sub>	0.399 <sub>1.23</sub>	0.292 <sub>0.117</sub>	0.404 <sub>0.128</sub>	90.45 <sub>0.89</sub>	0.250 <sub>0.085</sub>	<b>0.179</b> <sub>0.050</sub>	0.260 <sub>0.061</sub>	
	30%	88.62 <sub>1.14</sub>	0.427 <sub>0.069</sub>	0.313 <sub>0.081</sub>	0.481 <sub>0.056</sub>	86.38 <sub>3.83</sub>	0.336 <sub>0.067</sub>	0.151 <sub>0.058</sub>	0.380 <sub>0.062</sub>	
	40%	86.77 <sub>2.99</sub>	<b>0.481</b> <sub>0.080</sub>	<b>0.328</b> <sub>0.108</sub>	<b>0.506</b> <sub>0.049</sub>	82.93 <sub>3.88</sub>	<b>0.429</b> <sub>0.094</sub>	0.156 <sub>0.060</sub>	<b>0.442</b> <sub>0.069</sub>	
Dutch	FedAvg	-	81.53 <sub>0.44</sub>	0.075 <sub>0.008</sub>	0.069 <sub>0.008</sub>	0.170 <sub>0.009</sub>	82.01 <sub>0.32</sub>	0.070 <sub>0.004</sub>	0.062 <sub>0.008</sub>	0.185 <sub>0.009</sub>
		10%	<b>81.54</b> <sub>0.47</sub>	0.073 <sub>0.009</sub>	0.068 <sub>0.009</sub>	<b>0.170</b> <sub>0.008</sub>	<b>82.06</b> <sub>0.40</sub>	0.071 <sub>0.005</sub>	0.063 <sub>0.009</sub>	<b>0.185</b> <sub>0.011</sub>
		20%	81.47 <sub>0.42</sub>	0.076 <sub>0.009</sub>	0.071 <sub>0.010</sub>	0.169 <sub>0.009</sub>	81.78 <sub>0.20</sub>	0.079 <sub>0.011</sub>	<b>0.073</b> <sub>0.011</sub>	0.174 <sub>0.007</sub>
		30%	81.47 <sub>0.46</sub>	0.078 <sub>0.011</sub>	0.074 <sub>0.010</sub>	0.167 <sub>0.008</sub>	81.70 <sub>0.26</sub>	0.081 <sub>0.015</sub>	0.070 <sub>0.023</sub>	0.177 <sub>0.019</sub>
		40%	81.42 <sub>0.47</sub>	<b>0.082</b> <sub>0.011</sub>	<b>0.079</b> <sub>0.010</sub>	0.162 <sub>0.011</sub>	71.14 <sub>12.87</sub>	<b>0.083</b> <sub>0.018</sub>	0.067 <sub>0.016</sub>	0.115 <sub>0.065</sub>
	FLAME	-	81.52 <sub>0.41</sub>	0.074 <sub>0.009</sub>	0.069 <sub>0.008</sub>	<b>0.171</b> <sub>0.008</sub>	82.01 <sub>0.33</sub>	0.070 <sub>0.004</sub>	0.062 <sub>0.008</sub>	0.184 <sub>0.009</sub>
		10%	<b>81.54</b> <sub>0.46</sub>	0.073 <sub>0.009</sub>	0.068 <sub>0.009</sub>	0.170 <sub>0.008</sub>	<b>82.04</b> <sub>0.41</sub>	0.071 <sub>0.005</sub>	0.063 <sub>0.009</sub>	<b>0.185</b> <sub>0.010</sub>
		20%	81.46 <sub>0.44</sub>	0.076 <sub>0.009</sub>	0.071 <sub>0.010</sub>	0.169 <sub>0.009</sub>	81.78 <sub>0.19</sub>	0.079 <sub>0.010</sub>	<b>0.073</b> <sub>0.011</sub>	0.174 <sub>0.007</sub>
		30%	81.45 <sub>0.49</sub>	0.078 <sub>0.011</sub>	0.074 <sub>0.010</sub>	0.166 <sub>0.009</sub>	81.48 <sub>0.38</sub>	<b>0.086</b> <sub>0.015</sub>	0.069 <sub>0.024</sub>	0.179 <sub>0.022</sub>
		40%	81.39 <sub>0.50</sub>	<b>0.083</b> <sub>0.010</sub>	<b>0.080</b> <sub>0.010</sub>	0.162 <sub>0.011</sub>	71.46 <sub>12.45</sub>	0.080 <sub>0.013</sub>	0.063 <sub>0.019</sub>	0.116 <sub>0.064</sub>
	KRUM	-	<b>81.49</b> <sub>0.29</sub>	0.081 <sub>0.009</sub>	0.075 <sub>0.008</sub>	0.156 <sub>0.023</sub>	82.00 <sub>0.66</sub>	0.080 <sub>0.008</sub>	0.066 <sub>0.010</sub>	0.181 <sub>0.016</sub>
		10%	81.28 <sub>0.35</sub>	0.083 <sub>0.012</sub>	0.077 <sub>0.012</sub>	0.159 <sub>0.020</sub>	81.97 <sub>0.58</sub>	0.080 <sub>0.007</sub>	0.072 <sub>0.003</sub>	0.176 <sub>0.014</sub>
20%		81.16 <sub>0.70</sub>	0.084 <sub>0.016</sub>	0.077 <sub>0.014</sub>	0.162 <sub>0.018</sub>	81.79 <sub>0.38</sub>	<b>0.083</b> <sub>0.010</sub>	0.070 <sub>0.004</sub>	0.179 <sub>0.018</sub>	
30%		80.98 <sub>1.10</sub>	0.086 <sub>0.012</sub>	0.077 <sub>0.010</sub>	0.164 <sub>0.016</sub>	<b>82.03</b> <sub>0.79</sub>	0.070 <sub>0.011</sub>	0.058 <sub>0.020</sub>	<b>0.186</b> <sub>0.020</sub>	
40%		80.72 <sub>1.22</sub>	<b>0.091</b> <sub>0.018</sub>	<b>0.086</b> <sub>0.020</sub>	<b>0.166</b> <sub>0.015</sub>	81.73 <sub>0.68</sub>	0.082 <sub>0.014</sub>	<b>0.075</b> <sub>0.016</sub>	0.177 <sub>0.017</sub>	
TM	-	<b>81.63</b> <sub>0.58</sub>	0.073 <sub>0.013</sub>	0.068 <sub>0.014</sub>	0.169 <sub>0.010</sub>	<b>82.13</b> <sub>0.25</sub>	0.073 <sub>0.002</sub>	0.065 <sub>0.008</sub>	<b>0.183</b> <sub>0.011</sub>	
	10%	81.55 <sub>0.42</sub>	0.080 <sub>0.018</sub>	0.075 <sub>0.016</sub>	0.164 <sub>0.011</sub>	81.96 <sub>0.30</sub>	0.080 <sub>0.013</sub>	0.074 <sub>0.008</sub>	0.179 <sub>0.009</sub>	
	20%	81.68 <sub>0.65</sub>	0.080 <sub>0.014</sub>	0.075 <sub>0.015</sub>	0.162 <sub>0.019</sub>	81.60 <sub>0.16</sub>	0.081 <sub>0.016</sub>	<b>0.077</b> <sub>0.014</sub>	0.165 <sub>0.004</sub>	
	30%	81.42 <sub>1.04</sub>	0.083 <sub>0.015</sub>	0.080 <sub>0.017</sub>	<b>0.171</b> <sub>0.017</sub>	81.46 <sub>0.38</sub>	0.084 <sub>0.013</sub>	0.075 <sub>0.018</sub>	0.171 <sub>0.010</sub>	
	40%	79.27 <sub>3.41</sub>	<b>0.122</b> <sub>0.069</sub>	<b>0.119</b> <sub>0.072</sub>	0.128 <sub>0.054</sub>	81.25 <sub>0.52</sub>	<b>0.086</b> <sub>0.014</sub>	0.070 <sub>0.025</sub>	0.172 <sub>0.023</sub>	

Table 3: Comparison of different aggregators on accuracy and Accuracy parity (AP) in a **non-IID** setting on **UTKFace**. Here,  $r = 8$  with concentration parameters  $\rho \in \{0.25, 0.5, 0.75\}$ . Also, **Acc**: Accuracy, **Adv**: Adversary and we report **mean<sub>std</sub>** across four independent runs. **Bold** values indicate: highest accuracy and highest fairness violation, highlighting the strongest performance and most severe fairness degradation, respectively.

Method	% Adv	$\rho = 0.25$		$\rho = 0.5$		$\rho = 0.75$	
		Acc ( $\uparrow$ )	$\Delta_{AP}$ ( $\downarrow$ )	Acc ( $\uparrow$ )	$\Delta_{AP}$ ( $\downarrow$ )	Acc ( $\uparrow$ )	$\Delta_{AP}$ ( $\downarrow$ )
FedAvg	–	<b>68.24</b> <sub>5.46</sub>	0.035 <sub>0.002</sub>	<b>71.66</b> <sub>0.33</sub>	0.023 <sub>0.003</sub>	<b>72.35</b> <sub>0.24</sub>	0.030 <sub>0.006</sub>
	10%	61.44 <sub>8.83</sub>	0.066 <sub>0.017</sub>	69.26 <sub>0.57</sub>	0.052 <sub>0.012</sub>	69.96 <sub>0.06</sub>	0.041 <sub>0.009</sub>
	20%	62.43 <sub>3.34</sub>	0.091 <sub>0.048</sub>	65.78 <sub>0.29</sub>	0.085 <sub>0.007</sub>	66.04 <sub>0.67</sub>	0.093 <sub>0.025</sub>
	30%	57.63 <sub>4.27</sub>	0.141 <sub>0.046</sub>	64.49 <sub>1.19</sub>	0.122 <sub>0.026</sub>	62.64 <sub>0.20</sub>	0.136 <sub>0.013</sub>
	40%	56.25 <sub>3.28</sub>	<b>0.203</b> <sub>0.059</sub>	59.35 <sub>0.65</sub>	<b>0.161</b> <sub>0.024</sub>	60.75 <sub>0.61</sub>	<b>0.162</b> <sub>0.025</sub>
FLAME	–	<b>69.65</b> <sub>6.82</sub>	0.033 <sub>0.005</sub>	<b>72.68</b> <sub>1.12</sub>	0.021 <sub>0.002</sub>	<b>72.73</b> <sub>0.29</sub>	0.032 <sub>0.007</sub>
	10%	66.98 <sub>4.44</sub>	0.036 <sub>0.012</sub>	69.47 <sub>0.13</sub>	0.051 <sub>0.018</sub>	70.13 <sub>0.48</sub>	0.028 <sub>0.001</sub>
	20%	63.12 <sub>1.67</sub>	0.073 <sub>0.022</sub>	68.49 <sub>5.36</sub>	0.096 <sub>0.035</sub>	67.25 <sub>1.28</sub>	0.080 <sub>0.043</sub>
	30%	61.30 <sub>1.05</sub>	0.085 <sub>0.017</sub>	63.91 <sub>0.61</sub>	0.117 <sub>0.033</sub>	67.39 <sub>1.30</sub>	0.055 <sub>0.022</sub>
	40%	55.24 <sub>3.12</sub>	<b>0.201</b> <sub>0.034</sub>	65.32 <sub>0.03</sub>	<b>0.144</b> <sub>0.044</sub>	60.98 <sub>1.60</sub>	<b>0.104</b> <sub>0.003</sub>
KRUM	–	62.16 <sub>8.48</sub>	0.031 <sub>0.006</sub>	<b>74.53</b> <sub>0.42</sub>	0.022 <sub>0.004</sub>	<b>71.31</b> <sub>0.99</sub>	0.034 <sub>0.005</sub>
	10%	<b>71.96</b> <sub>1.90</sub>	0.034 <sub>0.005</sub>	72.21 <sub>0.46</sub>	0.038 <sub>0.003</sub>	68.38 <sub>3.82</sub>	0.038 <sub>0.020</sub>
	20%	71.51 <sub>1.71</sub>	0.045 <sub>0.010</sub>	70.15 <sub>2.86</sub>	0.061 <sub>0.005</sub>	67.80 <sub>3.20</sub>	0.043 <sub>0.010</sub>
	30%	64.36 <sub>0.82</sub>	0.077 <sub>0.018</sub>	67.49 <sub>4.74</sub>	0.080 <sub>0.031</sub>	63.24 <sub>4.46</sub>	0.067 <sub>0.030</sub>
	40%	63.98 <sub>1.91</sub>	<b>0.086</b> <sub>0.000</sub>	64.72 <sub>1.15</sub>	<b>0.125</b> <sub>0.009</sub>	62.22 <sub>2.26</sub>	<b>0.097</b> <sub>0.004</sub>
TM	–	<b>70.38</b> <sub>1.67</sub>	0.031 <sub>0.004</sub>	<b>71.84</b> <sub>0.07</sub>	0.022 <sub>0.009</sub>	<b>72.87</b> <sub>0.33</sub>	0.022 <sub>0.005</sub>
	10%	67.86 <sub>2.25</sub>	0.068 <sub>0.019</sub>	68.69 <sub>0.86</sub>	0.043 <sub>0.008</sub>	69.71 <sub>0.70</sub>	0.044 <sub>0.005</sub>
	20%	62.98 <sub>1.60</sub>	0.083 <sub>0.014</sub>	66.69 <sub>0.51</sub>	0.073 <sub>0.001</sub>	67.95 <sub>0.27</sub>	0.057 <sub>0.009</sub>
	30%	60.22 <sub>3.49</sub>	0.107 <sub>0.031</sub>	65.43 <sub>0.60</sub>	0.086 <sub>0.004</sub>	65.15 <sub>0.13</sub>	0.092 <sub>0.005</sub>
	40%	55.90 <sub>2.93</sub>	<b>0.135</b> <sub>0.028</sub>	63.68 <sub>1.20</sub>	<b>0.126</b> <sub>0.024</sub>	63.41 <sub>0.49</sub>	<b>0.110</b> <sub>0.018</sub>

Table 4: Comparing LoRA-FL with EAB-FL (Meerza & Liu, 2024) on Adult (Dua & Graff, 2017) with **race** as the sensitive attribute and 20% adversaries. Results for EAB-FL are taken from (Meerza & Liu, 2024, Table 1).

Attack	Accuracy (FedAvg/ Attack) ( $\uparrow$ )	$\Delta_{EO}$ (FedAvg/ Attack) ( $\downarrow$ )	$\Delta_{DP}$ (FedAvg/ Attack) ( $\downarrow$ )
EAB-FL (Meerza & Liu, 2024)	83.00 / 80.00	0.25 / 0.41	0.27 / 0.44
LoRA-FL (ours)	85.02 <sub>0.187</sub> / 84.52 <sub>0.219</sub>	0.09 <sub>0.025</sub> / 0.16 <sub>0.065</sub>	0.09 <sub>0.013</sub> / 0.13 <sub>0.022</sub>

shows that LoRA-FL remains effective in image-based FL with convolutional architectures, even under varying degrees of client heterogeneity controlled by  $\rho$ . As the fraction of adversarial clients increases, LoRA-FL consistently amplifies Accuracy Parity (AP) violations across all values of  $\rho$ , while accuracy degrades more gradually. Notably, lower values of  $\rho$ , corresponding to more non-IID client distributions, lead to systematically larger AP gaps for both FedAvg and FLAME, indicating that statistical heterogeneity further exacerbates fairness vulnerabilities in CNN-based models.

**Comparison with EAB-FL (Meerza & Liu, 2024).** From Table 4, on Adult (with **race** as the sensitive attribute and 20% adversaries), LoRA-FL amplifies fairness gaps at par with EAB-FL: LoRA-FL increases  $\Delta_{EO}$  and  $\Delta_{DP}$  by 77.8% and 44.4% over its baseline, while EAB-FL yields relative increases of 64.0% and 63.0%, respectively. Crucially, LoRA-FL incurs only a 0.50 accuracy drop (85.02  $\rightarrow$  84.52), compared to EAB-FLs significant drop (83.00  $\rightarrow$  80.00). These results demonstrate that LoRA-FL degrades fairness more effectively with negligible impact on utility.

**Other Ablations.** We conduct an additional ablation study on: (i) the number of agents, (ii) the algorithm of SPECTRE (Hayase et al., 2021) when implemented on LoRA-FL, and (iii) adapter rank. Additionally, we investigate the effect of omitting Phase 1 in LoRA-FL (i.e., when the adversary does not train for accuracy-specific adapters), with results presented in Appendix A. The key takeaways are as follows: *First*, removing Phase 1 disrupts the accuracy-fairness balance that LoRA-FL maintains, causing significant accuracy loss

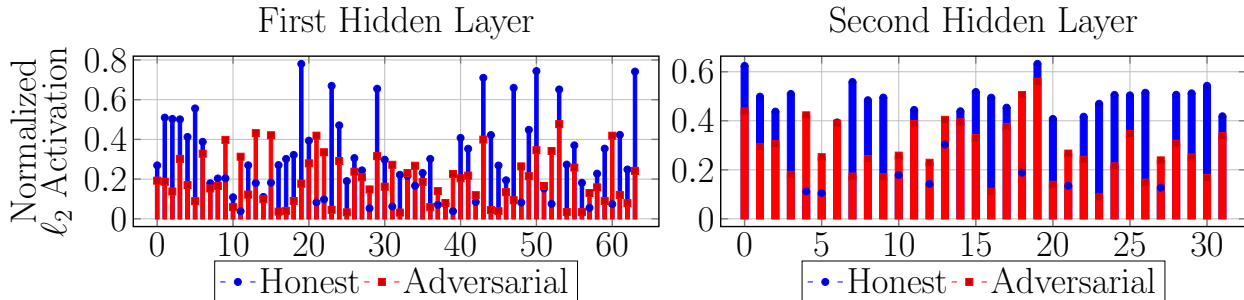


Figure 3: Activation values of neurons in the first and second fully connected layers under honest (blue) and 40% adversarial (red) setup. The plots reveal a higher correlation between activations corresponding to ‘Male’ and ‘Female’ inputs in the honest setting, which diminishes under adversarial perturbations, indicating a disruption in representational consistency.

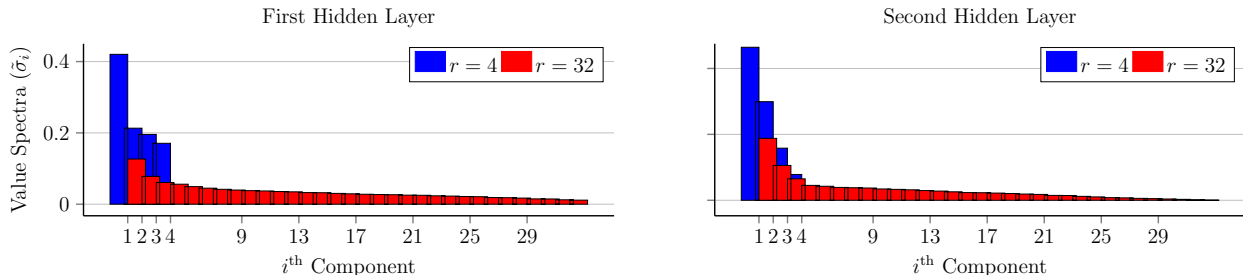


Figure 4: Singular value spectra of  $AB^T$  for two hidden layers on Adult. Although 32-dimensional, the adapters concentrate mass along 3–4 dominant directions, yielding an effective rank of approximately 4.

while still leading to substantial fairness degradation. *Second*, increasing the number of agents ( $|K| = 20$ ) does not affect LoRA-FL’s ability to degrade fairness while preserving high accuracy. *Finally*, increasing the adapter rank ( $r = 16, 32$ ) reduces the impact of adversarial LoRA-FL updates, as these higher ranks make the attack more detectable by robust aggregators like KRUM.

## 6 Interpreting the Role of Adapters in LoRA-FL

Here, we investigate three key questions regarding the role of adapters in LoRA-FL:

- Q1. How do adapters help adversarial updates evade filtering by parameter-based robust aggregators such as KRUM?
- Q2. How do adapters amplify bias in the models predictions?, and
- Q3. To what extent do adapters of different ranks span similar subspaces?

**Q1: Escaping the KRUM Trap.** Observe that any (adversarial) parameter update to the local model  $\theta$  that remains sufficiently close to the global model  $\Theta$  can evade filtering by robust aggregation methods such as KRUM. In Phase 2 of Algorithm 1, we exploit this by injecting unfairness into the low-rank adapters, thereby biasing local predictions while keeping parameter deviations within an acceptable range of  $\Theta$ .

*Frequency of Adversary Selection.* To assess the role of low-rank adaptation in the success of LoRA-FL, we vary the adapter rank  $r \in \{2, 4, 8, 16, 32\}$  and measure the frequency with which KRUM selects adversarial updates over  $T = 50$  communication rounds (same setup as Adult). As the rank increases, KRUM becomes significantly more effective at filtering adversarial updates: the selection frequency drops monotonically from  $0.94 \pm 0.01$  at rank 2 and  $0.71 \pm 0.02$  at rank 4, to  $0.52 \pm 0.03$  at rank 8,  $0.23 \pm 0.06$  at rank 16, and only  $0.01 \pm 0.02$  at rank 32. This highlights that a low-rank structure is critical for evading KRUM.

Motivated by this observation, we seek to understand why lower-rank adapters succeed where higher-rank ones fail. We track the  $\ell_2$  distance  $\|\Theta_t - \theta_{i,t}\|_2$  between an adversarial agent  $i$  and the global model over  $t \in [T]$  for different ranks (Fig. 2). At lower ranks (8 and below), adversarial updates remain close to or below the deviation exhibited by honest agents, whereas at higher ranks (16 and 32) the deviation grows substantially. Consequently, KRUM can reliably detect and filter higher-rank adversarial updates, as their parameter shifts exceed the tolerance required to evade robust aggregation.

**Q2: Interpreting the Role of Adapters in Amplifying Bias.** We evaluate the impact of LoRA-FL on demographic representational alignment in an FL setup with (i) only honest and (ii) 40% adversarial agents. Using the Adult dataset test split, we extract neuron activations from the two fully connected layers for Male and Female examples. For each layer  $l$  of width  $d$ , where  $d$  represents the number of neurons, we compute the  $\ell_2$ -mean activation vectors  $\mu_l^M, \mu_l^F \in \mathbb{R}^d$  for the respective gender subsets. These vectors are then normalized element-wise by the maximum activation value across both groups:  $\tilde{\mu}_l^M = \frac{\mu_l^M}{\text{norm}}$ ,  $\tilde{\mu}_l^F = \frac{\mu_l^F}{\text{norm}}$ , where  $\text{norm} := \max(\max_i(\mu_l^M)_i, \max_i(\mu_l^F)_i)$  represents the maximum activation across both vectors. Finally, we compute the element-wise product  $\mathbf{s}_l = \tilde{\mu}_l^M \odot \tilde{\mu}_l^F$ , which serves as the per-neuron co-activation score. The score  $\mathbf{s}_l$  measures the neuron-wise correlation between the activations for the two demographic groups. Higher values of  $\mathbf{s}_l$  indicate that neurons respond similarly to both groups, while lower values suggest more divergent responses.

Figure 3 presents this correlation for the two hidden layers of  $\Theta$  under (i) honest FL and (ii) 40% adversarial FL, where each neuron is plotted with its corresponding co-activation score, highlighting the degree of alignment between the groups at the neuron level. In the honest condition,  $\mathbf{s}_l$  remains uniformly high across neurons, demonstrating that adapters preserve representational consistency between Male and Female inputs. Under adversarial perturbations, however,  $\mathbf{s}_l$  drops markedly, revealing that adversaries induce divergent neuron activation patterns and disparate processing of demographic attributes. We see that LoRA-FL amplifies bias by degrading neuron-level alignment between gendered representations, disrupting the models ability to maintain consistent internal processing across demographic groups.

**Q3: Subspace Similarity for Different  $r$ .** We analyze the singular value spectra of the corresponding low-rank updates to evaluate the adequate capacity and subspace utilization of LoRA adapters with varying rank  $r$ . Specifically, given learned matrices  $A \in \mathbb{R}^{d \times r}$  and  $B \in \mathbb{R}^{k \times r}$  from LoRA adapters trained on the **Adult** dataset, we compute the singular values of the matrix product  $AB^T \in \mathbb{R}^{d \times k}$  via singular value decomposition (SVD) (Golub & Van Loan, 2013). That is,  $AB^T = U\Sigma V^T$ , where  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_{\min(d,k)})$  contains the singular values in descending order. We normalize the singular values such that they sum to one, i.e., we analyze  $\tilde{\sigma}_i = \frac{\sigma_i}{\sum_j \sigma_j}$ , for all  $i$ .

Figure 4 presents the spectra of the normalized singular values  $\tilde{\sigma}_i$  for the two hidden layers, for adapters with rank  $r = 4$  and  $r = 32$ . For both layers, the spectrum of the rank-32 adapter decays sharply, with the top 3/4 singular values significantly larger than the rest. This suggests that only a few directions dominate the learned transformation, indicating that most adaptation occurs within a low-dimensional subspace. This observation aligns with the empirical success of low-rank LoRA-FL attacks: the rank-4 adapter approximates the key directions of the rank-32 counterpart, enabling comparable degradation in fairness while maintaining low parameter deviation from  $\Theta$ .

## 7 Conclusion & Future Work

We present LoRA-FL, a low-rank fairness manipulation attack that exposes a fundamental limitation of robustness mechanisms in federated learning. Our evaluation spans both IID and non-IID data distributions, multiple aggregation rules including FedAvg, KRUM, TM, and FLAME, and model architectures ranging from MLPs for tabular data to convolutional networks for image classification. Across these configurations, we observe that strong predictive performance and convergence do not imply equitable treatment across sensitive groups. An important direction for future work is extending defenses to reason about structured, low-rank update geometry may help close the gap exploited by LoRA-FL.

## References

- Alekh Agarwal, Alina Beygelzimer, Miroslav Dudik, John Langford, and Hanna Wallach. A reductions approach to fair classification. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 60–69, Stockholm, Sweden, 10–15 Jul 2018. PMLR. URL <http://proceedings.mlr.press/v80/agarwal18a.html>.
- Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Machine bias: Theres software used across the country to predict future criminals. and its biased against blacks. *ProPublica*, 2016. URL <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>. Accessed: 2025-04-29.
- Solon Barocas and Andrew D Selbst. Big data’s disparate impact. *Cal. L. Rev.*, 104:671, 2016.
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning: Limitations and Opportunities*. MIT Press, 2023.
- Arjun Nitin Bhagoji, Supriyo Chakraborty, Prateek Mittal, and Seraphin Calo. Analyzing federated learning through an adversarial lens. In *International conference on machine learning*, pp. 634–643. PMLR, 2019.
- Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector machines. In *Proceedings of the 29th International Conference on Machine Learning*, pp. 1467–1474, 2012.
- M. Bilal Zafar, I. Valera, M. Gomez Rodriguez, and K. P. Gummadi. Fairness Constraints: Mechanisms for Fair Classification. *ArXiv e-prints*, July 2015.
- Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. Machine learning with adversaries: Byzantine tolerant gradient descent. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, 2017a. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/f4b9ec30ad9f68f89b29639786cb62ef-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/f4b9ec30ad9f68f89b29639786cb62ef-Paper.pdf).
- Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. Machine learning with adversaries: Byzantine tolerant gradient descent. *Advances in neural information processing systems*, 30, 2017b.
- Xiaoyu Cao, Minghong Fang, Jia Liu, and Neil Zhenqiang Gong. Fltrust: Byzantine-robust federated learning via trust bootstrapping. *arXiv:2012.13995*, 2020.
- Hongyan Chang and Reza Shokri. Bias propagation in federated learning. In *The Eleventh International Conference on Learning Representations*.
- Alexandra Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5 2:153–163, 2017.
- Yanbo Dai and Songze Li. Chameleon: Adapting to peer images for planting durable backdoors in federated learning. In *International Conference on Machine Learning*, pp. 6712–6725. PMLR, 2023.
- Chuntao Ding, Xu Cao, Jianhang Xie, Linlin Fan, Shangguang Wang, and Zhichao Lu. Lora-c: Parameter-efficient fine-tuning of robust cnn for iot devices. *arXiv preprint arXiv:2410.16954*, 2024.
- Dheeru Dua and Casey Graff. Uci machine learning repository, 2017. URL <https://archive.ics.uci.edu/ml/datasets/adult>.
- Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pp. 214–226, 2012.

- Yahya H Ezzeldin, Shen Yan, Chaoyang He, Emilio Ferrara, and A Salman Avestimehr. Fairfed: Enabling group fairness in federated learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 7494–7502, 2023.
- Alessandro Fabris, Nina Baranowska, Matthew J Dennis, David Graus, Philipp Hacker, Jorge Saldivar, Fredrik Zuiderveen Borgesius, and Asia J Biega. Fairness and bias in algorithmic hiring: A multidisciplinary survey. *ACM Transactions on Intelligent Systems and Technology*, 16(1):1–54, 2025.
- Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Gong. Local model poisoning attacks to {Byzantine-Robust} federated learning. In *29th USENIX security symposium (USENIX Security 20)*, pp. 1605–1622, 2020.
- Sorelle A Friedler, Carlos Scheidegger, Suresh Venkatasubramanian, Sonam Choudhary, Evan P Hamilton, and Derek Roth. A comparative study of fairness-enhancing interventions in machine learning. In *Proceedings of the conference on fairness, accountability, and transparency*, pp. 329–338, 2019.
- Jiashi Gao, Ziwei Wang, Xiangyu Zhao, Xin Yao, and Xuetao Wei. Pfattack: Stealthy attack bypassing group fairness in federated learning. *arXiv preprint arXiv:2410.06509*, 2024.
- Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013.
- Andrew Hard, Kanishka Rao, Rajiv Mathews, Swaroop Ramaswamy, Françoise Beaufays, Sean Augenstein, Hubert Eichner, Chloé Kiddon, and Daniel Ramage. Federated learning for mobile keyboard prediction. *arXiv preprint arXiv:1811.03604*, 2018.
- Moritz Hardt, Eric Price, and Nathan Srebro. Equality of opportunity in supervised learning. In *NIPS*, 2016a.
- Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. *Advances in neural information processing systems*, 29, 2016b.
- Jonathan Hayase, Wei Wei, Kyle Galyardt, and Sewoong Oh. Spectre: Defending against backdoor attacks using robust statistics. In *International Conference on Machine Learning*, pp. 4129–4139. PMLR, 2021.
- Jialuo He, Wei Chen, and Xiaojin Zhang. Fedaa: A reinforcement learning perspective on adaptive aggregation for fair and robust federated learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(16):17085–17093, Apr. 2025. doi: 10.1609/aaai.v39i16.33878. URL <https://ojs.aaai.org/index.php/AAAI/article/view/33878>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Bo Li, Yining Wang, Aarti Singh, and Yevgeniy Vorobeychik. Data poisoning attacks on factorization-based collaborative filtering. *Advances in neural information processing systems*, 29, 2016.
- Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning through personalization. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 6357–6368. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/li21h.html>.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *International Conference on Learning Representations (ICLR)*, 2019. arXiv:1711.05101.
- David Madras, Elliot Creager, Toniann Pitassi, and Richard S. Zemel. Learning adversarially fair and transferable representations. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, pp. 3381–3390, 2018. URL <http://proceedings.mlr.press/v80/madras18a.html>.

- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.
- Syed Irfan Ali Meerza and Jian Liu. Eab-fl: exacerbating algorithmic bias through model poisoning attacks in federated learning. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pp. 458–466, 2024.
- Ninareh Mehrabi, Muhammad Naveed, Fred Morstatter, and Aram Galstyan. Exacerbating algorithmic bias through fairness attacks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 8930–8938, 2021.
- Sérgio Moro, P. Cortez, and Paulo Rita. A data-driven approach to predict the success of bank telemarketing. *Decis. Support Syst.*, 62:22–31, 2014. URL <https://api.semanticscholar.org/CorpusID:14181100>.
- Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, pp. 807–814, 2010.
- Kai Wang Ng, Guo-Liang Tian, and Man-Lai Tang. *Dirichlet and Related Distributions: Theory, Methods and Applications*. Wiley Series in Probability and Statistics. John Wiley & Sons, Chichester, UK, 2011. ISBN 978-0-470-68819-9.
- Thien Duc Nguyen, Phillip Rieger, Huili Chen, Hossein Yalame, Helen Möllering, Hossein Fereidooni, Samuel Marchal, Markus Miettinen, Azalia Mirhoseini, Shaza Zeitouni, et al. {FLAME}: Taming backdoors in federated learning. In *31st USENIX Security Symposium (USENIX Security 22)*, pp. 1415–1432, 2022.
- Manisha Padala and Sujit Gujar. Fnnc: Achieving fairness through neural networks. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence. IJCAI*. <https://www.ijcai.org/proceedings/2020/0315.pdf> Go to original source, 2020.
- European Parliament and Council of the European Union. General data protection regulation. <https://eur-lex.europa.eu/eli/reg/2016/679/oj>, 2016. Regulation (EU) 2016/679.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: an imperative style, high-performance deep learning library. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, 2019.
- Dana Pessach and Erez Shmueli. A review on fairness in machine learning. *ACM Computing Surveys (CSUR)*, 55(3):1–44, 2022.
- Ricardo Trainotti Rabonato and Lilian Berton. A systematic review of fairness in machine learning. *AI and Ethics*, 5(3):1943–1954, 2025.
- Benjamin IP Rubinstein, Blaine Nelson, Ling Huang, Anthony D Joseph, Shing-hon Lau, Satish Rao, Nina Taft, and J Doug Tygar. Antidote: understanding and defending against poisoning of anomaly detectors. In *Proceedings of the 9th ACM SIGCOMM Conference on Internet Measurement*, pp. 1–14, 2009.
- Teresa Salazar, Helder Ara-Áejo, Alberto Cano, and Pedro Henriques Abreu. A survey on group fairness in federated learning: challenges, taxonomy of solutions and directions for future research. *arXiv preprint arXiv:2410.03855*, 2024.
- Virat Shejwalkar and Amir Houmansadr. Manipulating the byzantine: Optimizing model poisoning attacks and defenses for federated learning. In *NDSS*, 2021.
- Micah J Sheller, Brandon Edwards, G Anthony Reina, Jason Martin, Sarthak Pati, Aikaterini Kotrotsou, Mikhail Milchenko, Weilin Xu, Daniel Marcus, Rivka R Colen, et al. Federated learning in medicine: facilitating multi-institutional collaborations without sharing patient data. *Scientific reports*, 10(1):12598, 2020.

- Jinhyun So, Başak Güler, and A Salman Avestimehr. Byzantine-resilient secure federated learning. *IEEE Journal on Selected Areas in Communications*, 2020.
- David Solans, Battista Biggio, and Carlos Castillo. Poisoning attacks on algorithmic fairness. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 162–177. Springer, 2020.
- Minh-Hao Van, Wei Du, Xintao Wu, and Aidong Lu. Poisoning attacks on fair machine learning. In *International Conference on Database Systems for Advanced Applications*, pp. 370–386. Springer, 2022.
- Ganghua Wang, Ali Payani, Myungjin Lee, and Ramana Rao Kompella. Mitigating group bias in federated learning: Beyond local fairness. *Transactions on Machine Learning Research*.
- Geming Xia, Jian Chen, Chaodong Yu, and Jun Ma. Poisoning attacks in federated learning: A survey. *Ieee Access*, 11:10708–10722, 2023.
- Jie Xu, Benjamin S Glicksberg, Chang Su, Peter Walker, Jiang Bian, and Fei Wang. Federated learning for healthcare informatics. *Journal of healthcare informatics research*, 5:1–19, 2021.
- Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2):1–19, 2019.
- Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. Byzantine-robust distributed learning: Towards optimal statistical rates. In *International conference on machine learning*, pp. 5650–5659. Pmlr, 2018.
- Xubo Yue, Maher Nouiehed, and Raed Al Kontar. Gifair-fl: A framework for group and individual fairness in federated learning. *INFORMS Journal on Data Science*, 2(1):10–23, 2023.
- Zhifei Zhang, Yang Song, and Hai Qi. Utkface: Age progression/regression by conditional adversarial autoencoder, 2017. URL <https://susanqq.github.io/UTKFace/>. Large Scale Face Dataset.
- Indre Ilić, Faisal Kamiran, and Toon Calders. Handling conditional discrimination. In *2011 IEEE 11th International Conference on Data Mining*, pp. 992–1001, 2011.

Table A.1: Hyperparameters

Parameter	Symbol	Value	Description
Number of agents	$K$	10 or 20	Total number of agents in the system
Number of rounds	$T$	50	Total number of communication rounds
Local epochs	$E$	10 (Adult, Bank and Dutch), 4 (UTKFace)	Number of local training epochs for honest agents
Adversarial epochs	$E_A$	10 (Adult and Dutch), 20 (Bank)	Number of local training epochs for adversarial agents, 3 (UTKFace)
Agents per round	$m$	60% of $K$	Number of sampled agents per round
agent optimizer	$OPT$	AdamW	Optimizer used by honest agents
Adversarial optimizer	$OPT_{REG}$	AdamW	Optimizer used by adversarial agents for Phase 1
Adversarial optimizer	$OPT_F$	AdamW	Optimizer used by adversarial agents for Phase 2
Scaling factor	$\alpha$	1.0	Scaling applied to the low-rank update
Aggregator function	$Agg$	FedAvg or KRUM	Aggregation rule
Batch Size	$B$	512 (Adult, Bank and Dutch), 1024 (UTKFace)	
Honest agents Learning Rate	$\eta$	5e-4 (Adult, Bank and Dutch), 1e-4 (UTKFace)	AdamW LR for Honest agents
Adversarial agents Learning Rate	$\eta_{REG}$	5e-4 (Adult, Bank and Dutch), 3e-4 (UTKFace)	AdamW LR for Adversarial agents
Adversarial agents Learning Rate	$\eta_F$	5e-4	AdamW LR for Adversarial agents
Rank	$r$	4 (Adult and Dutch), 2 (Bank and UTKFace)	Rank of the Adversarial Adapters
Scaling Factor	$\alpha$	1	Controls the scale by which the adapters are fused in MLP
Scaling Factor	$\alpha'$	4	Controls the scale by which the adapters are fused in CNN

## A Training Details & Additional Experiments

### A.1 Hyperparameter Details

Table A.1 summarizes the key hyperparameters used for LoRA-FL. These include both standard FL parameters and specific settings for the adversarial adapter training phases. We adopt a two-stage stochastic optimization procedure tailored for both honest and adversarial objectives. All models are trained using mini-batch stochastic gradient descent with the AdamW optimizer, using a batch size of 512 for Adult, Bank, and Dutch datasets and 1024 for UTKFace and a fixed learning rate of 5e-4 (Adult, Bank, and Dutch) and 3e-4 (UTKFace). For standard (honest) training, we minimize the binary cross-entropy loss over ten local epochs in Adult, Bank, and Dutch, and four local epochs in the case of the UTKFace dataset. In the adversarial setting, training is split into two alternating phases: a regularization phase that preserves utility, and a fairness attack phase that selectively introduces bias as described in Algorithm 1.

### A.2 LoRA-FL: Training without Phase 1

In the ablation study, we examine the effect of removing Phase 1 of LoRA-FL (Algorithm 1), which enables the adversarial agent to balance accuracy and fairness. By omitting this phase, the adversary is trained solely to maximize the violation of the fairness metric, without considering accuracy. The results from this setup are reported on the Adult dataset. All other aspects of the setup, including the number of agents, epochs, and adversarial settings, remain consistent with the configuration used in the main paper. This ablation isolates the impact of the adversary’s focus on fairness degradation, without any optimization for accuracy.

Table A.2 presents the results for the ablation. As shown in Table A.2, removing Phase 1 significantly disrupts the accuracy-fairness trade-off that LoRA-FL is designed to maintain. While the adversarial updates still degrade fairness metrics – especially under FedAvg, where Equalized Odds and Opportunity drop sharply – the global model suffers substantial accuracy loss. For instance, with 30% adversaries, accuracy drops to 77.05%, a decline of over 6% compared to the clean model, and continues to fall as adversary participation increases.

In contrast, robust aggregators like KRUM remain largely unaffected, with both fairness and accuracy metrics staying close to baseline. Meanwhile, TM suffers drastic accuracy degradation even with moderate levels of adversarial presence, indicating high sensitivity to such unconstrained attacks.

These results demonstrate that omitting Phase 1 removes the stealth component of LoRA-FL: although the attack remains effective at harming fairness, the resulting perturbations become too aggressive, making them

Table A.2: **Ablation Study: LoRA-FL without Phase 1.** The adversarial agents omit training for accuracy. We compare the performance on different aggregators on accuracy and fairness metrics: Demographic Parity (DP), Equalized Odds (EO), and Equal Opportunity (EOpp). Here,  $|K| = 10$  and  $r = 4$ .

Aggregator	% Adversary	Accuracy ( $\uparrow$ )	$\Delta_{EO}$ ( $\downarrow$ )	$\Delta_{EOpp}$ ( $\downarrow$ )	$\Delta_{DP}$ ( $\downarrow$ )
<b>Adult Dua &amp; Graff (2017)</b>					
FedAvg	–	83.33 <sub>0.38</sub>	0.122 <sub>0.008</sub>	0.095 <sub>0.010</sub>	0.186 <sub>0.007</sub>
KRUM	–	82.56 <sub>0.20</sub>	0.115 <sub>0.016</sub>	0.082 <sub>0.018</sub>	0.190 <sub>0.008</sub>
TM	–	83.10 <sub>0.92</sub>	0.115 <sub>0.023</sub>	0.088 <sub>0.017</sub>	0.183 <sub>0.034</sub>
FedAvg	10%	81.73 <sub>0.53</sub>	0.448 <sub>0.034</sub>	0.446 <sub>0.036</sub>	0.330 <sub>0.023</sub>
FedAvg	20%	78.41 <sub>1.51</sub>	0.594 <sub>0.097</sub>	0.594 <sub>0.097</sub>	0.449 <sub>0.028</sub>
FedAvg	30%	77.05 <sub>1.58</sub>	0.693 <sub>0.059</sub>	0.693 <sub>0.059</sub>	0.475 <sub>0.038</sub>
FedAvg	40%	76.26 <sub>1.18</sub>	0.693 <sub>0.033</sub>	0.693 <sub>0.033</sub>	0.476 <sub>0.030</sub>
KRUM	10%	82.35 <sub>0.32</sub>	0.117 <sub>0.002</sub>	0.076 <sub>0.012</sub>	0.193 <sub>0.004</sub>
KRUM	20%	82.38 <sub>0.46</sub>	0.123 <sub>0.023</sub>	0.096 <sub>0.026</sub>	0.188 <sub>0.017</sub>
KRUM	30%	82.38 <sub>0.31</sub>	0.121 <sub>0.012</sub>	0.088 <sub>0.011</sub>	0.188 <sub>0.010</sub>
KRUM	40%	82.49 <sub>0.26</sub>	0.115 <sub>0.008</sub>	0.082 <sub>0.017</sub>	0.182 <sub>0.012</sub>
TM	10%	35.12 <sub>11.17</sub>	0.094 <sub>0.091</sub>	0.032 <sub>0.025</sub>	0.103 <sub>0.109</sub>
TM	20%	35.86 <sub>18.87</sub>	0.121 <sub>0.203</sub>	0.097 <sub>0.167</sub>	0.138 <sub>0.233</sub>
TM	30%	34.20 <sub>14.34</sub>	0.142 <sub>0.235</sub>	0.089 <sub>0.152</sub>	0.149 <sub>0.245</sub>
TM	40%	52.21 <sub>15.56</sub>	0.475 <sub>0.251</sub>	0.431 <sub>0.296</sub>	0.365 <sub>0.139</sub>

more detectable by robust defenses and compromising the models utility. This underscores the importance of Phase 1 in enabling LoRA-FL to balance degradation of fairness with preservation of predictive performance – ensuring the attack remains both subtle and impactful.

### A.3 Agents

Table A.3 presents results for a larger agent pool ( $|K| = 20$ ), showing that our low-rank adapter attack remains effective at degrading fairness, while maintaining high accuracy. For the Adult dataset, we observe that as the fraction of adversarial agents increases, fairness metrics such as  $\Delta_{EO}$  and  $\Delta_{DP}$  degrade substantially. For instance,  $\Delta_{EO}$  rises from 0.177 (clean) to 0.621 with 40% adversaries under FedAvg, a  $3.5\times$  increase, while accuracy drops by 3.7%. Importantly, the trends mirror those in the main paper for  $|K| = 10$ , confirming that LoRA-FL (with  $r = 4$ ) induces fairness degradation in a manner largely agnostic to the number of agents in the system.

### A.4 Adapter Rank

Table A.4 shows that increasing the adapter rank ( $r = 16, 32$ ) significantly mitigates the impact of adversarial LoRA-FL updates for both FedAvg and KRUM. In contrast to the strong degradation observed at lower ranks (i.e., for  $r = 4$  in the main paper), higher-rank adapters result in only marginal drops in accuracy and fairness, even under 40% adversarial agents. Notably, with  $r = 32$ , KRUM retains near-baseline fairness levels across all metrics, highlighting that low-rank constraints are a key enabler of the attacks potency by making adversarial updates more easily obscured or entangled in the parameter space.

Moreover, we observe that for KRUM, the **standard deviations of the fairness metrics across adversary proportions** (0%40%) are substantially lower at  $r = 32$  than at  $r = 16$ , indicating more stable behavior due to KRUMs ability to effectively filter out high-rank adversarial updates. Specifically, the standard deviation drops from 0.036 to 0.005 for  $\Delta_{EO}$ , 0.041 to 0.009 for  $\Delta_{EOpp}$ , and 0.031 to 0.004 for  $\Delta_{DP}$ . These findings

Table A.3: **Ablation Study: Number of Clients.** Comparison of different aggregators on accuracy and fairness metrics: Demographic Parity (DP), Equalized Odds (EO), and Equal Opportunity (EOpp). Here,  $r = 4$  and  $|K| = 10$ , with results averaged over four independent runs. We observe that our attack remains effective regardless of the number of participating clients, indicating its robustness to varying client pool sizes.

Aggregator	% Adversary	Accuracy ( $\uparrow$ )	$\Delta_{EO}$ ( $\downarrow$ )	$\Delta_{EOpp}$ ( $\downarrow$ )	$\Delta_{DP}$ ( $\downarrow$ )
<b>Adult Dua &amp; Graff (2017)</b>					
FedAvg	–	85.848 $\pm$ 0.304	0.177 $\pm$ 0.0107	0.0842 $\pm$ 0.0408	0.200 $\pm$ 0.002
KRUM	–	85.698 $\pm$ 0.208	0.159 $\pm$ 0.009	0.0875 $\pm$ 0.002	0.199 $\pm$ 0.003
FedAvg	10%	85.649 $\pm$ 0.168	0.236 $\pm$ 0.0104	0.089 $\pm$ 0.008	0.232 $\pm$ 0.007
FedAvg	20%	85.026 $\pm$ 0.0462	0.317 $\pm$ 0.0376	0.100 $\pm$ 0.0027	0.278 $\pm$ 0.015
FedAvg	30%	83.926 $\pm$ 0.0724	0.488 $\pm$ 0.0532	0.117 $\pm$ 0.009	0.393 $\pm$ 0.101
FedAvg	40%	82.673 $\pm$ 0.266	0.621 $\pm$ 0.109	0.129 $\pm$ 0.005	0.377 $\pm$ 0.018
KRUM	10%	85.358 $\pm$ 0.185	0.176 $\pm$ 0.005	0.117 $\pm$ 0.0106	0.224 $\pm$ 0.005
KRUM	20%	84.700 $\pm$ 0.272	0.269 $\pm$ 0.0132	0.127 $\pm$ 0.0113	0.287 $\pm$ 0.030
KRUM	30%	83.982 $\pm$ 0.538	0.434 $\pm$ 0.118	0.174 $\pm$ 0.0192	0.304 $\pm$ 0.056
KRUM	40%	83.278 $\pm$ 0.758	0.529 $\pm$ 0.146	0.234 $\pm$ 0.0076	0.316 $\pm$ 0.0411

suggest that higher adapter rank amplifies parameter deviation, enabling distance-based defenses like KRUM to more reliably identify and discard malicious updates.

### A.5 Implementation of algorithm on SPECTRE aggregator

SPECTRE Hayase et al. (2021) is a robust aggregator that utilizes robust mean estimation via spectral signatures. The formal procedure is detailed in Algorithm A.1. Although originally designed to counter data poisoning, we apply SPECTRE to LoRA-FL to determine if it can detect and isolate the spectral spikes associated with malicious client updates. As shown in Table A.5, LoRA-FL effectively circumvents this defense. As rank decreases, LoRA-FL becomes more stealthier: the accuracy increases from  $78.70 \pm 5.68$  at rank 64 and  $81.64 \pm 1.07$  at rank 8, to  $82.02 \pm 0.52$  at rank 4 and  $84.88 \pm 0.73$  at rank 2.

---

#### Algorithm A.1 SPECTRE Aggregator

---

- 1: **Input:** Local updates  $\{\theta_k\}_{k=1}^K$ , estimated number of adversaries  $q$
  - 2:  $\mu \leftarrow \text{median}(\{\theta_k\}_{k=1}^K)$   $\triangleright$  Centering
  - 3:  $X \leftarrow [(\theta_k - \mu)]_{k=1}^K$
  - 4:  $U, \Sigma, V^\top \leftarrow \text{SVD}(X)$   $\triangleright$  Dimensionality reduction to top- $k$  components
  - 5:  $X_{proj} \leftarrow XV_{1:k}$
  - 6:  $\hat{\Sigma} \leftarrow \frac{1}{K} X_{proj}^\top X_{proj}$   $\triangleright$  Projected covariance matrix
  - 7:  $\mathbf{W} \leftarrow \exp(\alpha \hat{\Sigma}) / \text{Tr}(\exp(\alpha \hat{\Sigma}))$   $\triangleright$  Compute Quantum Entropy weights
  - 8:  $\tau_k \leftarrow (\theta_k - \mu)^\top \mathbf{W} (\theta_k - \mu)$  for all  $k$   $\triangleright$  Spectral outlier scores
  - 9:  $\mathcal{S} \leftarrow \{\text{Indices of } K - q \text{ updates with smallest } \tau_k\}$
  - 10: **Output:**  $\Theta_{SPEC} = \frac{1}{|\mathcal{S}|} \sum_{k \in \mathcal{S}} \theta_k$
- 

## B Robust Aggregators

Here, we provide the formal algorithms for the robust aggregators.

Table A.4: **Ablation Study: Adapter Rank.** Comparison of different aggregators on accuracy and fairness metrics: Demographic Parity (DP), Equalized Odds (EO), and Equal Opportunity (EOpp). Here, the dataset used is **Adult** and  $|K| = 10$ . Higher adapter rank increases parameter deviation, enabling KRUM to more effectively filter adversarial updates and stabilize fairness metrics.

Aggregator	% Adversary	Accuracy ( $\uparrow$ )	$\Delta_{EO}$ ( $\downarrow$ )	$\Delta_{EOpp}$ ( $\downarrow$ )	$\Delta_{DP}$ ( $\downarrow$ )
$r = 16$					
FedAvg	–	84.910	0.118	0.096	0.175
KRUM	–	84.900	0.131	0.109	0.189
FedAvg	10%	84.830	0.141	0.113	0.207
FedAvg	20%	84.100	0.203	0.181	0.257
FedAvg	30%	82.850	0.305	0.297	0.312
FedAvg	40%	78.010	0.401	0.401	0.402
KRUM	10%	84.730	0.170	0.161	0.214
KRUM	20%	84.540	0.183	0.162	0.228
KRUM	30%	84.510	0.214	0.195	0.255
KRUM	40%	83.100	0.236	0.233	0.278
$r = 32$					
FedAvg	–	84.910	0.118	0.096	0.175
KRUM	–	84.900	0.131	0.109	0.189
FedAvg	10%	85.010	0.136	0.102	0.202
FedAvg	20%	83.810	0.211	0.184	0.261
FedAvg	30%	82.340	0.297	0.287	0.319
FedAvg	40%	78.640	0.160	0.160	0.067
KRUM	10%	84.840	0.132	0.090	0.188
KRUM	20%	84.910	0.141	0.114	0.197
KRUM	30%	84.730	0.128	0.100	0.193
KRUM	40%	84.660	0.128	0.091	0.196

---

**Algorithm B.1**  $m$ -KRUM Aggregation Blanchard et al. (2017b)

**Require:** Agent updates  $\{\theta_1, \theta_2, \dots, \theta_K\}$ , number of adversarial agents  $\tilde{q} = \lfloor qK \rfloor$ , agent sample sizes  $\{|\mathcal{X}_1|, \dots, |\mathcal{X}_K|\}$ , number of agents to aggregate  $m_{\min}$

**Ensure:** Aggregated Global Model  $\Theta^{\text{KRUM}}$

- 1: Let  $m := \max(K - \tilde{q}, m_{\min})$   $\triangleright$  Ensure  $m$  is at least  $K - \tilde{q}$
- 2: **for**  $i = 1$  to  $K$  **do**
- 3:     **Define**  $d_{i,j} = \|\theta_i - \theta_j\|_2$  as the Euclidean distance between pair-wise agent updates  $\theta_i$  and  $\theta_j$
- 4:     Compute distances  $d_{i,j}$  for all  $j \neq i$
- 5:     Let  $\mathcal{N}_i \leftarrow$  indices of  $K - \tilde{q} - 2$  closest updates to  $\theta_i$
- 6:     Compute score  $s_i = \sum_{j \in \mathcal{N}_i} d_{i,j}^2$
- 7: **end for**
- 8: Select the set  $\mathcal{M} \subseteq \{1, \dots, K\}$  of  $m$  agents with the lowest scores  $s_i$
- 9: Compute weighted average:

$$\Theta^{\text{KRUM}} = \sum_{i \in \mathcal{M}} \frac{|\mathcal{X}_i|}{\sum_{j \in \mathcal{M}} |\mathcal{X}_j|} \cdot \theta_i \quad (11)$$

10: **return**  $\Theta^{\text{KRUM}}$

---

Table A.5: SPECTRE accuracy and fairness metrics on the **Bank** dataset for both **IID** and **non-IID** settings. Here,  $r = 2$  with concentration parameters  $\rho \in \{0.25, 0.5, 0.75\}$  for non-IID settings. Also, **Acc**: Accuracy, **Adv**: Adversary, and we report  $\text{mean}_{\text{std}}$  across **four** independent runs. **Bold** values indicate highest accuracy and highest fairness violation.

Setting	% Adv	Acc ( $\uparrow$ )	$\Delta_{EO}$ ( $\downarrow$ )	$\Delta_{EOpp}$ ( $\downarrow$ )	$\Delta_{DP}$ ( $\downarrow$ )
IID	–	<b>91.01</b> <sub>0.25</sub>	0.22 <sub>0.02</sub>	0.15 <sub>0.03</sub>	0.23 <sub>0.02</sub>
	10%	90.92 <sub>0.11</sub>	0.21 <sub>0.01</sub>	0.14 <sub>0.02</sub>	0.25 <sub>0.02</sub>
	20%	90.48 <sub>0.22</sub>	0.25 <sub>0.03</sub>	0.15 <sub>0.02</sub>	0.30 <sub>0.03</sub>
	30%	90.33 <sub>0.17</sub>	0.26 <sub>0.02</sub>	0.16 <sub>0.03</sub>	0.30 <sub>0.02</sub>
	40%	88.64 <sub>0.40</sub>	<b>0.33</b> <sub>0.02</sub>	<b>0.17</b> <sub>0.02</sub>	<b>0.38</b> <sub>0.01</sub>
$\rho = 0.5$	–	<b>91.67</b> <sub>0.50</sub>	0.29 <sub>0.04</sub>	<b>0.28</b> <sub>0.03</sub>	0.23 <sub>0.01</sub>
	10%	91.41 <sub>0.45</sub>	0.29 <sub>0.03</sub>	0.25 <sub>0.03</sub>	0.25 <sub>0.02</sub>
	20%	91.33 <sub>0.43</sub>	0.32 <sub>0.04</sub>	0.26 <sub>0.03</sub>	0.28 <sub>0.05</sub>
	30%	90.95 <sub>0.66</sub>	0.32 <sub>0.05</sub>	0.25 <sub>0.01</sub>	0.30 <sub>0.06</sub>
	40%	90.18 <sub>1.18</sub>	<b>0.36</b> <sub>0.06</sub>	0.27 <sub>0.03</sub>	<b>0.33</b> <sub>0.06</sub>
$\rho = 0.25$	–	<b>89.90</b> <sub>1.30</sub>	0.28 <sub>0.04</sub>	0.26 <sub>0.05</sub>	0.23 <sub>0.02</sub>
	10%	88.46 <sub>2.87</sub>	0.36 <sub>0.12</sub>	0.28 <sub>0.08</sub>	0.33 <sub>0.07</sub>
	20%	86.81 <sub>3.88</sub>	0.38 <sub>0.13</sub>	0.30 <sub>0.10</sub>	0.33 <sub>0.08</sub>
	30%	53.84 <sub>36.34</sub>	0.46 <sub>0.16</sub>	<b>0.30</b> <sub>0.11</sub>	0.40 <sub>0.11</sub>
	40%	54.57 <sub>34.16</sub>	<b>0.49</b> <sub>0.12</sub>	0.26 <sub>0.00</sub>	<b>0.46</b> <sub>0.11</sub>
$\rho = 0.75$	–	<b>90.74</b> <sub>0.45</sub>	0.28 <sub>0.05</sub>	0.21 <sub>0.02</sub>	0.28 <sub>0.07</sub>
	10%	90.58 <sub>0.30</sub>	0.31 <sub>0.05</sub>	0.22 <sub>0.04</sub>	0.28 <sub>0.05</sub>
	20%	90.31 <sub>0.31</sub>	0.35 <sub>0.04</sub>	0.24 <sub>0.05</sub>	0.32 <sub>0.07</sub>
	30%	89.82 <sub>0.51</sub>	0.37 <sub>0.04</sub>	<b>0.26</b> <sub>0.04</sub>	0.35 <sub>0.06</sub>
	40%	88.73 <sub>0.45</sub>	<b>0.42</b> <sub>0.05</sub>	0.25 <sub>0.02</sub>	<b>0.40</b> <sub>0.06</sub>

---

### Algorithm B.2 $f$ -Trimmed-Mean Aggregation

---

**Require:** Agent updates  $\{\theta_1, \theta_2, \dots, \theta_K\}$ , number of values to trim  $f$

**Ensure:** Aggregated Global Model  $\Theta^{\text{TM}}$

- 1: **Assert:**  $K > 2f$  ▷ At least  $2f + 1$  agents required
  - 2: Initialize empty model  $\Theta^{\text{TM}}$
  - 3: **for** each parameter key  $w$  in model **do**
  - 4:   **if**  $w$  is a BatchNorm parameter **then**
  - 5:     Set  $\Theta^{\text{TM}}[w] \leftarrow \theta_1[w]$  ▷ Skip aggregation for BN layers
  - 6:     **continue**
  - 7:   **end if**
  - 8:   Stack agent parameters:  $V_w \leftarrow [\theta_1[w], \dots, \theta_K[w]]$  as matrix of shape  $(K, \text{param\_size})$
  - 9:   Sort  $V_w$  along agent dimension for each coordinate
  - 10:   Trim  $f$  smallest and  $f$  largest values at each coordinate
  - 11:   Compute coordinate-wise mean of trimmed values:  $\bar{v}_w$
  - 12:   Reshape  $\bar{v}_w$  to original shape and set  $\Theta^{\text{TM}}[w] \leftarrow \bar{v}_w$
  - 13: **end for**
  - 14: **return**  $\Theta^{\text{TM}}$
- 

### B.1 $m$ -KRUM

We employ the  $m$ -KRUM aggregation algorithm Blanchard et al. (2017b) to achieve robustness in the presence of adversarial agents. Given a set of local model updates  $\theta_1, \dots, \theta_K$  from  $K$  agents, and an upper bound  $\tilde{q} = \lfloor qK \rfloor$  on the number of potentially malicious agents, KRUM computes a robustness score  $s_i$  for each agent update  $\theta_i$  by summing the squared Euclidean distances to its  $K - \tilde{q} - 2$  closest peers. The  $m$  agents with the lowest scores are selected for aggregation, where  $m \geq K - \tilde{q}$ . A weighted average of these selected updates, using their respective sample sizes  $|\mathcal{X}_i|$ , yields the final global model  $\Theta^{\text{KRUM}}$ . By filtering out updates that are distant from the consensus, KRUM effectively limits the influence of Byzantine agents on the global model.

## B.2 Trimmed-Mean

We also employ the  $f$ -Trimmed-Mean aggregation algorithm Yin et al. (2018). Given a set of agent model updates  $\theta_1, \theta_2, \dots, \theta_K$ , the algorithm assumes that up to  $f$  of these may be adversarial and removes the  $f$  largest and  $f$  smallest values for each model parameter dimension independently. Specifically, for each parameter  $w$ , we collect all corresponding agent values, sort them element-wise, discard the extreme  $2f$  values, and take the mean of the remaining  $K - 2f$  entries to compute the aggregated parameter  $\bar{w}$ . This process is repeated for all parameters in the model. Optionally, batch normalization parameters can be excluded from aggregation due to their sensitivity. The resulting model  $\Theta^{\text{TM}}$  offers a robust estimate that mitigates the influence of malicious or corrupted updates.

## B.3 FLAME

FLAME (Federated Learning with Adversarial Model Evaluation) aggregation algorithm Nguyen et al. (2022), which is designed to provide robustness under both client heterogeneity and adversarial behavior. Given a set of client model updates  $\theta_1, \theta_2, \dots, \theta_K$  in a training round, FLAME first computes a pairwise similarity matrix using cosine similarity to capture directional agreement between updates while being invariant to scale. Based on this similarity matrix, FLAME applies a clustering procedure to partition updates into groups corresponding to benign and potentially malicious behavior. The cluster containing the largest number of updates is identified as the benign cluster, under the assumption that the majority of participating clients behave honestly. Updates that fall outside the benign cluster are either discarded or assigned reduced weights, depending on their similarity to the cluster centroid. The aggregator then computes the global update by averaging the selected benign updates. To further mitigate the influence of residual adversarial behavior and noisy updates—particularly under non-IID data distributions—FLAME injects adaptive Gaussian noise, calibrated to the empirical variance of the benign cluster. This noise injection step is intended to smooth the aggregated update and limit the impact of any remaining malicious contributions. The resulting global model, denoted  $\Theta^{\text{FLAME}}$ , provides robustness against Byzantine clients while preserving convergence and stability in heterogeneous federated learning settings.