A PROVABLY ROBUST ALGORITHM FOR DIFFEREN TIALLY PRIVATE CLUSTERED FEDERATED LEARNING

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026

031

Paper under double-blind review

ABSTRACT

Federated Learning (FL), which is a decentralized machine learning (ML) approach, often incorporates differential privacy (DP) to enhance data privacy guarantees. However, differentially private federated learning (DPFL) introduces performance disparities across clients, particularly affecting minority groups. Some recent works have attempted to address large data heterogeneity in vanilla FL settings through clustering clients, but these methods remain sensitive and prone to errors further exacerbated by the DP noise, making them inappropriate for DPFL settings. We propose an algorithm for differentially private clustered FL, which is robust to the DP noise in the system and identifies clients' clusters correctly. To this end, we propose to cluster clients based on both their model updates and training loss values. Furthermore, when clustering clients' model updates, our proposed approach addresses the server's uncertainties by employing large batch sizes as well as Gaussian Mixture Models (GMM) to reduce the impact of DP and stochastic noise and avoid potential clustering errors. This idea is efficient especially in privacy-sensitive scenarios with more DP noise. We provide theoretical analysis justifying our approach, and evaluate it extensively across diverse data distributions and privacy budgets. Our experimental results show its effectiveness in addressing large data heterogeneity in DPFL systems with a small computational cost.

1 INTRODUCTION

Federated learning (FL) (McMahan et al., 2017) is a collaborative ML paradigm, which allows multiple clients to train a shared global model without sharing their data. However, in order for
FL algorithms to ensure rigorous privacy guarantees against data privacy attacks (Hitaj et al., 2017; Rigaki & García, 2020; Wang et al., 2019; Zhu et al., 2019; Geiping et al., 2020), they are reinforced with DP (Dwork et al., 2006b;a; Dwork, 2011; Dwork & Roth, 2014). This is done in the presence of a trusted server (McMahan et al., 2018; Geyer et al., 2017) and in its absence (Zhao et al., 2020; Duchi et al., 2013; 2018). In the latter case and for record-level DP, each client adds noise to its stochastic gradients locally and shares its noisy model update with the server at the end of each round.

A key challenge in FL settings is ensuring a similar performance across clients under heterogeneous 040 data distributions, where several existing works focus on accuracy parity across clients with a single 041 common model (Mohri et al., 2019; Michieli & Ozay, 2021). However, a single global model often 042 fails to adapt to the data heterogeneity across clients (Chu et al., 2023), especially with extreme 043 *covariate and label shifts.* To address this, multiple methods were proposed to achieve performance 044 parity in non-DP FL settings: agnostic federated learning (Mohri et al., 2019), client reweighting (Li 045 et al., 2020b;a; Zhang et al., 2023), multi-task learning (Smith et al., 2017; Li et al., 2021; Marfoq 046 et al., 2021; Wu et al., 2023), transfer learning (Li & Wang, 2019; Liu et al., 2020) and clustered 047 FL (Ghosh et al., 2020; Mansour et al., 2020; Ruan & Joe-Wong, 2021; Sattler et al., 2019; Werner 048 et al., 2023; Briggs et al., 2020), where the latter is the focus of this work. On the other hand, when augmenting FL with DP for getting rigorous privacy guarantees, DP can have disparate impacts on the accuracy of different subgroups of clients - even with small imbalances and loose privacy 051 guarantees (Farrand et al., 2020; Fioretto et al., 2022; Bagdasaryan & Shmatikov, 2019). In fact, groups with minority data experience a larger drop in model utility (larger privacy cost). Being due to 052 the inequitable gradient clipping in DPSGD (Abadi et al., 2016; Bagdasaryan & Shmatikov, 2019; Xu et al., 2021; Esipova et al., 2022), this behavior has become increasingly important to be addressed.

054 As mentioned, clustered FL was proposed as an efficient personalization technique in vanilla FL for 055 performance parity under extreme data heterogeneity across clusters of clients: subsets of clients 056 are grouped together by the server based on their loss values (Ghosh et al., 2020; Mansour et al., 057 2020; Ruan & Joe-Wong, 2021; Chu et al., 2023; Liu et al., 2022) or their gradients (model updates) 058 (Sattler et al., 2019; Werner et al., 2023; Briggs et al., 2020). As discussed in (Werner et al., 2023) in details, the aforementioned two categories of clustered FL approaches are vulnerable to errors in clustering due to their sensitivity to: 1. model initialization 2. randomness in clients' model updates 060 due to stochastic noise. DP noise exacerbates this vulnerability, especially in the first few rounds of 061 FL training. To address this, we propose a clustered DPFL algorithm which uses both clients' model 062 updates and losses values to cluster them, making it more robust to DP/stochastic noise. 063

A correct clustering of clients results in equity of privacy cost between the client groups (Esipova et al., 2022; Tran et al., 2020). Justified by our theoretical analysis, our proposed algorithm uses a full batch size in the first round to reduce the noise in clients' model updates at the end of this round. Then, the server soft clusters clients based on these less noisy model updates using a Gaussian Mixture Model (GMM). Depending on the "confidence" of the learned GMM, the server keeps using it to soft cluster clients during the next few rounds. Finally, the server switches the clustering strategy to *local* clustering of clients based on their loss values in the remaining rounds. These altogether make our method effective and robust. The highlights of our contributions are as follows:

- We propose a DP clustered FL algorithm, which combines information from both clients' model updates and their loss values. The algorithm is robust and achieves high-quality clustering of clients, even in the presence of DP noise in the system.
- We theoretically prove that increasing clients' batch sizes, particularly in the initial communication round, consistently improves the server's ability to cluster clients based on their model updates at the end of the first round.
- We theoretically prove that using sufficiently large client batch sizes in the first round, enables super-linear convergence rate for learning a GMM on clients' model updates, which leads to fast and accurate clustering of clients with low computational overhead.
- Extensive evaluation across diverse and heterogeneous datasets and scenarios demonstrates the effectiveness of our robust clustered DPFL (RC-DPFL) algorithm in detecting the clustering structure of clients, which leads to a utility improvement for minority clusters.

2 RELATED WORK

073

074 075

076

077

078

079

081

082

084

Performance parity in FL: Performance parity of the final trained model across clients is an important goal in FL. Addressing this goal, Mohri et al. (2019) proposed Agnostic FL (AFL) by using a min-max optimization approach. TERM (Li et al., 2020a) used tilted losses to up-weight clients with large losses. Finally, Li et al. (2020b) and Zhang et al. (2023) proposed q-FFL and PropFair, inspired by α -fairness (Lan et al., 2010) and proportional fairness (Bertsimas et al., 2011), respectively. Generating one common model for all clients, these techniques do not perform well when the data distribution across clients is highly heterogeneous, leading to low overall performance in the system. This leads us to use stronger personalization techniques, e.g., client clustering.

095 Clustered FL: Clustered FL has been originally proposed for personalization in vanilla non-DP FL 096 with highly heterogeneous data, where clients can be naturally partitioned into clusters. Existing clustered FL algorithms cluster clients based on their loss values (Mansour et al., 2020; Ghosh et al., 098 2020; Ruan & Joe-Wong, 2021) or their model updates (based on e.g., their euclidean distance (Werner et al., 2023; Briggs et al., 2020) or cosine similarity (Sattler et al., 2019)). As studied in (Werner et al., 2023), the algorithms are prone to clustering errors in the early rounds of FL training 100 (due to gradient stochasticity, model initialization or the form of loss functions far from their optima), 101 which can even propagate in the subsequent rounds. This vulnerability is exacerbated in DPFL 102 systems, due to the extra DP noise. Without addressing this vulnerability, Luo et al. (2024) proposed 103 a clustered DPFL algorithm with a limited applicability, which clusters clients based on the labels 104 that they do not have in their local data, and is inapplicable when clients have all possible labels. 105

Differential privacy, group fairness and performance parity: Gradient clipping and random noise addition used in DPSGD disproportionately affect underrepresented groups. Some works tried to address the tension between group fairness and DP in centralized settings (Tran et al., 2020)

108 (by using Lagrangian duality) and FL settings (Pentyala et al., 2022) (by using Secure Multiparty Computation (MPC)). Another work tried to remove the disparate impact of DP on model performance 110 of minority groups in centralized settings (Esipova et al., 2022), by preventing gradient misalignment 111 across different groups of data. Unlike the previous works on group fairness, this work adopts cross-model fairness, where the performance cost of adding privacy to a non-private model must 112 be fairly distributed between different groups. We adopt the same notion - which is also used in 113 (Chu et al., 2023). Considering a highly heterogeneous data split, the mentioned approaches are 114 not appropriate due to generating one single model for all groups. In contrast, we propose a robust 115 "clustered" DPFL algorithm, which identifies different groups of clients and learns a model for each. 116

117 118

119

148

149

3 DEFINITIONS, NOTATIONS AND ASSUMPTIONS

120 There are multiple definitions of DP. We adopt the following definition in this work:

Definition 3.1 ((ϵ, δ)-DP (Dwork et al., 2006a)). A randomized mechanism $\mathcal{M} : \mathcal{D} \to \mathcal{R}$ with domain \mathcal{D} and range \mathcal{R} satisfies (ϵ, δ)-DP if for any two adjacent inputs $d, d' \in \mathcal{D}$, which differ only by a 123 single record (by removal), and for any measurable subset of outputs $\mathcal{S} \subseteq \mathcal{R}$ it holds that $Pr[\mathcal{M}(d) \in \mathcal{S}] \leq e^{\epsilon} Pr[\mathcal{M}(d') \in \mathcal{S}] + \delta.$

Gaussian mechanism randomizes the output of a query f as $\mathcal{M}(d) \triangleq f(d) + \mathcal{N}(0, \sigma^2)$. The randomized output of the Gaussian mechanism satisfies (ϵ, δ) -DP for a continuum of pairs (ϵ, δ) : it is

127 (ϵ, δ)-DP for all $\epsilon < 1$ and $\sigma > \frac{\sqrt{2 \ln(1.25/\delta)}}{\epsilon} \Delta_2 f$, where $\Delta_2 f \triangleq \max_{d,d'} || f(d) - f(d') ||_2$ is the 129 l_2 -sensitivity of the query f with respect to its input dataset. Also, the ϵ and δ privacy parameters 130 resulting from running Gaussian mechanism depend on the quantity $z = \frac{\sigma}{\Delta_2 f}$ (called "noise scale"). 131 We consider a DPFL system (see Figure 1, left), where there are n clients with the same desired 132 privacy parameters (ϵ, δ), and each runs DPSGD. In the context of Definition 3.1, we consider record-133 level (ϵ, δ)-DP for every client i: the set of model updates sent by client i to the server satisfies 134 (ϵ, δ)-DP (Definition 3.1) for all adjacent datasets \mathcal{D}_i and \mathcal{D}'_i differing in one record (by removal).

Let $x \in \mathcal{X} \subseteq \mathbb{R}^d$ and $y \in \mathcal{Y} = \{1, \dots, C\}$ denote an input data point and its target label. Client *i* 135 holds dataset \mathcal{D}_i with N_i samples from distribution $P_i(x, y) = P_i(y|x)P_i(x)$. Let $h: \mathcal{X} \times \boldsymbol{\theta} \to \mathbb{R}^C$ 136 be the predictor function, which is parameterized by $\theta \in \mathbb{R}^p$. Also, let $\ell : \mathbb{R}^C \times \mathcal{Y} \to \mathbb{R}_+$ 137 be the loss function used (cross-entropy loss). Client i in the system has empirical train loss 138 $f_i(\theta) = \frac{1}{N_i} \sum_{(x,y) \in \mathcal{D}_i} [\ell(h(x,\theta), y)]$, with minimum value f_i^* . There are E communication rounds 139 indexed by e. During each round e, client i runs K local epochs with learning rate η_l . There are M 140 clusters of clients indexed by m, and the server holds M cluster models $\{\theta_m^e\}_{m=1}^{n}$ for them at the 141 beginning of round e. Clients i and j belonging to the same cluster have the same data distributions, 142 while there is high data heterogeneity across clusters. s(i) denotes the true cluster of client i and 143 $R^{e}(i)$ denotes the cluster assigned to it at the beginning of round e. Let's assume the batch size that 144 client i uses in the first round e = 1 is b_i^1 , which may be different from the batch size $b_i^{>1}$ that it uses 145 in the rest of the rounds e > 1. At the t-th gradient update during the round e, and given a current 146 model θ , client *i* uses batch $\mathcal{B}_i^{e,t}$ with size b_i^e , and computes the following DP noisy batch gradient: 147

$$\tilde{g}_{i}^{e,t}(\boldsymbol{\theta}) = \frac{1}{b_{i}^{e}} \left[\left(\sum_{j \in \mathcal{B}_{i}^{e,t}} \bar{g}_{ij}(\boldsymbol{\theta}) \right) + \mathcal{N}(0, \sigma_{i, \text{DP}}^{2} \mathbb{I}_{p}) \right],$$
(1)

150 where $\bar{g}_{ij}(\boldsymbol{\theta}) = \operatorname{clip}(\nabla \ell(h(x_{ij}, \boldsymbol{\theta}), y_{ij}), c)$, and c is a clipping threshold: for a given vector \mathbf{v} , $\operatorname{clip}(\mathbf{v}, c) = \min\{\|\mathbf{v}\|, c\} \cdot \frac{\mathbf{v}}{\|\mathbf{v}\|}$. Also, \mathcal{N} is the Gaussian noise distribution with variance $\sigma_{i, \text{DP}}^2$, 151 152 where $\sigma_{i,\text{DP}} = c \cdot z_i(\epsilon, \delta, b_i^1, b_i^{>1}, N_i, K, E)$. z_i is the noise scale needed for achieving (ϵ, δ) -DP by 153 client *i*, which can be determined with a privacy accountant, e.g., the Renyi-DP accountant (Mironov 154 et al., 2019) used in this work, which is capable of accounting composition of heterogeneous DP 155 mechanisms (Mironov, 2017). The privacy parameter δ is fixed to 10^{-4} in this work. For an arbitrary random $\mathbf{v} = (v_1, \ldots, v_p)^\top \in \mathbb{R}^{p \times 1}$, we define $\operatorname{Var}(\mathbf{v}) := \sum_{j=1}^p \mathbb{E}[(v_j - \mathbb{E}[v_j])^2]$, i.e., variance of 156 157 \mathbf{v} is the sum of the variances of its elements. Table 1 in the appendix summarizes the used notations. 158 Finally, we have the following assumption:

Assumption 3.2. The stochastic gradient $g_i^{e,t}(\theta) = \frac{1}{b_i^e} \sum_{j \in \mathcal{B}_i^{e,t}} g_{ij}(\theta)$ is an unbiased estimate of $\nabla f_i(\theta)$ with a bounded variance: $\forall \theta \in \mathbb{R}^p : Var(g_i^{e,t}(\theta)) \leq \sigma_{i,g}^2(b_i^e)$. The tight bound $\sigma_{i,g}^2(b_i^e)$ is a constant depending only on the used batch size b_i^e : the larger the batch size b_i^e , the smaller $\sigma_{i,g}^2(b_i^e)$.



Figure 1: Left: Considered threat model in this work, where client i has local train data \mathcal{D}_i and DP privacy parameters (ϵ, δ) , and does not trust any external parties, including the server. **Right:** Three main stages of the proposed RC-DPFL algorithm, with the key components being highlighted.

174 175 176

177

172

173

164

167

4 MOTIVATION, METHODOLOGY AND PROPOSED ALGORITHM

178 We start with the shortcomings of the existing *non-DP* clustered FL algorithms. As discussed in 179 (Werner et al., 2023), algorithms clustering clients based on their loss values (Mansour et al., 2020; Ghosh et al., 2020; Ruan & Joe-Wong, 2021), i.e., assign client *i* to cluster $R^e(i) = \arg \min_m f_i(\boldsymbol{\theta}_m^e)$ 181 at the beginning of round e, are prone to clustering errors in the first few rounds, mainly due to 182 random initialization of cluster models $\{\theta_m^e\}_{m=1}^M$. On the other hand, clustering clients based on their 183 model updates (gradients) (Werner et al., 2023; Briggs et al., 2020; Sattler et al., 2019) makes sense only when the updates are obtained on the same model initialization. Additionally, even if we assume 185 these algorithms can initially cluster clients perfectly in each round e, the clients' model updates 186 (gradients) will approach zero as the clusters' models converge to their optimum parameters. Hence, 187 clients from different clusters may appear to belong to the same cluster, which results in clustering mistakes. To illustrate these shortcomings, we provide a more detailed example in Appendix C. 188

189 For the above mentioned reasons, we next propose an algorithm which starts with clustering clients 190 based on their model updates for the first several rounds and then switches its strategy to cluster 191 clients based on their loss values. We also augment this idea with some other non-obvious techniques 192 to enhance the clustering accuracy in the first few rounds, when the most clustering uncertainty exists.

193 194 195

197

199 200

201

202 203

204

205

206

207

208

209

210

4.1 RC-DPFL ALGORITHM

196 Considering the points above, which were overlooked in the existing *non-private* algorithms, we propose our *differentially private* RC-DPFL algorithm with the following steps (see Figure 1, right):

- Initializing clusters uniformly $(\forall m \in [M] : \boldsymbol{\theta}_{1}^{n} = \boldsymbol{\theta}^{init})$, clients use **full batch sizes** in the first round to make their model updates $\{\Delta \tilde{\boldsymbol{\theta}}_{i}^{1}\}_{i=1}^{n}$ less noisy. Then, the server **soft clusters** them by running GMM on their model updates. The remaining clustering uncertainties are incorporated in the probabilities returned by GMM ($\pi_{i,m}$).
- During the subsequent rounds $e \in \{2, \ldots, E_c\}$, the server uses the learned GMM to softcluster clients: client i contributes to the training of each cluster (m) model proportional to the probability of its assignment to that cluster $(\pi_{i,m})$. The duration E_c for using the GMM depends on the "confidence level" of the GMM.
- After the first E_c rounds, some progress has been made in the training of the cluster models $\{\theta_m^{E_c}\}_{m=1}^M$. Now, is the right time to hard cluster clients based on their loss values in the remaining rounds to build more personalized models per cluster: $R^{e}(i) = \arg \min_{m} f_{i}(\boldsymbol{\theta}_{m}^{e})$.

211 In Section 4.2.1, we provide theoretical justification for why using full batch sizes in the initial 212 round improves the clustering quality of GMM considerably. Also, in Section 4.3 we analyze the 213 convergence rate for learning the GMM and show that the computational overhead of using GMM is also low. Note that even when clients have a limited memory budget, they can still perform DPSGD 214 with full batch size using gradient accumulation technique (see Appendix I). The technique causes no 215 extra computational overhead, as it just accumulates multiple gradient updates into one update.

Algorithm 1: RC-DPFL **Input:** Initial parameter θ^{init} , number of clusters M, batch size b, dataset sizes $\{N_1, \ldots, N_n\}$, noise scales $\{z_1, \ldots, z_n\}$, gradient norm bound c, local epochs K, global round E. **Output:** cluster models $\{\theta_m^E\}_{m=1}^M$ 1 Initialize $\theta_1^1 = \ldots = \theta_m^1 = \theta^{init};$ // "uniform" initializations ² for $e \in \{1, ..., E\}$ do if e = 1 then for each client $i \in \{1, ..., n\}$ in parallel do $b_i^1 \leftarrow N_i$; // full batch size $\Delta \tilde{\boldsymbol{\theta}}_i^1 \leftarrow \texttt{DPSGD} \left(\boldsymbol{\theta}_i^1, b_i^1, N_i, K, z_i, c \right)$ on server: if M is unknown then $M = \arg\max_{m} \operatorname{MSS}\left(\mathbf{GMM}(\Delta \tilde{\boldsymbol{\theta}}_{1}^{1}, \dots, \Delta \tilde{\boldsymbol{\theta}}_{n}^{1}; m)\right);$ // Appendix F.2 $\{\pi_1, \ldots, \pi_n, \text{MPO}\} = \mathbf{GMM}(\Delta \tilde{\boldsymbol{\theta}}_1^1, \ldots, \Delta \tilde{\boldsymbol{\theta}}_n^1; M);$ set $E_n(\text{MPO})$; // **1st stage:** GMM // E_c is set based on MPO set $E_c(MPO)$; continue; // go to next round (e=2)else if $e \in \{2, \ldots, E_c\}$ then for each client $i \in \{1, \ldots, n\}$ do $R^{e}(i) \leftarrow m \text{ with probability } \pi_{i}[m]; // 2nd stage: soft clustering$ else on server: broadcast cluster models $\{\boldsymbol{\theta}_m^e\}_{m=1}^M$ to all clients for each client $i \in \{1, \ldots, n\}$ do $R^{e}(i) = \arg\min_{m} f_{i}(\boldsymbol{\theta}_{m}^{e});$ // 3rd stage: "local" clustering for each client $i \in \{1, ..., n\}$ in parallel do // batch size b $b_i^e \leftarrow b$; $\Delta \tilde{\boldsymbol{\theta}}_{i}^{e} \leftarrow \texttt{DPSGD}\left(\boldsymbol{\theta}_{R^{e}(i)}^{e}, b_{i}^{e}, N_{i}, K, z_{i}, c\right)$ on server: for each client $i \in \{1, \dots, n\}$ do $\bigcup w_i^e \leftarrow \frac{N_i}{\sum_{j=1}^n \mathbbm{1}_{R^e(j)=R^e(i)N_j}}$ for $m \in \{1, ..., M\}$ do $\boldsymbol{\theta}_m^{e+1} \leftarrow \boldsymbol{\theta}_m^e + \sum_{i \in \{1, \dots, n\}} \mathbbm{1}_{R^e(i) = m} w_i^e \Delta \tilde{\boldsymbol{\theta}}_i^e \; ; \qquad \textit{//} i \; \text{contributes to} \; R^e(i)$

4.2 REDUCING GMM UNCERTAINTY VIA USING FULL BATCH SIZE IN THE FIRST ROUND

The DP noise in the model updates $\{\Delta \hat{\theta}_i^1\}_{i=1}^n$ makes it harder for the server to cluster clients by running GMM on the model updates. Thus, an efficient clustering algorithm should be robust to this extra DP noise. The following lemma, which is an extension of a similar result in (Malekmohammadi et al., 2024), shows that the noise in model update $\Delta \tilde{\theta}_i^e$ at the of round *e*, including stochastic and DP noise, heavily drops with the batch size b_i^e that client *i* uses during round *e*. This suggests to use large batch sizes in the first round to improve the quality of clustering on the server side.

Lemma 4.1. Let us assume $\theta_i^{e,0}$ is the model parameter that client *i* is assigned at the beginning of round *e*. At the end of round, the client generates the noisy DP model update $\Delta \tilde{\theta}_i^e(b_i^e)$ after K local epochs with step size η_l . The amount of noise in the resulting model update can be found as:

$$\sigma_i^{e^2}(b_i^e) := \operatorname{Var}(\Delta \tilde{\boldsymbol{\theta}}_i^e(b_i^e) | \boldsymbol{\theta}_i^{e,0}) \approx K \cdot N_i \cdot \eta_l^2 \cdot \frac{pc^2 z_i^2(\epsilon, \delta, b_i^1, b_i^{>1}, N_i, K, E)}{b_i^{e^3}}.$$
(2)

We have shown b_i^e as an argument of $\sigma_i^{e^2}(b_i^e)$ to emphasize on its dependence on b_i^e . The lemma means that the noise level in $\Delta \tilde{\theta}_i^e$ decreases fast with b_i^e (Malekmohammadi et al., 2024; Räisä et al.,



Figure 2: PCA visualization of updates $\{\Delta \hat{\theta}_i^1\}_{i=1}^n$ on 2D space. Left: $\epsilon_i = 10$, $b_i^e = 32$ for all *i* and *e*. **Right:** $\epsilon_i = 10$, $b_i^1 = b^1 = N = 6600$, i.e., full batch size (assuming $N_i = N = 6600$ for all clients), and $b_i^{>1} = 32$ for all *i*. The empty markers show the centers of the Gaussian components. The model updates are obtained from running DPFedAvg on CIFAR10 with covariate shift (rotation) between clusters, and under the same values as in Figure 3.



Figure 3: Plot of $\operatorname{Var}(\Delta \tilde{\theta}_i^1(b_i^1)|\theta_i^{init})$ (left) and $\operatorname{Var}(\Delta \tilde{\theta}_i^e(b_i^e)|\theta_i^{e,0})$ (e > 1) (right) v.s. both b_i^1 and $b_i^{>1}$ obtained from Equation (2) and Renyi-DP Accountant (Mironov et al., 2019) in a setting with $N_i = 6600, \epsilon = 5, \delta = 10^{-4}, c = 3, K = 1, E = 200, p = 11, 181, 642, \eta_l = 5 \times 10^{-4}$. There are two clear messages: 1) for all $e \in \{1, \dots, E\}$, $\operatorname{Var}(\Delta \tilde{\theta}_i^e(b_i^e)|\theta_i^{e,0})$ decreases with b_i^e quickly. This was observed in Lemma 4.1. 2) The effect of $b_i^{>1}$ in the left figure is more than the effect of b_i^1 in the right figure. The reason is that $b_i^{>1}$ is used in E - 1 rounds, while b_i^1 is used only in the first round. Also, see Figure 8 in the appendix for the plot of $z_i(\epsilon, \delta, b_i^1, b_i^{>1}, N_i, K, E)$ v.s. b_i^1 and $b_i^{>1}$.

2024). Let us consider e = 1 especially: If a client i can increase its batch size 10 times by using its full batch size in round e = 1, the variance of the noise in its model update $\Delta \theta_i^i(b_i^1)$ drops almost 1000 times. If all clients do so, it becomes much easier for the server to cluster them at the end of the first round, by learning a GMM on $\{\Delta \tilde{\theta}_i^1\}_{i=1}^n$, as their updates become more separable. An illustration of the considerable effect of using full batch sizes in the first round (i.e., $b_i^1 = N_i$) on the noise level in model updates $\{\Delta \tilde{\theta}_i^1\}_{i=1}^n$ is shown in Figure 2. Furthermore, instead of fixing $b_i^{>1}$ to some value, we have also demonstrated the effect of *both* batch sizes b_i^1 and $b_i^{>1}$ on the noise levels $\operatorname{Var}(\Delta \tilde{\theta}_i^1 | \theta_i^{init})$ (e = 1) and $\operatorname{Var}(\Delta \tilde{\theta}_i^e | \theta_i^{e,0})$ (e > 1) separately, in Figure 3. As a take away, Figure 3 left, suggests that in order to make $\{\Delta \tilde{\theta}_i^1\}_{i=1}^n$ less noisy, we have to make $\{b_i^1\}_{i=1}^n$ larger and make $\{b_{i=1}^{j}\}_{i=1}^{n}$ smaller, similar to what done in Figure 2 right. These interesting results are consistent with the observations in (De et al., 2022; Anil et al., 2021; Dörmann et al., 2021; Hoory et al., 2021; Li et al., 2022; Luo et al., 2021) that increasing the batch size can significantly improve the privacy-utility trade-off of DPSGD. In the next section, we will provide a theoretical justification for these observations, especially Figure 2.

324 4.2.1 EFFECT OF BATCH SIZES $\{b_i^1\}_{i=1}^n$ ON MODEL UPDATES $\{\Delta \tilde{\theta}_i^1\}_{i=1}^n$ 325

326 Hereafter, we focus on round e = 1, and show theoretically why increasing batch sizes $\{b_i^1\}_{i=1}^n$ improves the distinguishability of the model updates $\{\Delta \tilde{\theta}_i^1\}_{i=1}^n$. For simplicity, we assume clients 327 have the same dataset sizes and batch sizes: $\forall i : N_i = N, b_i^{-1} = b^1$. Also, remember that $\theta_i^{1,0} = \theta^{init}$. 328 According to Equation (2) and having uniform privacy parameters (ϵ, δ) , we have: $\forall i : \sigma_i^{1^2}(b^1) :=$ $\operatorname{Var}[\Delta \tilde{\theta}_i^1(b^1)|\theta^{init}] = \sigma^{1^2}(b^1)$. Hence, we can consider the model updates $\{\Delta \tilde{\theta}_i^1(b^1)\}_{i=1}^n$ as the samples from a mixture of M Gaussian distributions with mean, covariance matrix, prior probability 330 331 332 parameters: $\psi^*(b^1) = \{\mu_m^*(b^1), \Sigma_m^*(b^1), \alpha_m^*\}_{m=1}^M$, where $\forall m : \alpha_m^* > 0$ and $\mu_m^*(b^1) \neq \mu_{m'}^*(b^1)$ 333 $(m \neq m')$. Also, model update $\Delta \tilde{\theta}_i^1(b^1)$ comes from component m = s(i): 334

335

340 341 342

365

 $\mu_m^*(b^1) := \mathbb{E}\left[\Delta \tilde{\boldsymbol{\theta}}_i^1(b_i^1) \middle| \boldsymbol{\theta}^{\textit{init}}, b_i^1 = b^1, s(i) = m\right],$ $1^{2}(h^{1})$ Г ٦

(3)

$$\Sigma_{m}^{*}(b^{1}) := \mathbb{E}\left[\left(\Delta\tilde{\theta}_{i}^{1}(b_{i}^{1}) - \mu_{m}^{*}(b^{1})\right)\left(\Delta\tilde{\theta}_{i}^{1}(b_{i}^{1}) - \mu_{m}^{*}(b^{1})\right)^{\top} \middle| \theta^{init}, b_{i}^{1} = b^{1}, s(i) = m\right] = \frac{\sigma^{*}(b^{*})}{p} \mathbb{I}_{p}$$
(4)

343 where the last equality is from $\operatorname{Var}[\Delta \tilde{\theta}_i^1 | \theta^{init}, b_i^1 = b^1] = \mathbb{E}[\|\Delta \tilde{\theta}_i^1 - \mu^*_{s(i)}(b^1)\|^2] = \sigma^{1^2}(b^1)$ and 344 that the noises existing in each of the p elements of $\Delta \tilde{\theta}_i^1$ are *i.i.d* (hence, $\Sigma_m^*(b^1)$ is a diagonal 345 covariance matrix with equal diagonal elements). Intuitively, we expect more separation between the 346 true Gaussian components $\{\mathcal{N}(\mu_m^*(b^1), \Sigma_m^*(b^1))\}_{m=1}^M$, from which clients' updates $\{\Delta \tilde{\theta}_i^1\}_{i=1}^n$ are sampled, to make the model updates more distinguishable for server. In the following, we show that 347 the overlap between the Gaussian components $\{\mathcal{N}(\mu_m^*(b^1), \Sigma_m^*(b^1))\}_{m=1}^M$ decreases fast with b^1 . 348

349 **Lemma 4.2.** Let us assume $\Delta_{m,m'}(b^1) := \|\mu_m^*(b^1) - \mu_{m'}^*(b^1)\|$ when $\forall i : b_i^1 = b^1$. Then, the overlap 350 between the pair $\mathcal{N}(\mu_m^*(b^1), \Sigma_m^*(b^1))$ and $\mathcal{N}(\mu_{m'}^*(b^1), \Sigma_{m'}^*(b^1))$ is $O_{m,m'} = 2Q(\frac{\sqrt{p}\Delta_{m,m'}(b^1)}{2\sigma^1(b^1)})$, 351 where $\sigma^{1^2}(b^1) := \operatorname{Var}[\Delta \tilde{\theta}_i^1 | \theta^{init}, b_i^1 = b^1]$ and $Q(\cdot)$ is the tail distribution function of the standard 352 normal distribution. Furthermore, if we increase $b_i^1 = b^1$ to $b_i^1 = kb^1 \leq N$ (for all i), we have 353 $O_{m,m'} \leq 2Q(\frac{\sqrt{kp}\Delta_{m,m'}(b^1)}{2\rho\sigma^1(b^1)})$, where $1 \leq \rho \in \mathcal{O}(1)$ is a small constant. 354 355

356 The lemma states that using a large batch size in the first round results in a *fast* reduction of the 357 overlap between the underlying components, which leads to more distinguishability for $\{\Delta \hat{\theta}_i^1\}_{i=1}^n$ on 358 the server side (see Figure 2, right). One of the beneficial consequences of this well separation is that *RC-DPFL* becomes robust to the initialization of the GMM model. Furthermore, note that for a fixed batch size b^1 , the terms $\Delta_{m,m'}(b^1)$ and $\sigma^1(b^1)$ represent the "data heterogeneity level across clusters" 360 m and m''' and "privacy sensitivity of their clients", respectively. We define the "separation score" 361 $SS(m,m') = \frac{\sqrt{p}\Delta_{m,m'}(b^1)}{2\sigma^1(b^1)} = \frac{\Delta_{m,m'}(b^1)}{2\sigma^1(b^1)/\sqrt{p}}$ between two components m and m' as a measure of their 362 separability. The larger SS(m, m'), the smaller their overlap $O_{m,m'} = 2Q(SS(m, m'))$. Based on 364 the form of the Q function, an SS(m, m') above 3 can be considered as a complete separation.

366 4.2.2 CONFIDENCE OF GMM 367

As we observed in Lemma 4.2, the separation score SS(m,m') (the overlap $O_{m,m'}$) increases 368 (decreases) as b^1 increases. Remember that $SS(m, m') = \frac{\Delta_{m,m'}(b^1)}{2\sigma^1(b^1)/\sqrt{p}}$, and note that $\sigma^{1^2}(b^1)/p$ is the value of diagonal elements of energy in the second 369 370 value of diagonal elements of covariance matrices of Gausssian components (Equation (4)), which 371 the GMM aims to learn. Therefore, when the GMM is learned, we can use its parameters to get an 372 estimate SS(m, m') for every cluster pair m and m'. Then, we can define the **"minimum pairwise**" 373 separation score'' as MSS = $\min_{m,m'} \widehat{SS}(m,m') \in [0,+\infty)$ as a measure of confidence of the 374 learned GMM in its clusterings. The larger the MSS of a learned GMM, the more "confident" it is in its 375 clustering decisions. For instance, if we learn a GMM on Figure 2 left, it will have a much smaller MSS than when we learn a GMM on Figure 2 right. We can similarly define the estimated "maximum 376 pairwise overlap'' for a learned GMM as MPO = $2Q(MSS) \in [0, 1)$, as a measure of uncerntainty 377 of the learned GMM (the smaller the better. Q is a decreasing function).

378 4.3 CONVERGENCE RATE OF EM FOR LEARNING GMM 379 4.3 CONVERGENCE RATE OF EM FOR LEARNING GMM

Let us define the maximum pairwise overlap in $\psi^*(b^1) = \{\mu_m^*(b^1), \Sigma_m^*(b^1), \alpha_m^*\}_{m=1}^M$, as $O^{\max}(\psi^*(b^1)) = \max_{m,m'} O_{m,m'}(\psi^*(b^1))$. According to Lemma 4.2, when b^1 is large enough, $O^{\max}(\psi^*(b^1))$ decreases (like in Figure 2, right) and we can expect EM to converge to the true GMM parameters $\psi^*(b^1)$. Next, we analyze the local convergence rate of EM around the true solution.

Theorem 4.3. (*Ma et al.*, 2000) Given model updates $\{\Delta \tilde{\theta}_i^1(b^1)\}_{i=1}^n$, which are samples from a true mixture of Gaussians $\{\mathcal{N}(\mu_m^*(b^1), \Sigma_m^*(b^1)), \alpha_m^*\}_{m=1}^M$, if $O^{max}(\psi^*(b^1))$ is small enough, then:

386 387

384

399

 $\lim_{r \to \infty} \frac{\|\psi^{r+1} - \psi^*(b^1)\|}{\|\psi^r - \psi^*(b^1)\|} = o\left(\left[O^{\max}(\psi^*(b^1))\right]^{0.5-\gamma}\right),\tag{5}$

as *n* increases. ψ^r is the GMM parameters returned by EM after *r* iterations. γ is an arbitrary small positive number, and o(x) means it is a higher order infinitesimal as $x \to 0$: $\lim_{x\to 0} \frac{o(x)}{x} = 0$.

This means that convergence rate of EM around the true solution $\psi^*(b^1)$ is faster than how $\begin{bmatrix} O^{\max}(\psi^*(b^1)) \end{bmatrix}^{0.5-\gamma}$ decreases with b^1 . In Lemma 4.2, we showed that $O^{\max}(\psi^*(b^1))$ indeed drops fast as b^1 increases. Therefore, if clients have a large enough dataset size and use full batch sizes in the first round, convergence rate of EM approaches approximately 0. *Hence, as an important consequence, the computational complexity of learning the GMM in the first round decreases fast.*

400 4.4 APPLICABILITY OF RC-DPFL

401 402 Even when the number of the underlying clusters (M) is not known beforehand, we can find it with 403 high accuracy based on the confidence metric $MSS \in (0, +\infty)$ defined above (line 9 of Algorithm 1). 404 Intuitively, we choose the M which yields to the largest confidence level MSS for the resulting GMM. 405 We have provided further details about how to find M in these scenarios in Appendix F.2.

The strategy switching time E_c can also be set using the uncertainty metric MPO $\in [0, 1)$. Intuitively, if the learned GMM is not certain about its clustering decisions, RC-DPFL should not rely on its decisions for a large E_c , and vice versa. Hence, we can set E_c as a decreasing function of MPO. For instance, $E_c = (1 - \text{MPO})\frac{E}{2}$ means that if a GMM is completely confident about its clusterings, e.g., what happens in Figure 2 right, the server changes the clustering strategy to loss-based after the first half of the training time. This change happens earlier as the uncertainty increases (e.g., when ϵ is small), and RC-DPFL slowly gets close to the completely loss-based clustering.

Furthermore, we already know that in order to have a quality client clustering at the end of the first round, $\{b_i^{>1}\}_{i=1}^n$ should be small (from Figure 3. Also, see Appendix F.1 for a more detailed discussion). Finally, note that after the training progress made in the first E_c rounds, the loss-based hard clustering is performed "locally" at clients' side (Ghosh et al., 2020) (line 17 in Algorithm 1). Also, at this stage, the sensitivity of the local model selection of a client *i* to adding/removing a data point to its local dataset is effectively zero. Therefore, there is no privacy concern regarding the local loss-based clusterings performed in the last stage of RC-DPFL (see Appendix H for a formal privacy proof). These important features altogether make RC-DPFL a robust and applicable algorithm.

420 421

422

- 5 EVALUATION
- 423 424 5.1 EXPERIMENTAL SETUP

Datasets, models and baseline algorithms: We evaluate our proposed method on three benchamark datasets, including: MNIST (Deng, 2012), FMNIST (Xiao et al., 2017) and CIFAR10 (Krizhevsky, 2009), with heterogeneous data distributions from covariate shift (rotation) (Kairouz et al., 2021; Werner et al., 2023) and concept shift (label flip) (Werner et al., 2023), which are the commonly used data splits in the literature (see Appendix B). We consider four clusters of clients indexed by $m \in \{0, 1, 2, 3\}$ with $\{3, 6, 6, 6\}$ clients, where the smallest cluster is considered as the minority group. To the best of our knowledge, there was no prior work on DP clustered FL, so we compared to the DP version of existing algorithms. More specifically, we compare with the following baseline algorithms, which are combined with DPSGD: 1. DPFedAvg (Noble et al., 2021): clients run DPSGD
locally and send their model updates to the server 2. KM-CDPFL (Werner et al., 2023): FL with
myopic clustering, in which server clusters clients at the end of each round using running K-means
clustering on their DP model updates 3. f-CDPFL (Ghosh et al., 2020; Mansour et al., 2020): FL with
loss clustering, which clusters clients based on their train loss values on existing cluster models 4.
Oracle-CDPFL: an oracle algorithm which has the knowledge of true clusters from the first round.

Evaluation metrics and baselines: Given the set of n clients, fairness in a DPFL system can 439 440 be measured in terms of the disparate impact of DP on utility (performance drop) of different groups (Chu et al., 2023; Bagdasaryan & Shmatikov, 2019; Tran et al., 2021; Esipova et al., 2022): 441 $\mathcal{F}_{acc} = \max_{i,j \in [n]} |\Delta acc_i(\theta_i) - \Delta acc_j(\theta_j)|$, where θ_i is the model assigned to agent i at the end 442 of DP training, and $\Delta acc_i(\theta_i) = \max_{\theta^*} acc_i(\theta^*) - acc_i(\theta_i)$, where θ^* is any possible model. 443 Similarly, we can measure fairness in terms of the increment enforced to clients' train loss (Tran et al., 444 2021; Esipova et al., 2022; Chu et al., 2023): $\mathcal{F}_{loss} = \max_{i,j \in [n]} |\xi_i(\theta_i) - \xi_j(\theta_j)|$, where $\xi_i(\theta_i) = \xi_i(\theta_i)$ 445 $f_i(\theta_i) - \min_{\theta^*} f_i(\theta^*)$. These notions of fairness compare the cost of adding differential privacy on 446 different clients, and define client-level fairness as the equality of "performance drop" across clients. 447 Following (Chu et al., 2023), we estimate the model θ^* for each cluster by centrally training a model 448 with SGD based on the data of the clients belonging to that cluster. We also consider the following 449 evaluation metrics: average test accuracy (overall, majority, minority), worst accuracy across clients 450 (Mohri et al., 2019), and maximum accuracy disparity across clients: $\max_{i,j} |acc_i(\theta_i) - acc_j(\theta_j)|$. 451

5.2 Results

438

452

453

456

457

458

459

460

461

479

480

In our experiments, we aim to 1) compare RC-DPFL with other clustering approaches, 2) analyze its robustness to noise; and 3) evaluate its robustness to different types of data heterogeneity.

RQ1: How does RC-DPFL perform compared to other algorithms? We first explore how RC-DPFL performs in comparison with the defined baseline algorithms. Figure 4 shows the performance for MNIST and FMNIST in terms of per cluster performance and fairness metric \mathcal{F}_{acc} . Through these results, it is clear that RC-DPFL performance on-par with the oracle-DPFL baseline, which constitutes the ideal case. This is mainly attributed to the accurate clustering obtained in RC-DPFL as seen in subfigure (c), which compares the success rate of clustering using loss function (f-CDPFL) versus our approach. Also, KM-CDPFL incurs the highest unfairness \mathcal{F}_{acc} , due to clustering errors.



Figure 4: Comparison of RC-DPFL with the defined baselines on MNIST (Top row) and FMNIST (bottom row) with C1 being the minority cluster. All results are for $\epsilon = 5$.

RQ2: How does RC-DPFL perform under different levels of noise? Figure 5 on CIFAR10, shows the effect of varying levels of DP noise on the fairness of different algorithms ($\delta = 10^{-4}$) in terms of three different metrics. RC-DPFL performs close to the oracle algorithm in terms of all the three metrics, which shows its robustness to the DP noise in the system. RC-DPFL has the smallest gap between majority and minority groups in terms of the three disparity metrics and outperforms the two other baseline algorithms for improving the minority cluster. The gap is even larger on smaller values of ϵ , for instance, the minimum accuracy for $\epsilon = 2$ using RCDPFL is 5% and 9% higher than the best performing benchmark algorithm on FMNIST and MNIST respectively, while unfairness (\mathcal{F}_{acc}) is 6% and 13% lower. Detailed results for other datasets can be found in Tables 4–9 in the appendix, which include results for accuracy across groups and various fairness metrics.





RQ3: How does RC-DPFL perform under different types of data heterogeneity across clients? We evaluate how different types of distribution shift across client groups affect the clustered DPFL algorithms. To do so, we compare covariate shift and concept shift on CIFAR10 dataset. Concept shift, which can also be viewed as a label flipping attack, has a more significant impact on performance in the single model case, as labels vary across client groups. Figure 6 shows results with $\epsilon = 5.0$ in terms of per-cluster performance and fairness, as well as clustering accuracy for different values of ϵ . Through these results, we notice that it is easier to detect minorities with the loss values in the case of concept-shift. Nonetheless, we also note that mistakes become more costly in this case. Table 11, Table 10 in the appendix show a high variance for f-CPFL across experiments, especially smaller values of ϵ , while RC-DPFL is more consistent. Additionally, in terms of fairness metrics, and across different ϵ values, RC-DPFL still outperforms the baselines, and remains closer to the oracle case.



Figure 6: Comparison of RC-DPFL with the defined baselines on CIFAR10 with covariate shift (Top row) and concept-shift (bottom row) with C1 being the minority cluster. All results are for $\epsilon = 5$.

6 CONCLUSION

We proposed the first DP clustered FL algorithm, which addresses high data heterogeneity in privacysensitive FL environments. By clustering clients based on their model updates and training loss values,
and mitigating noise impacts with larger batch sizes, our approach enhances utility and fairness with
minimal computational overhead, while maintaining DP. Moreover, the robustness to noise, and the
ability to handle various types of distribution shifts shows the applicability of our approach.

540 REFERENCES

542 Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In Proceedings of the 2016 ACM SIGSAC 543 Conference on Computer and Communications Security, 2016. URL https://doi.org/10. 544 1145/2976749.2978318. 546 Rohan Anil, Badih Ghazi, Vineet Gupta, Ravi Kumar, and Pasin Manurangsi. Large-scale differen-547 tially private bert, 2021. URL https://arxiv.org/abs/2108.01624. 548 549 Eugene Bagdasaryan and Vitaly Shmatikov. Differential privacy has disparate impact on model accuracy. In Neural Information Processing Systems, 2019. 550 551 Dimitris Bertsimas, Vivek F Farias, and Nikolaos Trichakis. The price of fairness. Operations 552 research, 59(1):17-31, 2011. 553 554 Christopher Briggs, Zhong Fan, and Peter Andras. Federated learning with hierarchical clustering of 555 local updates to improve training on non-iid data, 2020. 556 Wenda Chu, Chulin Xie, Boxin Wang, Linyi Li, Lang Yin, Arash Nourian, Han Zhao, and Bo Li. Focus: Fairness via agent-awareness for federated learning on heterogeneous data, 2023. 558 559 Soham De, Leonard Berrada, Jamie Hayes, Samuel L. Smith, and Borja Balle. Unlocking high-560 accuracy differentially private image classification through scale, 2022. URL https://arxiv. 561 org/abs/2204.13650. 562 563 Li Deng. The mnist database of handwritten digit images for machine learning research [best of the 564 web]. IEEE Signal Processing Magazine, 2012. URL https://ieeexplore.ieee.org/ 565 document/6296535. 566 John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. Local privacy and statistical minimax 567 rates. 2013 51st Annual Allerton Conference on Communication, Control, and Computing (Aller-568 ton), pp. 1592-1592, 2013. URL https://api.semanticscholar.org/CorpusID: 569 1597053. 570 571 John C. Duchi, Martin J. Wainwright, and Michael I. Jordan. Minimax optimal procedures for locally 572 private estimation. Journal of the American Statistical Association, 113:182 – 201, 2018. URL 573 https://api.semanticscholar.org/CorpusID:15762329. 574 Cynthia Dwork. A firm foundation for private data analysis. Commun. ACM, 2011. URL https: 575 //doi.org/10.1145/1866739.1866758. 576 577 Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. Found. Trends 578 Theor. Comput. Sci., 2014. URL https://dl.acm.org/doi/10.1561/0400000042. 579 580 Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. Our 581 data, ourselves: Privacy via distributed noise generation. In Proceedings of the 24th Annual 582 International Conference on The Theory and Applications of Cryptographic Techniques, 2006a. 583 URL https://doi.org/10.1007/11761679_29. 584 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity 585 in private data analysis. In Proceedings of the Third Conference on Theory of Cryptography. 586 Springer-Verlag, 2006b. URL https://doi.org/10.1007/11681878_14. 588 Friedrich Dörmann, Osvald Frisk, Lars Nørvang Andersen, and Christian Fischer Pedersen. Not all 589 noise is accounted equally: How differentially private learning benefits from large sampling rates, 590 2021. URL https://arxiv.org/abs/2110.06255. Maria S. Esipova, Atiyeh Ashari Ghomi, Yaqiao Luo, and Jesse C. Cresswell. Disparate impact in 592 differential privacy from gradient misalignment. ArXiv, abs/2206.07737, 2022. URL https: //api.semanticscholar.org/CorpusID:249712405.

594 595 596	Tom Farrand, FatemehSadat Mireshghallah, Sahib Singh, and Andrew Trask. Neither private nor fair: Impact of data imbalance on utility and fairness in differential privacy. <i>Proceedings of the 2020 Workshop on Privacy-Preserving Machine Learning in Practice</i> , 2020. URL https:
597	<pre>//api.semanticscholar.org/CorpusID:221655207.</pre>
598	Ferdinando Fioretto, Cuong Tran, Pascal Van Hentenryck, and Kevu Zhu. Differential privacy and
599 600	fairness in decisions and learning tasks: A survey. In <i>Proceedings of the Thirty-First International</i>
601	Joint Conference on Artificial Intelligence. International Joint Conferences on Artificial Intelligence
602 603	Organization, jul 2022. doi: 10.24963/ijcal.2022//66. URL https://doi.org/10.24963% 2Fijcai.2022%2F766.
604	Jonas Geining Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. Inverting gradients
605 606	- how easy is it to break privacy in federated learning? <i>ArXiv</i> , 2020. URL https://api.
607	Semanereseneral.org/corpusit.211720317.
608 609	Robin C. Geyer, Tassilo Klein, and Moin Nabi. Differentially private federated learning: A client level perspective. <i>ArXiv</i> , 2017. URL https://arxiv.org/pdf/1712.07557.pdf.
610	
611	Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. An efficient framework for clustered federated learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin
612	(eds.), Advances in Neural Information Processing Systems, volume 33, pp. 19586–19597. Cur-
614	<pre>ran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/ paper/2020/file/e32cc80bf07915058ce90722ee17bb71-Paper.pdf.</pre>
615	Kaiming He X Zhang Shaoging Ren and Jian Sun. Deep residual learning for image recog-
616	nition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR).
617	2015. URL https://www.cv-foundation.org/openaccess/content cvpr
618	2016/papers/He_Deep_Residual_Learning_CVPR_2016_paper.pdf.
620	Briland Hitaj, Giuseppe Ateniese, and Fernando Pérez-Cruz. Deep models under the gan: Information
621	leakage from collaborative deep learning. Proceedings of the 2017 ACM SIGSAC Conference on
622 623	Computer and Communications Security, 2017. URL https://api.semanticscholar.org/CorpusID:5051282.
624	Olderer Harres Ander Dates A date Tradition On Contract Aller Orderer Res Laboration Nuller A
625 626 627	Uri Stemmer, Ayelet Benjamini, Avinatan Hassidim, and Yossi Matias. Learning and evaluating a differentially private pre-trained language model. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pp. 1178–1189, 2021. URL https://aclanthology.org/
620	2021.findings-emnlp.102/.
629 630	Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Phagaii Kallista Banawitz, Zachary Charles, Graham Cormada, Bachal Cummings, et al. Ad
632 633	vances and open problems in federated learning. <i>Foundations and trends</i> ® <i>in machine learning</i> , 14(1–2):1–210, 2021.
634	Alex Krizhevsky Learning multiple layers of features from tiny images 2009 LIRL https:
635 636	//www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf.
637	Tian Lan, David Kao, Mung Chiang, and Ashutosh Sabharwal. An axiomatic theory of fairness in
638 639	network resource allocation. In 2010 Proceedings IEEE INFOCOM, pp. 1–9, March 2010. doi: 10.1109/INFCOM.2010.5461911. ISSN: 0743-166X.
640	
641 642	ArXiv, abs/1910.03581, 2019. URL https://api.semanticscholar.org/CorpusID:
643	203951869.
644	Tian Li, Ahmad Beirami, Maziar Sanjabi, and Virginia Smith. Tilted empirical risk minimization. In
645	International Conference on Learning Representations, 2020a.
646 647	Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. Fair resource allocation in federated learning. In <i>International Conference on Learning Representations</i> , 2020b.

648 649 650 651 652	Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning through personalization. In Marina Meila and Tong Zhang (eds.), <i>Proceedings of the 38th International Conference on Machine Learning</i> , volume 139 of <i>Proceedings of Machine Learning Research</i> , pp. 6357–6368. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/li21h.html.
653 654 655	Xuechen Li, Florian Tramèr, Percy Liang, and Tatsunori Hashimoto. Large language models can be strong differentially private learners, 2022. URL https://arxiv.org/abs/2110.05679.
656 657 658	Yang Liu, Yan Kang, Chaoping Xing, Tianjian Chen, and Qiang Yang. A secure federated trans- fer learning framework. <i>IEEE Intelligent Systems</i> , 35:70–82, 2020. URL https://api. semanticscholar.org/CorpusID:219013245.
659 660 661	Ziyu Liu, Shengyuan Hu, Zhiwei Steven Wu, and Virginia Smith. On privacy and personalization in cross-silo federated learning, 2022. URL https://arxiv.org/abs/2206.07902.
662 663 664 665	Guixun Luo, Naiyue Chen, Jiahuan He, Bingwei Jin, Zhiyuan Zhang, and Yidong Li. Privacy- preserving clustering federated learning for non-iid data. <i>Future Generation Computer Systems</i> , 154:384–395, 2024. URL https://www.sciencedirect.com/science/article/ pii/S0167739X24000050.
666 667 668	Zelun Luo, Daniel J. Wu, Ehsan Adeli, and Li Fei-Fei. Scalable differential privacy with sparse network finetuning. In 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 5057–5066, 2021.
669 670 671 672	Jinwen Ma, Lei Xu, and Michael I. Jordan. Asymptotic convergence rate of the em algorithm for gaussian mixtures. <i>Neural Computation</i> , 2000. URL https://api.semanticscholar.org/CorpusID:10273602.
673 674	Saber Malekmohammadi, Yaoliang Yu, and Yang Cao. Noise-aware algorithm for heterogeneous differentially private federated learning, 2024. URL https://arxiv.org/abs/2406.03519.
675 676 677	Yishay Mansour, Mehryar Mohri, Jae Ro, and Ananda Theertha Suresh. Three approaches for personalization with applications to federated learning, 2020.
678 679 680	Othmane Marfoq, Giovanni Neglia, Aurélien Bellet, Laetitia Kameni, and Richard Vidal. Federated multi-task learning under a mixture of distributions. In <i>Neural Information Processing Systems</i> , 2021. URL https://api.semanticscholar.org/CorpusID:236470180.
681 682 683	Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In <i>Artificial intelligence and statistics</i> , pp. 1273–1282. PMLR, 2017.
684 685 686 687	H. Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially private recurrent language models. In <i>ICLR</i> , 2018. URL https://arxiv.org/pdf/1710.06963.pdf.
688 689 690	Umberto Michieli and Mete Ozay. Are all users treated fairly in federated learning systems? In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 2318–2322, 2021.
691 692 693	Ilya Mironov. Rényi differential privacy. In 2017 IEEE 30th Computer Security Foundations Symposium (CSF), pp. 263–275. IEEE, August 2017. doi: 10.1109/csf.2017.11. URL http://dx.doi.org/10.1109/CSF.2017.11.
695 696	Ilya Mironov, Kunal Talwar, and Li Zhang. Rényi differential privacy of the sampled gaussian mechanism, 2019. URL https://arxiv.org/abs/1908.10530.
697 698	Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In <i>International Conference on Machine Learning</i> , pp. 4615–4625. PMLR, 2019.
700 701	Maxence Noble, Aurélien Bellet, and Aymeric Dieuleveut. Differentially private federated learning on heterogeneous data. In <i>International Conference on Artificial Intelligence and Statistics</i> , 2021. URL https://proceedings.mlr.press/v151/noble22a/noble22a.pdf.

702 Sikha Pentyala, Nicola Neophytou, Anderson Nascimento, Martine De Cock, and Golnoosh Farnadi. 703 Privfairfl: Privacy-preserving group fairness in federated learning, 2022. 704 705 Maria Rigaki and Sebastián García. A survey of privacy attacks in machine learning. ArXiv, 2020. URL https://api.semanticscholar.org/CorpusID:220525609. 706 707 Yichen Ruan and Carlee Joe-Wong. Fedsoft: Soft clustered federated learning with proximal local 708 updating. CoRR, abs/2112.06053, 2021. URL https://arxiv.org/abs/2112.06053. 709 Ossi Räisä, Joonas Jälkö, and Antti Honkela. Subsampling is not magic: Why large batch sizes work 710 for differentially private stochastic optimisation, 2024. URL https://arxiv.org/abs/ 711 2402.03990. 712 713 Felix Sattler, Klaus-Robert Müller, and Wojciech Samek. Clustered federated learning: Model-714 agnostic distributed multitask optimization under privacy constraints. IEEE Transactions 715 on Neural Networks and Learning Systems, 32:3710-3722, 2019. URL https://api. 716 semanticscholar.org/CorpusID:203736521. 717 Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet Talwalkar. Federated multi-718 task learning. In Neural Information Processing Systems, 2017. URL https://api. 719 semanticscholar.org/CorpusID:3586416. 720 721 Cuong Tran, Ferdinando Fioretto, and Pascal Van Hentenryck. Differentially private and fair deep learning: A lagrangian dual approach. ArXiv, abs/2009.12562, 2020. URL https://api. 722 semanticscholar.org/CorpusID:221970859. 723 724 Cuong Tran, My-Hoa Nathalie Dinh, and Ferdinando Fioretto. Differentially private empirical risk 725 minimization under the fairness lens. In Neural Information Processing Systems, 2021. URL 726 https://api.semanticscholar.org/CorpusID:248498291. 727 Zhibo Wang, Mengkai Song, Zhifei Zhang, Yang Song, Qian Wang, and Hairong Qi. Beyond inferring 728 class representatives: User-level privacy leakage from federated learning. IEEE INFOCOM, 2019. 729 URL https://api.semanticscholar.org/CorpusID:54436587. 730 731 Mariel Werner, Lie He, Sai Praneeth Karimireddy, Michael Jordan, and Martin Jaggi. Provably 732 personalized and robust federated learning, 2023. 733 Yue Wu, Shuaicheng Zhang, Wenchao Yu, Yanchi Liu, Quanquan Gu, Dawei Zhou, Haifeng Chen, and 734 Wei Cheng. Personalized federated learning under mixture of distributions. ArXiv, abs/2305.01068, 735 2023. URL https://api.semanticscholar.org/CorpusID:258436670. 736 737 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking 738 machine learning algorithms. CoRR, 2017. URL http://arxiv.org/abs/1708.07747. 739 Depeng Xu, Wei Du, and Xintao Wu. Removing disparate impact on model accuracy in differen-740 tially private stochastic gradient descent. Proceedings of the 27th ACM SIGKDD Conference 741 on Knowledge Discovery & Data Mining, 2021. URL https://api.semanticscholar. 742 org/CorpusID:236980106. 743 744 Guojun Zhang, Saber Malekmohammadi, Xi Chen, and Yaoliang Yu. Proportional fairness in federated learning. Transactions on Machine Learning Research, 2023. URL https: 745 //openreview.net/forum?id=ryUHgEdWCQ. 746 747 Yang Zhao, Jun Zhao, Mengmeng Yang, Teng Wang, Ning Wang, Lingjuan Lyu, Dusit Tao Niyato, 748 and Kwok-Yan Lam. Local differential privacy-based federated learning for internet of things. IEEE 749 Internet of Things Journal, 8:8836-8853, 2020. URL https://api.semanticscholar. 750 org/CorpusID:215828540. 751 Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. In Neural Information 752 Processing Systems, 2019. URL https://api.semanticscholar.org/CorpusID: 753 195316471. 754

755