

Appendix for A Provably Robust Algorithm for Differentially Private Clustered Federated Learning

A NOTATIONS

Table 1: Used notations

n	number of clients, which are indexed by i
x_{ij}, y_{ij}	j -th data point of client i and its label
\mathcal{D}_i, N_i	local train set of client i and its size
$\mathcal{D}_{i,aug}$	augmented local train set of client i
$\mathcal{B}_i^{e,t}$	the train data batch used by client i in round e and at the t -th gradient update
b_i^e	batch size of client i in round e : $ \mathcal{B}_i^{e,t} = b_i^e$
b_i^1	batch size of client i in the first round $e = 1$
$b_i^{>1}$	set of batch sizes of client i in the rounds $e > 1$
ϵ, δ	desired DP privacy parameters
E	total number of global communication rounds in the DPFL system, indexed by e
θ_m^e	model parameter for cluster m , at the beginning of global round e
K	number of local train epochs performed by clients during each global round e
h	predictor function, e.g., CNN model, with parameter θ
ℓ	cross entropy loss
$s(i)$	the true cluster of client i
$R^e(i)$	the cluster assigned to client i in round e
$\theta_i^{e,0}$	the model parameter passed to client i at the beginning of round e to start its local training
$\Delta \tilde{\theta}_i^e$	the noisy model update of client i at the end of round e , starting from $\theta_i^{e,0}$
$\sigma_i^{e,2}$	conditional variance of the noisy model update $\Delta \tilde{\theta}_i^e$ of client i : $\text{Var}(\Delta \tilde{\theta}_i^e \theta_i^{e,0})$
$\mu_m^*(b^1)$	the center of the m -th cluster (when all clients use batch size b^1 in the first round)
$\Sigma_m^*(b^1)$	the covariance matrix of the m -th cluster (when all clients use batch size b^1 in the first round)
α_m^*	the prior probability of the m -th cluster

B EXPERIMENTAL SETUP

B.1 DATASETS

Data split: We use three datasets MNIST, FMNIST and CIFAR10, and consider a distributed setting with 21 clients. In order to create majority and minority clusters, we consider 4 clusters with different number of clients $\{3, 6, 6, 6\}$ (21 clients in total). The first cluster with the minimum number of clients is the “minority” cluster, and the last three are the “majority” ones. The data distribution $P(x, y)$ varies across clusters. We use two methods for making such data heterogeneity: 1. **covariate shift** 2. **concept shift**. In covariate shift we assume that features marginal distribution $P(x)$ differs from one cluster to another cluster. In order to create this variation, we first allocate samples to all clients in an *uniform* way. Then we rotate the data points (images) belonging to the clients in cluster k by $k * 90$ degrees. For concept shift, we assume that conditional distribution $P(y|x)$ differs from one cluster to another cluster, and we first allocate data samples to clients in a uniform way, and flip the labels of the points allocated to clients: we flip y_{ij} (label of the j -th data point of client i , which belongs to cluster k) to $(y_{ij} + k) \bmod 10$. The local datasets are balanced—all users have the same amount of training samples. The local data is split into train and test sets with ratios 80%, and 20%, respectively. In the reported experimental results, all users participate in each communication round.

Table 2: CNN model for classification on MNIST/FMNIST datasets

Layer	Output Shape	# of Trainable Parameters	Activation	Hyper-parameters
Input	(1, 28, 28)	0		
Conv2d	(16, 28, 28)	416	ReLU	kernel size =5; strides=(1, 1)
MaxPool2d	(16, 14, 14)	0		pool size=(2, 2)
Conv2d	(32, 14, 14)	12,832	ReLU	kernel size =5; strides=(1, 1)
MaxPool2d	(32, 7, 7)	0		pool size=(2, 2)
Flatten	1568	0		
Dense	10	15,690	ReLU	
Total		28,938		

B.2 MODELS AND OPTIMIZATION

We use a simple 2-layer CNN model with ReLU activation, the detail of which can be found in Table 2 for MNIST and FMNIST. Also, we use the residual neural network (ResNet-18) defined in He et al. (2015), which is a large model. To update the local models allocated to each client during each round, we apply DP-SGD (Abadi et al., 2016) with a noise scale z which depends on some parameters, as in Equation (2)

Table 3: Details of the experiments and the used datasets in the main body of the paper. ResNet-18 is the residual neural networks defined in He et al. (2015). CNN: Convolutional Neural Network defined in Table 2.

Datasets	Train set size	Test set size	Data Partition method	# of clients	Model	# of parameters
MNIST	48000	12000	cov. shift	{3, 6, 6, 6}	CNN	28,938
FMNIST	50000	10000	cov. shift	{3, 6, 6, 6}	CNN	28,938
CIFAR10	50000	10000	cov. and con. shift	{3, 6, 6, 6}	ResNet-18	11,181,642

B.3 DP PARAMETERS

For each dataset, 5 different values of ϵ from set $\{2, 3, 4, 5, 10\}$ are used. We fix δ for all experiments to 10^{-4} . We also set the clipping threshold c equal to 3, as it results in better test accuracy for the considered datasets, as reported in (Abadi et al., 2016). We use the Renyi DP (RDP) privacy accountant (TensorFlow privacy implementation) during the training time. This accountant is able to handle the difference in the batch size between the first round $e = 1$ and the next rounds $e > 1$ by accounting the composition of the corresponding *heterogeneous* private mechanisms.

B.4 ALGORITHMS TO COMPARE AND TUNING HYPERPARAMETERS

We compare our RC-DPFL algorithm, which benefits from robust clustering, with four baseline algorithms, including: 1) “DPFedAvg” (Noble et al., 2021), which learns one global model for all clients 2) An extension of the IFCA algorithm (Ghosh et al., 2020) to DPFL systems, which we call “f-CDPFL” 3) An extension of the gradient based clustering algorithm (algorithm 1 in (Werner et al., 2023)) to DPFL systems, which we call “KM-CDPFL” 4) An oracle algorithm, which has the knowledge of the true underlying clients’ clusters, which we call “O-CDPFL”. For each algorithm and each dataset, we find the best learning rate from a grid: the one which is small enough to avoid divergence of the DP federated optimization, and results in the lowest average loss (across clients) at the end of FL training on a “validation set” with size 10,000 samples. Here are the grids we use for each dataset:

- MNIST: {1e-3, 2e-3, 5e-3, 1e-2};
- FMNIST: {1e-3, 2e-3, 5e-3, 1e-2};
- CIFAR10: {1e-3, 2e-3, 5e-3, 1e-2}.

B.5 GAUSSIAN MIXTURE MODEL

We use the Gaussian Mixture Model of Scikitlearn, which can be found here: <https://scikit-learn.org/dev/modules/generated/sklearn.mixture.GaussianMixture.html>. The GMM model has three hyper-parameters:

- 1) parameter initialization, which we set to “k-means++”. This is because this type of initialization leads to both low time to initialize and low number of EM iterations for the GMM to converge
- 2) Type of the covariance matrix, which we set to “spherical”, i.e., each component has a diagonal covariance matrix with a single value as its diagonal elements. This is in accordance with Equation (18) and that we know *the covariance matrices should be diagonal*.
- 3) Finally, the number of components (clusters) is either known or it is unknown. In the latter case, we have explained in Appendix F.2 how we can find the true number of clusters by using the confidence level (MSS) of the GMM model.

B.6 HYPERPARAMETERS OF THE RC-DPFL

As explained in the paper, RC-DPFL has four hyperparameters, which we know how to set, as explained in Section 4.4:

- 1) b_i^1 is set locally by clients to N_i (their dataset size), i.e., full batch size for easier clustering at the end of the first round
- 2) $b_i^{>1}$, which is the batch size used during the rounds $e > 1$, and has to be set to a small value, as observed in Figure 3 right, and as explained in Appendix F.1. As observed in Figure 10, RC-DPFL is not sensitive to this parameter, as long as a small value is chosen for it. For the results in the paper, we have set $b_i^{>1} = 32$
- 3) E_c which is the time of changing the clustering strategy, and as explained in Section 4.4, it has to be related to the uncertainty level (MPO) of the learned GMM such that E_c decreases as MPO of the learned GMM increases. In our experiments, we have used $E_c = (1 - \text{MPO}) \frac{E}{2}$. In this way, if the learned GMM is uncertain about its clusterings, RC-DPFL does trust it for many rounds and instead, uses loss-based hard clustering from the early rounds by using a smaller E_c . Therefore, RC-DPFL gradually reduces to the loss-based clustering method in the baseline “f-CDPFL” as the MPO of the learned GMM increases.
- 4) Finally, if the true number of clusters M is not known, we can find it by using the very intuitive method explained in Appendix F.2.

Therefore, all the hyperparameters of RC-DPFL can be set efficiently and easily.

C EXAMPLES

In this section, we explain an example to elaborate that why clustering clients based on their losses (model updates) is prone to errors in the first (last) rounds. For example, consider Figure 7, where there are $M = 2$ clusters (red and blue) and $n = 4$ clients. The clients in the red cluster have loss functions $f_1(\theta) = 4(\theta + 6)^2$ and $f_2(\theta) = 4(\theta + 5)^2$ with optimum cluster parameter $\theta_1^\infty = -5.5$. Also, the clients in the blue cluster have loss functions $f_3(\theta) = 4(\theta - 5)^2$ and $f_4(\theta) = 4(\theta - 6)^2$ with optimum cluster parameter $\theta_2^\infty = 5.5$. Clustering algorithms, which cluster clients based on their loss values on clusters’ models, are vulnerable to model initialization. For example, in Figure 7, if we initialize the clusters’ parameters with $\theta_1^0 = -11$ and $\theta_2^0 = 0$ (shown in the figure), all four clients will initially select cluster 2, since they have smaller losses on its parameter. At $\theta_2^0 = 0$, the average of clients’ gradients (model updates) is zero, so all clients will remain stuck at θ_2^0 and will always select cluster 2.

On the other hand, clustering clients based on their model updates (gradients) (Werner et al., 2023; Briggs et al., 2020; Sattler et al., 2019) have clearly issues. One of these issues appears after some rounds of training. For instance, even if we assume these algorithms can initially cluster clients “perfectly” in each round e , the clients’ model updates (gradients) will approach zero as the clusters’ models converge to their optimum parameters. Hence, clients from different clusters may appear to

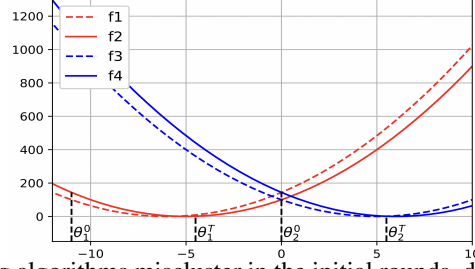


Figure 7: Loss-based clustering algorithms miscluster in the initial rounds, due to model initialization. Also, even with the assumption of perfect clustering of clients in the first rounds, clustering algorithms based on gradients (model updates) lead to clustering errors in the last rounds, due to the gradients approaching to zero.

belong to the same cluster, which results in clustering mistakes. For example, as shown in Figure 1, right, let us assume after T rounds of “correct” clustering of clients, the clusters’ parameters get to $\theta_1^T = -4.5$ and $\theta_2^T = 5.5$ (shown in the figure). At this parameters, clients 1 and 2 (which have been “correctly” assigned to cluster 1 so far) will have gradients $f'_1(\theta_1^T) = 12$ and $f'_2(\theta_1^T) = 4$. Similarly, clients 3 and 4 (which have been “correctly” assigned to cluster 2 so far) will have $f'_3(\theta_2^T) = 4$ and $f'_4(\theta_2^T) = -4$. We see that f'_2 is closer to f'_3 and f'_4 than to f'_1 , and in the next round it will wrongly be assigned to wrong cluster 2. This happens while the clients are clearly distinguishable based on their losses, as some progress in training has been made after T rounds: $f_1(\theta_1^T) = 9$, while $f_1(\theta_2^T) = 23^2$, which clearly means that client 1 correctly belongs to cluster 1. Therefore, after making some progress in training the clusters’ models, it makes more sense to use a loss-based clustering strategy than using a strategy based on clients’ gradients (model updates).

D PROOFS

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:*

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

Proof. The proof is restated from (Malekmohammadi et al., 2024). We consider two illustrative scenarios:

Scenario 1: the clipping threshold c is effective for all samples in a batch: in this case we have: $\forall j \in \mathcal{B}_i^{e,t} : c < \|g_{ij}(\theta)\|$. Also, we know that the two sources of randomness (i.e. stochastic and Gaussian noise) are independent, thus their variances can be summed up. Let us assume that $E[\tilde{g}_{ij}(\theta)] = G_i(\theta)$ for all samples j . From Equation (1), we can find the mean of each *batch gradient* $\tilde{g}_i^{e,t}(\theta)$ (of client i in round e and gradient step t) as follows:

$$\mathbb{E}[\tilde{g}_i^{e,t}(\theta)] = \frac{1}{b_i^e} \sum_{j \in \mathcal{B}_i^{e,t}} \mathbb{E}[\tilde{g}_{ij}(\theta)] = \frac{1}{b_i^e} \sum_{j \in \mathcal{B}_i^{e,t}} G_i(\theta) = G_i(\theta). \quad (6)$$

Also, from Equation (1), we can find the variance of each *batch gradient* $\tilde{g}_i^{e,t}(\theta)$ (of client i in round e and gradient step t) as follows:

$$\begin{aligned}
\sigma_{i,\tilde{g}}^2(b_i^e) &:= \text{Var}[\tilde{g}_i^{e,t}(\theta)] = \text{Var}\left[\frac{1}{b_i^e} \sum_{j \in \mathcal{B}_i^{e,t}} \bar{g}_{ij}(\theta)\right] + \frac{p\sigma_{i,\text{DP}}^2}{b_i^{e^2}} \\
&= \frac{1}{b_i^{e^2}} \left(\mathbb{E}\left[\left\|\sum_{j \in \mathcal{B}_i^{e,t}} \bar{g}_{ij}(\theta)\right\|^2\right] - \left\|\mathbb{E}\left[\sum_{j \in \mathcal{B}_i^{e,t}} \bar{g}_{ij}(\theta)\right]\right\|^2 \right) + \frac{pc^2 z_i^2(\epsilon_i, \delta_i, b_i^1, b_i^{>1}, N_i, K, E)}{b_i^{e^2}} \\
&= \frac{1}{b_i^{e^2}} \left(\mathbb{E}\left[\left\|\sum_{j \in \mathcal{B}_i^{e,t}} \bar{g}_{ij}(\theta)\right\|^2\right] - \left\|\sum_{j \in \mathcal{B}_i^{e,t}} G_i(\theta)\right\|^2 \right) + \frac{pc^2 z_i^2(\epsilon_i, \delta_i, b_i^1, b_i^{>1}, N_i, K, E)}{b_i^{e^2}} \\
&= \frac{1}{b_i^{e^2}} \left(\underbrace{\mathbb{E}\left[\left\|\sum_{j \in \mathcal{B}_i^{e,t}} \bar{g}_{ij}(\theta)\right\|^2\right]}_{\mathcal{A}} - b_i^{e^2} \|G_i(\theta)\|^2 \right) + \frac{pc^2 z_i^2(\epsilon_i, \delta_i, b_i^1, b_i^{>1}, N_i, K, E)}{b_i^{e^2}}, \quad (7)
\end{aligned}$$

where:

$$\begin{aligned}
\mathcal{A} &= \mathbb{E}\left[\left\|\sum_{j \in \mathcal{B}_i^{e,t}} \bar{g}_{ij}(\theta)\right\|^2\right] = \sum_{j \in \mathcal{B}_i^{e,t}} \mathbb{E}\left[\|\bar{g}_{ij}(\theta)\|^2\right] + \sum_{m \neq n \in \mathcal{B}_i^{e,t}} 2\mathbb{E}\left[\bar{g}_{im}(\theta)^\top \bar{g}_{in}(\theta)\right] \\
&= \sum_{j \in \mathcal{B}_i^{e,t}} \mathbb{E}\left[\|\bar{g}_{ij}(\theta)\|^2\right] + \sum_{m \neq n \in \mathcal{B}_i^{e,t}} 2\mathbb{E}\left[\bar{g}_{im}(\theta)\right]^\top \mathbb{E}\left[\bar{g}_{in}(\theta)\right] \\
&= b_i^e c^2 + 2 \binom{b_i^e}{2} \|G_i(\theta)\|^2. \quad (8)
\end{aligned}$$

The last equation has used Equation (6) and that we clip the norm of sample gradients $\bar{g}_{ij}(\theta)$ with an “effective” clipping threshold c . By replacing \mathcal{A} into eq. 7, we can rewrite it as:

$$\begin{aligned}
\sigma_{i,\tilde{g}}^2(b_i^e) &:= \text{Var}[\tilde{g}_i^{e,t}(\theta)] = \frac{1}{b_i^{e^2}} \left(\mathbb{E}\left[\left\|\sum_{j \in \mathcal{B}_i^{e,t}} \bar{g}_{ij}(\theta)\right\|^2\right] - b_i^{e^2} \|G_i(\theta)\|^2 \right) + \frac{pc^2 z_i^2(\epsilon_i, \delta_i, b_i^1, b_i^{>1}, N_i, K, E)}{b_i^{e^2}} \\
&= \frac{1}{b_i^{e^2}} \left(b_i^e c^2 + \left(2 \binom{b_i^e}{2} - b_i^{e^2}\right) \|G_i(\theta)\|^2 \right) + \frac{pc^2 z_i^2(\epsilon_i, \delta_i, b_i^1, b_i^{>1}, N_i, K, E)}{b_i^{e^2}} \\
&= \frac{c^2 - \|G_i(\theta)\|^2}{b_i^e} + \frac{pc^2 z_i^2(\epsilon_i, \delta_i, b_i^1, b_i^{>1}, N_i, K, E)}{b_i^{e^2}} \approx \frac{pc^2 z_i^2(\epsilon_i, \delta_i, b_i^1, b_i^{>1}, N_i, K, E)}{b_i^{e^2}} \quad (9)
\end{aligned}$$

The last approximation is valid because $p \gg 1$ (it is the number of model parameters).

Scenario 2: the clipping threshold c is ineffective for all samples in a batch: when the clipping is ineffective for all samples, i.e., $\forall j \in \mathcal{B}_i^{e,t} : c > \|g_{ij}(\theta)\|$, we have a noisy version of the batch gradient $g_i^{e,t}(\theta) = \frac{1}{b_i^e} \sum_{j \in \mathcal{B}_i^{e,t}} g_{ij}(\theta)$, which is unbiased with variance bounded by $\sigma_{i,g}^2(b_i^e)$ (see Assumption 3.2). We note that $\sigma_{i,g}^2(b_i^e)$ is a constant that depends on the used batch size b_i^e . The larger the batch size b_i^e used during round e , the smaller the constant. Hence, in this case:

$$\mathbb{E}[\tilde{g}_i^{e,t}(\theta)] = \mathbb{E}[g_i^{e,t}(\theta)] = \nabla f_i(\theta), \quad (10)$$

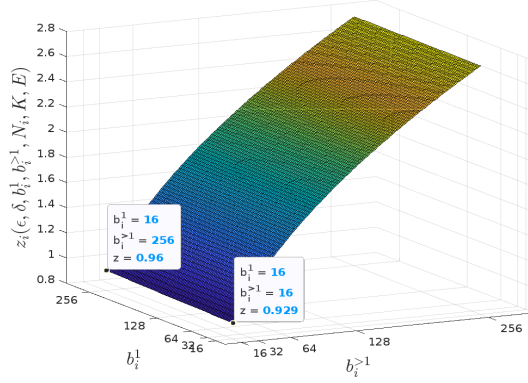


Figure 8: Plot of $z_i(\epsilon, \delta, b_i^1, b_i^{>1}, N_i, K, E)$ v.s. b_i^1 and $b_i^{>1}$ obtained from Renyi-DP Accountant (Mironov et al., 2019) in a setting with $N_i = 6600, \epsilon = 5, \delta = 10^{-4}, K = 1, E = 200$. It is clearly observed that the effect of $b_i^{>1}$ is much more than the effect of b_i^1 . The reason is that $b_i^{>1}$ is used in $E - 1$ rounds, while b_i^1 is used only in the first round. So it is the value of $b_i^{>1}$ that affects z_i the most.

and

$$\begin{aligned} \sigma_{i,\tilde{g}}^2(b_i^e) &= \text{Var}[\tilde{g}_i^{e,t}(\theta)] = \text{Var}[g_i^{e,t}(\theta)] + \frac{p\sigma_{i,\text{DP}}^2}{b_i^{e^2}} \leq \sigma_{i,g}^2(b_i^e) + \frac{p\sigma_{i,\text{DP}}^2}{b_i^{e^2}} \\ &= \sigma_{i,g}^2(b_i^e) + \frac{pc^2 z_i^2(\epsilon, \delta, b_i^1, b_i^{>1}, N_i, K, E)}{b_i^{e^2}} \\ &\approx \frac{pc^2 z_i^2(\epsilon, \delta, b_i^1, b_i^{>1}, N_i, K, E)}{b_i^{e^2}}. \end{aligned} \quad (11)$$

The approximation is valid because $p \gg 1$ (number of model parameters). Also, note that $\sigma_{i,g}^2(b_i^e)$ decreases with b_i^e . Therefore, we got to the same result as in Equation (9).

As observed in see Figure 8, z_i grows with b_i^1 and $b_i^{>1}$ *sub-linearly* (especially with b_i^1). Therefore, the variance of the client i 's DP batch gradients $\tilde{g}_i^{e,t}(\theta)$ during communication round e , decreases with b_i^e fast. The larger the batch size b_i^e , the less the noise existing in its batch gradients during the same round.

With the findings above, we now investigate the effect of batch size b_i^e on **the noise level in clients' model updates at the end of round e** . During the global communication round e , a participating client i performs $E_i^e = K \cdot \lceil \frac{N_i}{b_i^e} \rceil$ batch gradient updates locally with step size η_i :

$$\theta_i^{e,k} = \theta_i^{e,k-1} - \eta_i \tilde{g}_i(\theta_i^{e,k-1}), \quad k = 1, \dots, E_i^e. \quad (12)$$

Hence,

$$\Delta \tilde{\theta}_i^e = \theta_i^{e,E_i^e} - \theta_i^{e,0} \quad (13)$$

In each update, it adds a Gaussian noise from $\mathcal{N}(0, \frac{c^2 z_i^2(\epsilon, \delta, b_i^1, b_i^{>1}, N_i, K, E)}{b_i^{e^2}} \mathbb{I}_p)$ to its batch gradients independently (see Equation (1)). Hence:

$$\text{Var}[\Delta \tilde{\theta}_i^e | \theta_i^{e,0}] = E_i^e \cdot \eta_i^2 \cdot \sigma_{i,\tilde{g}}^2(b_i^e), \quad (14)$$

where $\sigma_{i,\tilde{g}}^2(b_i^e)$ was computed in Equation (9) and Equation (11), and was a decreasing function of b_i^e . Therefore:

$$\text{Var}[\Delta \tilde{\theta}_i^e | \theta_i^{e,0}] \approx K \cdot N_i \cdot \eta_i^2 \cdot \frac{pc^2 z_i^2(\epsilon, \delta, b_i^1, b_i^{>1}, N_i, K, E)}{b_i^{e^3}}. \quad (15)$$

□

E EFFECT OF BATCH SIZE INCREMENT ON CLUSTERING

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 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 $O_{m,m'} = 2Q(\frac{\sqrt{p}\Delta_{m,m'}(b^1)}{2\sigma^1(b^1)})$, where $\sigma^{1^2}(b^1) := \text{Var}[\Delta\tilde{\theta}_i^1 | \theta^{init}, b_i^1 = b^1]$ and $Q(\cdot)$ is the tail distribution function of the standard normal distribution. Furthermore, if we increase $b_i^1 = b^1$ to $b_i^1 = kb^1 \leq N$ (for all i), we have $O_{m,m'} \leq 2Q(\frac{\sqrt{kp}\Delta_{m,m'}(b^1)}{2\rho\sigma^1(b^1)})$, where $1 \leq \rho \in \mathcal{O}(1)$ is a small constant.*

Proof. We first find the overlap between two arbitrary Gaussian distributions. Without loss of generality, let's assume we are in 1-dimensional space and that we have two Gaussian distributions both with variance σ^2 and with means $\mu_1 = 0$ and $\mu_2 = \mu$ ($\|\mu_1 - \mu_2\| = \mu$), respectively. Based on symmetry of the distributions, the two components start to overlap at $x = \frac{\mu}{2}$. Hence, we can find the overlap between the two gaussians as follows:

$$O := 2 \int_{\frac{\mu}{2}}^{\infty} \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{x^2}{2\sigma^2}} dx = 2 \int_{\frac{\mu}{2\sigma}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = 2Q(\frac{\mu}{2\sigma}), \quad (16)$$

where $Q(\cdot)$ is the tail distribution function of the standard normal distribution. Now, let's consider the 2-dimensional space, and consider two similar symmetric distributions centered at $\mu_1 = (0, 0)$ and $\mu_2 = (\mu, 0)$ ($\|\mu_1 - \mu_2\| = \mu$) and with $\Sigma_1 = \Sigma_2 = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$. The overlap between the two gaussians can be found as:

$$O = 2 \int_{-\infty}^{\infty} \int_{\frac{\mu}{2}}^{\infty} \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}} dx dy = 2 \int_{\frac{\mu}{2}}^{\infty} \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{x^2}{2\sigma^2}} dx \cdot \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{y^2}{2\sigma^2}} dy = 2Q(\frac{\mu}{2\sigma}). \quad (17)$$

If we compute the overlap for two similar symmetric p -dimensional distributions with $\|\mu_1 - \mu_2\| = \mu$ and variance σ^2 in every direction, we will get to the same result $2Q(\frac{\mu}{2\sigma})$.

In the lemma, when using batch size b^1 , we have two Gaussian distributions $\mathcal{N}(\mu_m^*(b^1), \Sigma_m^*(b^1))$ and $\mathcal{N}(\mu_{m'}^*(b^1), \Sigma_{m'}^*(b^1))$, where

$$\Sigma_m^*(b^1) = \Sigma_{m'}^*(b^1) = \begin{bmatrix} \frac{\sigma^{1^2}(b^1)}{p} & & \\ & \ddots & \\ & & \frac{\sigma^{1^2}(b^1)}{p} \end{bmatrix}. \quad (18)$$

Therefore, from Equation (17), we can immediately conclude that the overlap between the two Gaussians, which we denote with $O_{m,m'}(b^1)$, is:

$$O_{m,m'}(b^1) = 2Q(\frac{\sqrt{p}\Delta_{m,m'}(b^1)}{2\sigma^1(b^1)}), \quad (19)$$

which proves the first part of the lemma.

Now, let's see the effect of increasing batch size. First, note that we had:

$$\begin{aligned} \Delta\tilde{\theta}_i^1 &= \theta_i^{1,E_i^1} - \theta_i^{1,0}, \\ \theta_i^{1,k} &= \theta_i^{1,k-1} - \eta_i \tilde{g}_i(\theta_i^{1,k-1}), \quad k = 1, \dots, E_i^1, \end{aligned} \quad (20)$$

where $E_i^1 = K \cdot \lceil \frac{N}{b^1} \rceil$ is the total number of gradients steps taken by client i during communication round $e = 1$. Therefore, considering that DP batch gradients are clipped with a bound c , we have:

$$\|\mathbb{E}[\Delta\tilde{\theta}_i^1(b^1)]\| \leq E_i^1 \cdot \eta_l \cdot c. \quad (21)$$

When we increase batch size b_i^1 for all clients from b^1 to kb^1 , the upperbound in Equation (21) gets k times smaller. In fact by doing so, the number of local gradient updates that client i performs during round $e = 1$, which is equal to E_i^1 , decreases k times. As such, we can write:

$$\Delta\tilde{\theta}_i^1(b^1) = k \cdot \Delta\tilde{\theta}_i^1(kb^1) + v_i, \quad (22)$$

where $v_i \in \mathbb{R}^p$ is a vector capturing the discrepancies between $\Delta\tilde{\theta}_i^1(b^1)$ and $k \cdot \Delta\tilde{\theta}_i^1(kb^1)$. Therefore, we have:

$$\begin{aligned} \mu_m^*(b^1) &= \mathbb{E}[\Delta\tilde{\theta}_i^1(b^1)|s(i) = m] = \mathbb{E}[k \cdot \Delta\tilde{\theta}_i^1(kb^1) + v_i|s(i) = m] \\ &= k \cdot \mathbb{E}[\Delta\tilde{\theta}_i^1(kb^1)] + \mathbb{E}[v_i|s(i) = m] = k \cdot \mu_m^*(kb^1) + \mathbb{E}[v_i|s(i) = m]. \end{aligned} \quad (23)$$

Therefore, we have:

$$\|\mu_m^*(b^1) - \mu_{m'}^*(b^1)\| = \left\| k\mu_m^*(kb^1) - k\mu_{m'}^*(kb^1) + \left(\mathbb{E}[v_i|s(i) = m] - \mathbb{E}[v_i|s(i) = m'] \right) \right\|. \quad (24)$$

Based on our experiments, the last term above, in parenthesis, is small and we can have the following approximation for the equation above:

$$\|\mu_m^*(b^1) - \mu_{m'}^*(b^1)\| \approx \|k\mu_m^*(kb^1) - k\mu_{m'}^*(kb^1)\|, \quad (25)$$

or equivalently:

$$\|\mu_m^*(kb^1) - \mu_{m'}^*(kb^1)\| \approx \frac{\|\mu_m^*(b^1) - \mu_{m'}^*(b^1)\|}{k}. \quad (26)$$

Figure 9 (left) shows the validity of the approximation above with some experimental results. On the other hand, from Equation (2) and also noting that a client, with dataset size N and batch size b^1 , takes $\frac{N}{b^1}$ gradient steps during each epoch of the first round, we have:

$$\forall m \in [M] : \sigma_m^2(b^1) = \sigma^2(b^1) \approx K \cdot N \cdot \eta_l^2 \cdot \frac{pc^2 z^2(\epsilon, \delta, b^1, b^{>1}, N, K, E)}{b^{1^3}}. \quad (27)$$

When we change the batch size used during the first communication round $e = 1$ from b^1 to kb^1 and we fix the batch size of rounds $e > 1$, then the noise scale z changes from $z(\epsilon, \delta, b^1, b^{>1}, N_i, K, E)$ to $z(\epsilon, \delta, kb^1, b^{>1}, N_i, K, E)$. Confirmed by our experimental analysis (see Figure 9, right), the amount of change in z due to this is small, as we have changed the batch size only in the first round $e = 1$ from b^1 to kb^1 , while the batch sizes in the other $E - 1$ rounds are unchanged and $E \gg 1$. Therefore, supported by the results in Figure 9, we can always establish an upper bound on the amount of change in z as b^1 increases: $z(\epsilon, \delta, kb^1, b^{>1}, N, K, E) \leq \rho z(\epsilon, \delta, b^1, b^{>1}, N, K, E)$, where ρ is a small constant (e.g. $\rho = 2.5$ in Figure 9). So we have:

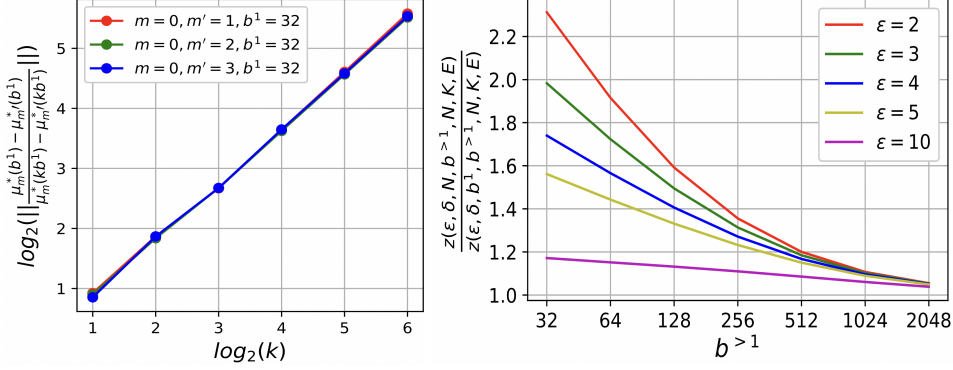


Figure 9: **Left:** Distance between the centers of different clusters, i.e., the distance between $\mu_m^*(b^1)$ and $\mu_{m'}^*(b^1)$, decreases k times as b^1 increases k times. The three curves in the plot are obtained on CIFAR10 with 4 clusters $m \in \{0, 1, 2, 3\}$ obtained from covariate shift (rotation). The curves are overlapping all with slope 0.95, which is very close to 1. This shows the validity of the approximation in Equation (26). **Right:** Effect of changing batch size b^1 to full batch size in the first round on the noise scale z . In the denominator, b^1 is equal to $b^{>1}$. Results are obtained from Renyi-DP accountant (Mironov et al., 2019) with $N = 50000$, $K = 1$ and $E = 200$. For each value of ϵ , we have shown the results for seven values of $b^{>1}$.

$$\begin{aligned}
 \forall m \in [M] : \sigma_m^2(kb^1) &= \sigma^2(kb^1) \approx K \cdot N \cdot \eta_l^2 \cdot \frac{pc^2 z^2(\epsilon, \delta, kb^1, b^{>1}, N, K, E)}{(kb^1)^3} \\
 &\leq K \cdot N \cdot \eta_l^2 \cdot \frac{pc^2 \rho^2 z^2(\epsilon, \delta, b^1, b^{>1}, N, K, E)}{(kb^1)^3} \\
 &= \frac{\rho^2 \sigma^2(b^1)}{k^3}.
 \end{aligned} \tag{28}$$

From Equation (26) and Equation (28), we have:

$$O_{m,m'}(kb^1) = 2Q\left(\frac{\sqrt{p}\Delta_{m,m'}(kb^1)}{2\sigma(kb^1)}\right) \leq 2Q\left(\frac{\sqrt{p}\frac{\Delta_{m,m'}(b^1)}{k}}{2\frac{\rho\sigma(b^1)}{k^{\frac{3}{2}}}}\right) = 2Q\left(\frac{\sqrt{kp}\Delta_{m,m'}(b^1)}{2\rho\sigma(b^1)}\right), \tag{29}$$

which completes the proof. \square

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:

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

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 $x \rightarrow 0$: $\lim_{x \rightarrow 0} \frac{o(x)}{x} = 0$.

Proof. The proof directly follows from the proof of Theorem 1 in Ma et al. (2000) by considering $\{\Delta\tilde{\theta}_i^1(b^1)\}_{i=1}^n$ as the samples of Gaussian mixture $\{\mathcal{N}(\mu_m^*(b^1), \Sigma_m^*(b^1)), \alpha_m^*\}_{m=1}^M$. \square

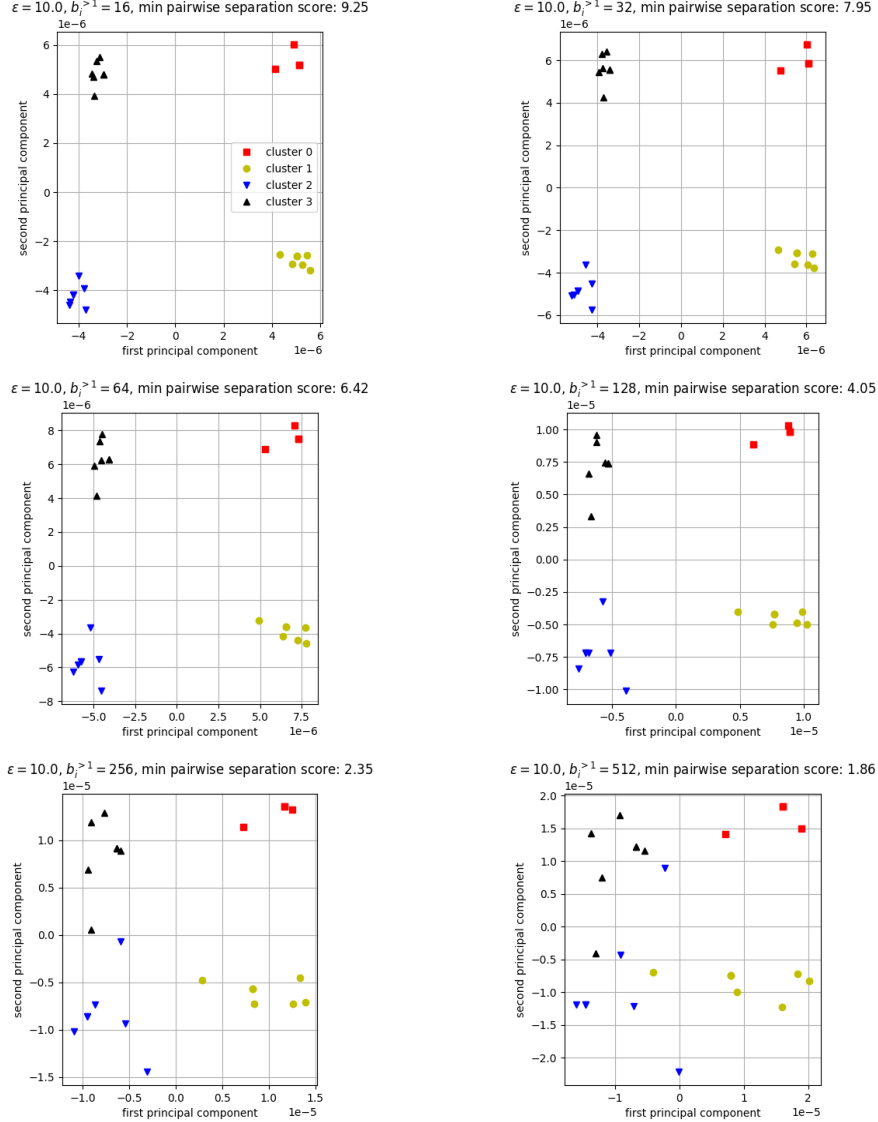


Figure 10: The effect of increasing the batch size after the first round, i.e., $b_i^{>1}$, on the model updates $\{\Delta\tilde{\theta}_i^1\}_{i=1}^n$ at the end of the first round. All clients have used full batch sizes in the first round, i.e., $\forall i : b_i^1 = N_i$. The level of noise in $\{\Delta\tilde{\theta}_i^1\}_{i=1}^n$ affects the quality and confidence of the client clustering that the server performs at the end of the first round. As can be observed, for a fixed $\epsilon = 10$, the model updates scatter further in space as $b_i^{>1}$ increases and different clusters get less separated. This leads to an decrement in the confidence level or the “minimum pairwise separation score (MSS)” of the resulting GMM, as $b_i^{>1}$ increases (see Section 4.2.2 for explanations about the score). The MSS score for each $b_i^{>1}$ is mentioned on the top of the corresponding plot. All the results are obtained on CIFAR10 with covariate shift (rotation) across clusters.

F SETTING HYPER-PARAMETERS OF RC-DPFL

F.1 THE EFFECT OF $b_i^{>1}$ ON CLUSTERING

As we observed in Lemma 4.1 and Figure 3 left, $\text{Var}(\Delta\tilde{\theta}_i^1(b_i^1)|\theta^{init})$, which affects the clustering done at the end of the first round, is an increasing function of $b_i^{>1}$. More generally, we can describe the effect of increasing $b_i^{>1}$ in three things: 1) increasing the noise variance in $\text{Var}(\Delta\tilde{\theta}_i^1(b_i^1)|\theta^{init})$ (as shown in Figure 3, right) 2) decreasing the noise variance in $\text{Var}(\Delta\tilde{\theta}_i^1(b_i^e)|\theta_i^{e,0})$ (as shown in Figure 3, right) 3) decreasing the number of gradients steps during each round e for $e > 1$.

While first one is only limited two the first round $e = 1$, while the last two affect the remaining $E - 1$ rounds $e > 1$ and have conflicting effects on the final accuracy. However, an important point about the problem of clustered DPFL is that finding the true structure of clusters in the first round is a prerequisite for making progress in the next rounds. Therefore, the first effect of increasing the noise variance in $\text{Var}(\Delta\tilde{\theta}_i^1(b_i^1)|\theta^{init})$ is more important and is the most undesirable result. We have demonstrated this effect in Figure 10, which shows that how increasing $b_i^{>1}$ adversely affects the clustering done at the end of the first round. Note how MSS of the learned GMM increases as $b_i^{>1}$ increases. Therefore, in order to have a reliable client clustering at the end of the first round, we need to keep the value of $b_i^{>1}$ as small as possible. Following this observation, we have fixed $b_i^{>1}$ to 32 in all our experimental results, and RC-DPFL was able to outperform the existing baseline algorithms.

F.2 FINDING THE NUMBER OF CLUSTERS (M)

Knowing the number of clusters is broadly accepted and applied in the clustered FL literature (Ghosh et al., 2020; Ruan & Joe-Wong, 2021; Briggs et al., 2020). This is the assumption of our baseline algorithms too. Yet, techniques to determine the number of clusters can enable our approach to be more widely adopted. In this section, we show that how we can find the true number of clusters (M) when it is not given. Our method, which results in a high accuracy, relies on the MSS score (confidence level) defined in Section 4.2.2: $\text{MSS} = \min_{m,m'} \hat{S}(m, m') \in [0, +\infty)$, where $\hat{S}(m, m') \approx \frac{\sqrt{p}\Delta_{m,m'}(b^1)}{2\sigma^1(b^1)} = \frac{\Delta_{m,m'}(b^1)}{2\sigma^1(b^1)/\sqrt{p}}$ (please see the detailed explanations in Section 4.2.2). Consider the Figure 2 right as an example. There is a good separation between the $M = 4$ existing clusters, thanks to using full batch size b^1 in the first round. Fitting a GMM with 4 components to the model updates results in the highest MSS for the learned GMM model: remember that MSS was the maximum pairwise separation score between the different components of the learned GMM. In contrast, if we fit a GMM with 3 components (less than the true number of components) to the same model updates in the figure, then two clusters will be merged into one component (for examples clusters 0 and 1) leading to a high radius ($\sigma^1(b^1)/\sqrt{p}$) for one of the three components of the resulting GMM. This leads to a low MSS (confidence level) for the resulting GMM. Similarly, if we fit a GMM with 5 components, one of the four clusters (for example cluster 1) will be split between two of the 5 components (call them m and m'), which leads to a low inter-component distance ($\Delta_{m,m'}(b^1)$) for the pair of components. This also leads to a low MSS for the resulting GMM. However, fitting a GMM with $M = 4$ components leads to a well separation between all the true components and maximizes the resulting MSS. Based on this very intuitive observation, we propose the following method for setting m at the end of the first round. We select the number of clusters/components, which leads to the maximum MSS for the resulting GMM. More specifically:

$$M = \arg \max_{m \in S} \text{MSS}(\text{GMM}(\Delta\tilde{\theta}_1^1, \dots, \Delta\tilde{\theta}_n^1; m)), \text{ (line 9 of Algorithm 1)} \quad (30)$$

where S is a set of candidate values for M : At the end of the first round and on the server, we learn one GMM for each candidate value in S on the same received model updates $\{\Delta\tilde{\theta}_i^1\}_{i=1}^n$. Finally, we choose the value resulting in the GMM with the highest MSS (confidence). It is noteworthy that we know from Lemma 4.2 that learning the GMM does not incur much computational cost.

We have evaluated this method on multiple data splits and different privacy budgets (ϵ) on CIFAR10, MNIST and FMNIST, and it worked perfectly, as shown in Figure 11. As can be observed, the method has made only one mistake for $\epsilon = 4$ (seed 1) and two mistakes for $\epsilon = 3$ (seeds 0 and 1), out of 20 total experiments. Even in those three cases, it has predicted M as 5, which is closest to the

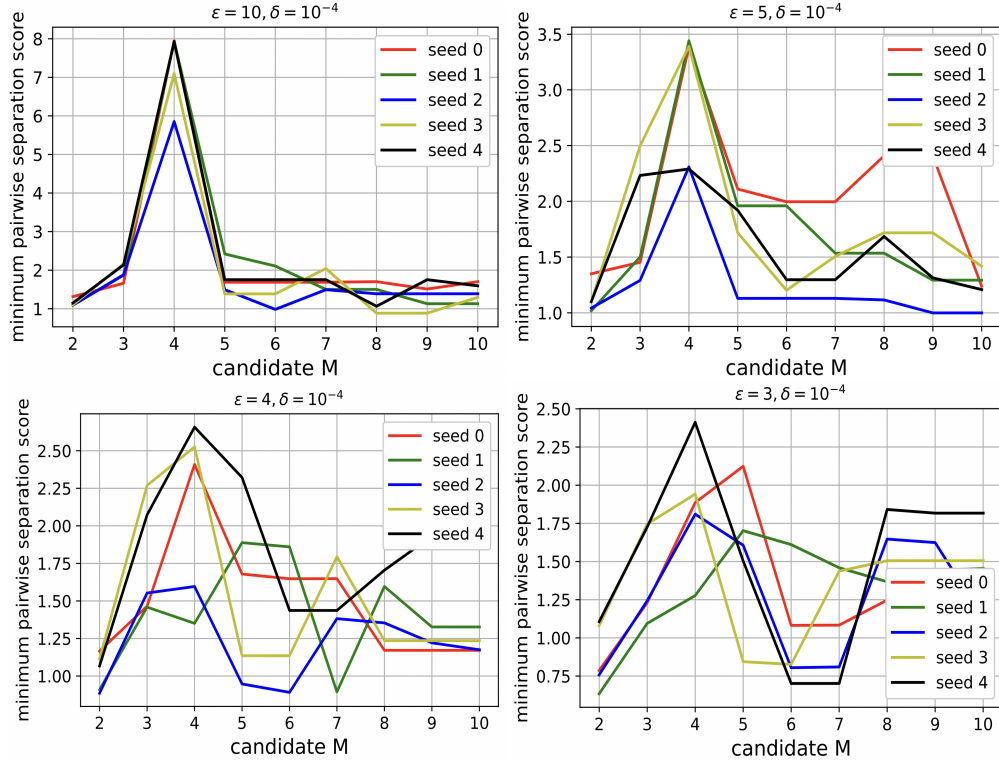


Figure 11: The minimum pairwise separation score (MSS) or confidence of the GMM learned on $\{\Delta\hat{\theta}_i^1\}_{i=1}^n$ peaks at the true cluster number, which is equal to 4 in all the plots above. Each figure is for a different value of ϵ (mentioned on top of each figure), and are obtained on CIFAR10 with covariate shift (rotation) across clusters, and 5 different random data splits (5 seeds). All the results are obtained with full batch sizes in the first round and $b_i^{>1} = 32$ for all i . We can use this observation as a method to find the true number of clusters (M) when it is not given. For larger ϵ , this method work perfectly and even when ϵ is too small, e.g., $\epsilon = 3$, this method works well and predicts the true number of clusters correctly most of the times: 3 out of the 5 curves in the bottom right plot have a peak at $M = 4$ (the true cluster number), and the other 2 curves predict 5 as the true number, which is the closest and the best alternative for the true value $M = 4$.

true value ($M = 4$) and does not lead to much performance drop. This is because having $M = 5$ splits an existing cluster into two and it is better than predicting for example $M = 3$, which results in "mixing" two clusters with heterogeneous data. **Same method could predict the number of underlying clusters with 100% accuracy for the MNIST and FMNIST datasets for all values of ϵ .** Finally, note that none of the existing baseline algorithms has such an easy and applicable strategy for finding M . This shows another useful feature of the proposed RC-DPFL.

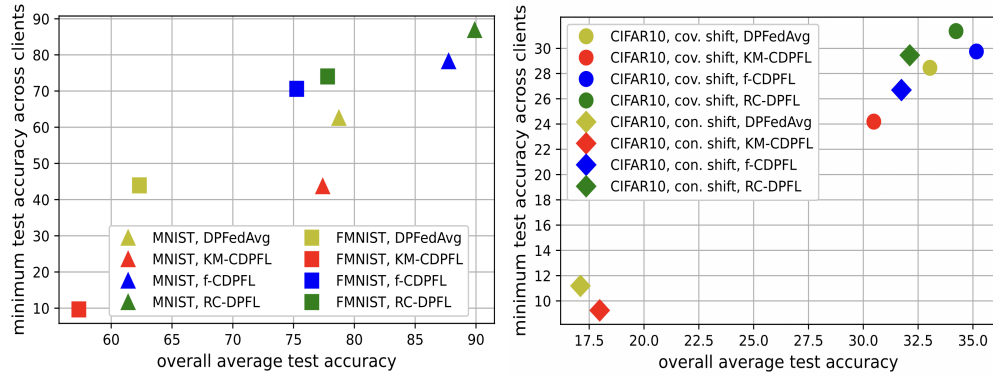


Figure 12: Pareto frontier for overall utility and fairness (in terms of the minimum test accuracy across clients). All the results are for $\epsilon = 4$. The results on MNIST and FMNIST (left) are both on data splits with covariate shift. For other values of ϵ , we observe similar results, as observed in Table 4 to Table 11.

G DETAILED EXPERIMENTAL RESULTS

G.1 MNIST DATA SPLIT WITH COVARIATE SHIFT

Table 4: Average test accuracy of all (All), majority (Maj) and minority (Min) clients on the split of MNIST dataset with covariate shift (rotation). The results are averaged over 3 seeds.

algorithm		$\epsilon = 10$	$\epsilon = 5$	$\epsilon = 4$	$\epsilon = 3$	$\epsilon = 2$	
MNIST	DPFedAvg	All:	80.92 \pm 0.83	79.37 \pm 0.76	78.73 \pm 1.04	72.88 \pm 0.79	72.29 \pm 0.77
		Maj:	82.47 \pm 0.01	81.16 \pm 0.01	80.35 \pm 0.05	75.45 \pm 0.02	75.02 \pm 0.05
		Min:	67.16 \pm 1.37	64.04 \pm 2.06	63.29 \pm 1.83	52.80 \pm 2.24	52.50 \pm 2.76
	f-CDPFL	All:	90.64 \pm 1.45	89.12 \pm 1.69	87.75 \pm 0.99	84.56 \pm 3.91	85.01 \pm 2.71
		Maj:	92.72\pm1.42	91.25\pm0.01	85.93 \pm 0.02	81.75 \pm 0.3	88.80 \pm 0.03
		Min:	86.98 \pm 5.77	79.12 \pm 8.94	81.36 \pm 6.03	88.12 \pm 1.02	76.07 \pm 17.79
	KM-CDPFL	All:	81.42 \pm 3.03	77.80 \pm 0.47	77.40 \pm 2.06	71.91 \pm 2.29	68.89 \pm 1.20
		Maj:	81.55 \pm 0.01	79.20 \pm 0.01	80.60 \pm 0.07	74.38 \pm 0.02	73.11 \pm 0.02
		Min:	78.78 \pm 0.08	59.01 \pm 5.07	52.57 \pm 9.51	42.33 \pm 5.2	35.13 \pm 1.81
	RC-DPFL	All:	92.28\pm0.40	90.7\pm0.47	89.89\pm0.43	89.00\pm0.44	87.92\pm0.36
		Maj:	92.06 \pm 0.01	90.63 \pm 1.42	89.95\pm0.01	88.78\pm0.02	87.76\pm0.03
		Min:	90.45\pm0.19	89.33\pm0.24	87.44\pm0.54	88.21\pm0.67	85.83\pm0.26
	Oracle-CDPFL	All:	92.71 \pm 0.27	91.06 \pm 0.06	91.09 \pm 0.53	89.33 \pm 0.44	88.94 \pm 0.47
		Maj:	92.70 \pm 0.01	91.10 \pm 0.03	90.91 \pm 0.02	89.10 \pm 0.05	88.71 \pm 1.42
		Min:	91.11 \pm 0.25	90.28 \pm 0.02	89.87 \pm 0.46	89.10 \pm 0.089	88.23 \pm 0.52

Table 5: Fairness evaluation in terms of different metrics on the split of MNIST dataset with covariate shift (rotation). The results are averaged over 3 seeds.

MNIST

algorithm		$\epsilon = 10$	$\epsilon = 5$	$\epsilon = 4$	$\epsilon = 3$	$\epsilon = 2$
DPFedAvg	\mathcal{F}_{acc} :	18.41 \pm 0.80	20.81 \pm 2.18	20.74 \pm 1.46	27.93 \pm 1.68	27.57 \pm 1.89
	\mathcal{F}_{loss} :	1.18 \pm 0.16	1.49 \pm 0.26	1.54 \pm 0.25	1.53 \pm 0.36	1.61 \pm 0.32
	Min Acc:	66.33 \pm 1.28	63.26 \pm 2.25	62.58 \pm 2.20	51.55 \pm 2.88	51.46 \pm 3.04
	Accuracy Disparity:	18.56 \pm 0.57	21.06 \pm 2.26	21.00 \pm 1.50	28.25 \pm 1.75	27.93 \pm 1.87
f-CDPFL	\mathcal{F}_{acc} :	8.70 \pm 5.52	13.39 \pm 9.28	13.60 \pm 3.07	10.15 \pm 5.75	17.81 \pm 15.77
	\mathcal{F}_{loss} :	0.60 \pm 0.21	0.81 \pm 0.19	0.86 \pm 0.11	0.57 \pm 0.11	0.89 \pm 0.17
	Min Acc:	84.81 \pm 4.98	78.76 \pm 9.02	78.25 \pm 2.71	79.75 \pm 6.44	72.20 \pm 15.79
	Accuracy Disparity:	9.05 \pm 5.33	13.70 \pm 9.12	13.83 \pm 3.07	10.21 \pm 5.70	17.88 \pm 15.70
KM-CDPFL	\mathcal{F}_{acc} :	18.67 \pm 6.54	41.43 \pm 10.58	46.69 \pm 9.40	46.68 \pm 4.56	54.17 \pm 2.67
	\mathcal{F}_{loss} :	1.48 \pm 0.31	2.73 \pm 0.24	3.33 \pm 0.23	2.33 \pm 0.24	3.26 \pm 0.58
	Min Acc:	73.33 \pm 6.95	48.66 \pm 10.04	43.71 \pm 9.25	40.58 \pm 5.24	32.86 \pm 2.57
	Accuracy Disparity:	18.81 \pm 6.45	41.58 \pm 10.67	47.08 \pm 9.58	46.98 \pm 4.43	54.48 \pm 2.69
RC-DPFL	\mathcal{F}_{acc} :	4.04\pm0.90	3.13\pm0.36	4.27\pm0.80	2.95\pm0.84	4.40\pm1.22
	\mathcal{F}_{loss} :	0.55\pm0.15	0.46\pm0.16	0.68\pm0.06	0.35\pm0.14	0.52\pm0.15
	Min Acc:	89.46\pm0.28	88.71\pm0.14	86.91\pm0.52	87.43\pm0.81	85.23\pm0.54
	Accuracy Disparity:	4.41\pm0.81	3.36\pm0.23	4.53\pm0.59	2.99\pm0.64	4.48\pm1.00
Oracle-CDPFL	\mathcal{F}_{acc} :	3.35 \pm 0.81	3.10 \pm 0.73	2.93 \pm 0.49	2.92 \pm 0.75	2.85 \pm 0.66
	\mathcal{F}_{loss} :	0.43 \pm 0.12	0.36 \pm 0.08	0.43 \pm 0.08	0.32 \pm 0.05	0.36 \pm 0.08
	Min Acc:	90.33 \pm 0.45	89.36 \pm 0.37	89.13 \pm 0.28	87.73 \pm 0.81	87.43 \pm 0.86
	Accuracy Disparity:	3.71 \pm 0.79	3.13 \pm 0.51	3.21 \pm 0.38	2.91 \pm 0.74	2.93 \pm 0.47

G.2 FMNIST DATA SPLIT WITH COVARIATE SHIFT

Table 6: Average test accuracy of all (All), majority (Maj) and minority (Min) clients on the split of FMNIST dataset with covariate shift (rotation). The results are averaged over 3 seeds.

algorithm		$\epsilon = 10$	$\epsilon = 5$	$\epsilon = 4$	$\epsilon = 3$	$\epsilon = 2$	
FMNIST	DPFedAvg	All:	62.62 \pm 0.39	62.22 \pm 0.35	62.32 \pm 0.93	59.58 \pm 0.25	59.78 \pm 0.25
		Maj:	66.27 \pm 0.05	65.86 \pm 0.05	65.86 \pm 0.03	62.90 \pm 7.10	62.96 \pm 0.03
		Min:	42.32 \pm 1.55	41.61 \pm 1.18	44.84 \pm 6.62	38.15 \pm 0.73	43.23 \pm 3.91
	f-CDPFL	All:	78.19 \pm 1.53	74.37 \pm 0.46	75.26 \pm 2.38	70.10 \pm 1.30	70.42 \pm 3.84
		Maj:	76.12 \pm 0.05	73.43 \pm 0.03	72.48 \pm 0.03	67.82 \pm 0.04	67.35 \pm 0.06
		Min:	76.21 \pm 4.92	73.43 \pm 6.62	77.89\pm0.96	69.85\pm7.48	69.88 \pm 7.63
	KM-CDPFL	All:	60.17 \pm 2.15	57.65 \pm 1.09	57.34 \pm 0.19	54.82 \pm 2.95	53.55 \pm 0.47
		Maj:	57.79 \pm 0.02	58.75 \pm 0.03	63.22 \pm 7.10	53.31 \pm 0.03	58.30 \pm 0.06
		Min:	59.81 \pm 4.34	31.72 \pm 20.48	15.95 \pm 2.90	32.33 \pm 2.78	19.70 \pm 7.14
	RC-DPFL	All:	80.26\pm0.22	78.15\pm0.23	77.80\pm0.50	74.04\pm1.54	74.84\pm0.43
		Maj:	80.50\pm0.02	78.41\pm0.02	78.12\pm0.06	75.11\pm0.06	74.93\pm0.04
		Min:	79.03\pm1.57	77.23\pm1.30	76.21 \pm 1.96	69.35 \pm 7.17	73.29\pm1.55
	Oracle-CDPFL	All:	80.46 \pm 0.10	78.44 \pm 0.18	78.31 \pm 0.22	75.19 \pm 0.16	75.11 \pm 0.17
		Maj:	80.76 \pm 0.02	78.42 \pm 0.03	78.31 \pm 0.3	75.02 \pm 0.05	75.21 \pm 0.03
		Min:	79.76 \pm 0.96	78.35 \pm 0.62	77.81 \pm 1.02	75.33 \pm 0.53	75.00 \pm 1.01

Table 7: Fairness evaluation in terms of different metrics on the split of FMNIST dataset with covariate shift (rotation). The results are averaged over 3 seeds.

FMNIST

algorithm		$\epsilon = 10$	$\epsilon = 5$	$\epsilon = 4$	$\epsilon = 3$	$\epsilon = 2$	
DPFedAvg	\mathcal{F}_{acc} :	25.07 \pm 2.16	24.86 \pm 1.44	23.31 \pm 6.22	26.53 \pm 2.35	22.97 \pm 4.59	
	\mathcal{F}_{loss} :	1.26 \pm 0.05	1.02 \pm 0.08	0.99 \pm 0.24	0.83 \pm 0.08	0.77 \pm 0.18	
	Min Acc:	41.60 \pm 1.55	40.69 \pm 1.34	43.93 \pm 6.45	37.30 \pm 1.02	42.22 \pm 4.08	
	Accuracy Disparity:	27.86 \pm 0.98	28.11 \pm 1.27	25.83 \pm 5.02	29.58 \pm 0.70	25.12 \pm 3.31	
	f-CDPFL	\mathcal{F}_{acc} :	11.59 \pm 1.16	11.43 \pm 1.22	12.00 \pm 0.80	18.13 \pm 2.15	12.87 \pm 1.59
		\mathcal{F}_{loss} :	0.32 \pm 0.09	0.33 \pm 0.05	0.37 \pm 0.03	0.42 \pm 0.07	0.33 \pm 0.04
		Min Acc:	72.31 \pm 4.33	66.73 \pm 2.89	70.61 \pm 4.12	61.19 \pm 2.26	64.93 \pm 6.69
		Accuracy Disparity:	9.36 \pm 4.22	13.11 \pm 2.64	9.60 \pm 3.85	15.08 \pm 2.01	11.08 \pm 4.86
	KM-CDPFL	\mathcal{F}_{acc} :	43.55 \pm 2.16	59.20 \pm 6.77	63.41 \pm 1.25	48.46 \pm 4.18	54.50 \pm 6.54
		\mathcal{F}_{loss} :	2.44 \pm 0.25	3.36 \pm 0.99	3.33 \pm 0.35	1.89 \pm 0.42	2.28 \pm 0.40
		Min Acc:	33.33 \pm 2.08	12.41 \pm 7.07	9.66 \pm 0.23	21.03 \pm 2.39	14.03 \pm 7.47
		Accuracy Disparity:	43.83 \pm 2.48	62.78 \pm 7.73	66.08 \pm 0.87	51.40 \pm 2.53	57.91 \pm 8.10
	RC-DPFL	\mathcal{F}_{acc} :	7.37\pm2.70	7.25\pm3.29	7.73\pm2.95	10.55\pm5.18	6.92\pm2.69
		\mathcal{F}_{loss} :	0.28\pm0.06	0.24\pm0.05	0.30\pm0.08	0.22\pm0.07	0.27\pm0.02
		Min Acc:	77.31\pm1.07	75.16\pm1.04	74.03\pm1.29	67.66\pm6.64	71.76\pm1.62
		Accuracy Disparity:	5.18\pm1.43	5.5\pm1.14	6.46\pm1.10	9.21\pm6.10	5.66\pm1.29
	Oracle-CDPFL	\mathcal{F}_{acc} :	7.37 \pm 2.90	6.27 \pm 2.96	6.50 \pm 3.19	6.51 \pm 3.01	6.55 \pm 3.41
		\mathcal{F}_{loss} :	0.24 \pm 0.03	0.24 \pm 0.07	0.24 \pm 0.06	0.20 \pm 0.08	0.20 \pm 0.08
		Min Acc:	77.70 \pm 0.64	76.61 \pm 0.59	75.98 \pm 0.87	73.11 \pm 0.33	72.71 \pm 0.68
		Accuracy Disparity:	5.05 \pm 1.32	3.91 \pm 0.49	4.61 \pm 0.77	4.08 \pm 0.02	4.75 \pm 0.35

G.3 CIFAR10 DATA SPLIT WITH COVARIATE SHIFT

Table 8: Average test accuracy of all (All), majority (Maj) and minority (Min) clients on the split of CIFAR10 dataset with covariate shift (rotation). The results are averaged over 3 seeds.

algorithm		$\epsilon = 10$	$\epsilon = 5$	$\epsilon = 4$	$\epsilon = 3$	$\epsilon = 2$	
CIFAR10 (covariate shift)	DPFedAvg	All:	36.50 \pm 0.10	33.64 \pm 0.01	33.04 \pm 0.07	28.55 \pm 0.01	28.12 \pm 0.18
		Maj:	37.39 \pm 0.02	34.18 \pm 0.02	33.69 \pm 0.05	29.09 \pm 0.01	29.00 \pm 0.04
		Min:	32.10 \pm 0.03	30.03 \pm 0.14	29.68 \pm 0.15	25.84 \pm 0.24	25.48 \pm 0.47
	f-CDPFL	All:	40.20\pm0.50	35.72\pm0.01	35.16\pm0.01	30.38\pm0.21	29.17\pm0.07
		Maj:	41.00\pm0.04	36.62\pm0.02	36.03\pm0.05	30.81\pm0.08	30.05\pm0.11
		Min:	36.80 \pm 1.53	31.16 \pm 0.61	30.83 \pm 1.32	27.58 \pm 1.47	26.05 \pm 0.13
	KM-CDPFL	All:	34.86 \pm 0.23	31.46 \pm 0.12	30.48 \pm 0.84	26.48 \pm 0.01	26.10 \pm 0.08
		Maj:	36.05 \pm 0.12	32.47 \pm 0.15	32.02 \pm 0.04	27.16 \pm 0.01	26.79 \pm 0.12
		Min:	27.69 \pm 0.23	25.43 \pm 0.08	25.20 \pm 0.26	22.43 \pm 0.01	21.93 \pm 0.13
	RC-DPFL	All:	39.97 \pm 0.09	35.20 \pm 0.25	34.23 \pm 0.23	29.92 \pm 0.50	27.78 \pm 0.69
		Maj:	40.52 \pm 0.12	35.97 \pm 0.04	35.05 \pm 0.05	30.52 \pm 0.05	28.80 \pm 0.05
		Min:	37.75\pm0.10	32.76\pm0.26	32.26\pm0.54	27.85\pm1.12	26.13\pm0.05
	Oracle-CDPFL	All:	40.43 \pm 0.17	36.39 \pm 0.20	35.46 \pm 0.43	30.79 \pm 0.19	29.34 \pm 0.61
		Maj:	41.07 \pm 0.05	36.95 \pm 0.05	36.47 \pm 0.11	31.12 \pm 0.05	30.41 \pm 3.55
		Min:	37.84 \pm 0.32	34.49 \pm 0.33	33.07 \pm 0.39	29.63 \pm 0.76	27.38 \pm 1.12

Table 9: Fairness evaluation in terms of different metrics on the split of CIFAR10 dataset with covariate shift (rotation). The results are averaged over 3 seeds.

algorithm		$\epsilon = 10$	$\epsilon = 5$	$\epsilon = 4$	$\epsilon = 3$	$\epsilon = 2$	
CIFAR10 (covariate shift)	DPFedAvg	\mathcal{F}_{acc} :	9.31 \pm 0.76	9.32 \pm 1.06	9.81 \pm 0.59	9.13 \pm 2.13	9.11 \pm 2.15
		\mathcal{F}_{loss} :	0.16 \pm 0.01	0.14 \pm 0.01	0.12 \pm 0.01	0.08 \pm 0.01	0.07 \pm 0.01
		Min Acc:	31.65 \pm 0.02	28.81 \pm 0.05	28.45 \pm 0.28	24.79 \pm 0.19	24.65 \pm 0.33
		Accuracy Disparity:	8.57 \pm 0.17	8.14 \pm 0.45	8.09 \pm 1.10	7.63 \pm 0.81	7.07 \pm 0.33
	f-CDPFL	\mathcal{F}_{acc} :	7.53 \pm 0.73	7.90 \pm 0.53	7.79 \pm 0.75	7.71 \pm 1.36	7.61 \pm 0.9
		\mathcal{F}_{loss} :	0.16 \pm 0.04	0.19 \pm 0.02	0.18 \pm 0.05	0.12 \pm 0.01	0.09\pm0.02
		Min Acc:	35.97 \pm 1.47	30.31 \pm 0.56	29.75 \pm 0.76	27.43 \pm 1.04	25.01 \pm 0.93
		Accuracy Disparity:	7.71 \pm 1.50	8.29 \pm 1.13	7.57 \pm 0.76	6.55 \pm 1.30	6.89 \pm 1.52
	KM-CDPFL	\mathcal{F}_{acc} :	10.88 \pm 0.28	9.30 \pm 0.97	8.90 \pm 0.421	9.18 \pm 0.61	8.76 \pm 0.64
		\mathcal{F}_{loss} :	0.32 \pm 0.07	0.24 \pm 0.01	0.22 \pm 0.02	0.13 \pm 0.01	0.12 \pm 0.02
		Min Acc:	26.87 \pm 0.72	24.65 \pm 0.08	24.19 \pm 0.82	21.05 \pm 0.23	20.87 \pm 0.12
		Accuracy Disparity:	14.21 \pm 0.65	11.63 \pm 0.07	10.79 \pm 2.05	9.71 \pm 0.23	9.35 \pm 0.13
	RC-DPFL	\mathcal{F}_{acc} :	7.4\pm0.68	7.49\pm0.74	7.05\pm1.10	7.04\pm1.76	7.03\pm1.17
		\mathcal{F}_{loss} :	0.10\pm0.03	0.10\pm0.01	0.10\pm0.01	0.10\pm0.01	0.68 \pm 0.01
		Min Acc:	37.29\pm0.11	31.95\pm0.28	31.37\pm0.33	27.81\pm0.98	25.37\pm0.14
		Accuracy Disparity:	6.35\pm0.34	6.21\pm0.56	5.19\pm1.15	5.87\pm0.48	6.05\pm0.84
	Oracle-CDPFL	\mathcal{F}_{acc} :	6.78 \pm 0.46	6.53 \pm 0.61	6.43 \pm 0.72	6.89 \pm 0.07	6.62 \pm 0.28
		\mathcal{F}_{loss} :	0.10 \pm 0.02	0.09 \pm 0.01	0.10 \pm 0.02	0.07 \pm 0.01	0.08 \pm 0.01
		Min Acc:	37.41 \pm 0.39	33.29 \pm 0.42	31.91 \pm 0.16	28.25 \pm 0.76	26.13 \pm 0.98
		Accuracy Disparity:	6.01 \pm 0.28	5.29 \pm 0.22	4.95 \pm 0.53	4.63 \pm 0.65	4.75 \pm 0.25

G.4 CIFAR10 DATA SPLIT WITH CONCEPT SHIFT

Table 10: Average test accuracy of all (All), majority (Maj) and minority (Min) clients on the split of CIFAR10 dataset with concept shift (label flip). The results are averaged over 3 seeds.

algorithm		$\epsilon = 10$	$\epsilon = 5$	$\epsilon = 4$	$\epsilon = 3$	$\epsilon = 2$	
CIFAR10 (concept shift)	DPFedAvg	All:	17.98 \pm 0.31	17.19 \pm 0.20	17.10 \pm 0.14	17.06 \pm 0.19	15.28 \pm 0.27
		Maj:	19.26 \pm 0.02	18.17 \pm 0.02	17.85 \pm 0.11	17.80 \pm 0.19	15.86 \pm 0.25
		Min:	12.51 \pm 0.58	12.86 \pm 0.32	12.56 \pm 0.43	12.65 \pm 0.40	11.55 \pm 0.88
	f-CDPFL	All:	35.27 \pm 2.83	33.21 \pm 0.26	31.74 \pm 1.56	29.70 \pm 1.84	25.18\pm1.32
		Maj:	35.12 \pm 3.30	33.56 \pm 0.02	32.42\pm0.92	30.11 \pm 2.22	25.89\pm0.65
		Min:	36.17\pm0.37	32.12 \pm 0.21	27.68 \pm 5.58	27.25\pm4.70	20.95\pm5.37
	KM-CDPFL	All:	19.67 \pm 0.70	17.80 \pm 0.30	17.98 \pm 0.58	17.70 \pm 0.67	15.50 \pm 0.30
		Maj:	20.95 \pm 0.72	19.49 \pm 0.03	19.27 \pm 0.65	18.99 \pm 0.80	16.44 \pm 0.30
		Min:	12.01 \pm 0.91	10.46 \pm 0.55	10.31 \pm 0.46	9.91 \pm 0.52	9.89 \pm 1.01
	RC-DPFL	All:	37.16\pm0.26	33.41\pm0.49	32.12\pm0.44	30.25\pm1.45	23.85 \pm 1.40
		Maj:	37.45\pm0.23	33.62\pm0.02	32.39 \pm 0.45	30.98\pm0.91	24.72 \pm 0.80
		Min:	35.38 \pm 0.44	32.22\pm0.60	30.54\pm0.51	25.90 \pm 4.70	18.61 \pm 5.20
	Oracle-CDPFL	All:	37.40 \pm 0.30	33.93 \pm 0.20	32.98 \pm 0.31	32.36 \pm 0.24	26.39 \pm 0.23
		Maj:	37.64 \pm 0.27	33.81 \pm 0.05	33.20 \pm 0.28	32.67 \pm 0.28	26.59 \pm 0.12
		Min:	35.98 \pm 0.56	32.91 \pm 0.28	31.65 \pm 0.45	30.56 \pm 0.12	25.16 \pm 0.88

Table 11: Fairness evaluation in terms of different metrics on the split of CIFAR10 dataset with concept shift (label flip). The results are averaged over 3 seeds.

algorithm		$\epsilon = 10$	$\epsilon = 5$	$\epsilon = 4$	$\epsilon = 3$	$\epsilon = 2$	
CIFAR10 (concept shift)	DPFedAvg	\mathcal{F}_{acc} :	12.14 \pm 0.92	9.25 \pm 1.55	10.03 \pm 1.38	9.59 \pm 2.26	8.59 \pm 0.99
		\mathcal{F}_{loss} :	0.29 \pm 0.01	0.24 \pm 0.03	0.21 \pm 0.02	0.20 \pm 0.01	0.13 \pm 0.02
		Min Acc:	11.11 \pm 0.45	11.63 \pm 0.25	11.19 \pm 0.31	11.55 \pm 0.76	11.17 \pm 0.86
		Accuracy Disparity:	11.89 \pm 0.27	9.49 \pm 0.05	10.53 \pm 1.06	9.83 \pm 1.62	6.95 \pm 1.82
	f-CDPFL	\mathcal{F}_{acc} :	11.40 \pm 4.66	8.38 \pm 0.34	9.44 \pm 1.13	11.54 \pm 5.45	9.54 \pm 0.98
		\mathcal{F}_{loss} :	0.14 \pm 0.07	0.08 \pm 0.01	0.16 \pm 0.12	0.22 \pm 0.09	0.09\pm0.05
		Min Acc:	32.15 \pm 4.49	31.07\pm0.17	26.69 \pm 5.44	23.45 \pm 4.48	19.97\pm5.13
		Accuracy Disparity:	7.33 \pm 4.44	4.26 \pm 0.17	8.45 \pm 5.56	11.01 \pm 4.73	8.59\pm5.34
	KM-CDPFL	\mathcal{F}_{acc} :	20.99 \pm 1.67	16.66 \pm 2.41	16.41 \pm 2.92	15.95 \pm 3.23	11.15 \pm 2.63
		\mathcal{F}_{loss} :	0.59 \pm 0.07	0.46 \pm 0.01	0.47 \pm 0.02	0.45 \pm 0.03	0.24 \pm 0.01
		Min Acc:	10.61 \pm 1.30	9.43 \pm 0.61	9.25 \pm 0.39	8.89 \pm 0.66	9.07 \pm 1.103
		Accuracy Disparity:	21.35 \pm 2.44	17.85 \pm 1.47	18.21 \pm 2.03	18.71 \pm 2.63	12.75 \pm 1.81
	RC-DPFL	\mathcal{F}_{acc} :	8.09\pm1.18	8.20\pm0.20	7.88\pm0.70	8.37\pm1.08	8.03\pm0.03
		\mathcal{F}_{loss} :	0.09\pm0.01	0.07\pm0.02	0.06\pm0.01	0.16\pm0.08	0.12 \pm 0.07
		Min Acc:	34.69\pm0.11	30.97 \pm 0.82	29.45\pm0.71	24.57\pm4.99	17.61 \pm 5.09
		Accuracy Disparity:	4.45\pm0.79	4.21\pm0.08	5.06\pm0.81	9.53\pm5.49	10.07 \pm 4.82
	Oracle-CDPFL	\mathcal{F}_{acc} :	8.77 \pm 1.29	8.88 \pm 0.38	8.00 \pm 0.44	7.80 \pm 0.45	8.28 \pm 0.46
		\mathcal{F}_{loss} :	0.09 \pm 0.01	0.07 \pm 0.01	0.06 \pm 0.01	0.08 \pm 0.01	0.05 \pm 0.01
		Min Acc:	34.83 \pm 0.63	30.59 \pm 0.32	30.91 \pm 0.53	29.63 \pm 0.30	23.83 \pm 0.96
		Accuracy Disparity:	4.97 \pm 0.88	4.91 \pm 0.22	4.49 \pm 0.21	5.17 \pm 0.51	5.17 \pm 0.96

H PRIVACY GUARANTEES OF RC-DPFL

The privacy guarantee of RC-DPFL for each client i in the system comes from the fact that the client runs DP-SGD with a fixed DP noise variance $\sigma_{i,DP}^2 = c^2 \cdot z_i^2(\epsilon, \delta, b_i^1, b_i^{>1}, N_i, K, E)$ in each of its batch gradient computations. We provide a formal privacy guarantee for the algorithm to show the record-level DP privacy guarantees provided to each client i with respect to its local dataset \mathcal{D}_i .

Theorem H.1. *The set of model updates $\{\Delta\tilde{\theta}_i^e\}_{e=1}^E$, which are uploaded to the server by each client $i \in \{1, \dots, n\}$ during the training time, satisfies (ϵ, δ) -DP with respect to the client's local dataset \mathcal{D}_i , where the parameters ϵ and δ depend on the amount of DP noise $\sigma_{i,DP}^2$ used by the client.*

Proof. We remember that the sensitivity of the batch gradient in Equation (1) to every data sample is c . Therefore, based on Proposition 7 in (Mironov, 2017) each of the batch gradient computations by client i (in the first round $e = 1$ as well as the next rounds $e > 1$) is $(\alpha, \frac{\alpha c^2}{2\sigma_{i,DP}^2})$ -RDP w.r.t the local dataset \mathcal{D}_i . Therefore, if the client runs E_i^{tot} total number of gradient updates during the training time, which results in the model updates $\{\Delta\tilde{\theta}_i^e\}_{e=1}^E$ uploaded to the server, the set of model updates will be $(\alpha, \frac{E_i^{\text{tot}} \alpha c^2}{2\sigma_{i,DP}^2})$ -RDP w.r.t \mathcal{D}_i , according to Proposition 1 in (Mironov, 2017). Finally, according to Proposition 3 in the same work, this guarantee is equivalent to $(\frac{E_i^{\text{tot}} \alpha c^2}{2\sigma_{i,DP}^2} + \frac{\log(1/\delta)}{\alpha-1}, \delta)$ -DP (for any $\delta > 1$). The RDP-based guarantee is computed over a bunch of orders α and the best result among them is chosen as the privacy guarantee. Therefore, the proof is complete and the set $\{\Delta\tilde{\theta}_i^e\}_{e=1}^E$ satisfies (ϵ, δ) -DP w.r.t \mathcal{D}_i , with $\epsilon = \frac{E_i^{\text{tot}} \alpha c^2}{2\sigma_{i,DP}^2} + \frac{\log(1/\delta)}{\alpha-1}$ derived above, and $\delta > 0$. Tighter bounds for ϵ can be derived by using the numerical procedure, proposed in (Mironov et al., 2019), for accounting sampled Gaussian mechanism. \square

As a complement to the above theorem, which focuses on the privacy guarantee of the uploaded model updates, we note that the local model selection by clients in the third stage of RC-DPFL has also an effectively zero "record-level" privacy leakage. We explain about this important point as follows. The sensitivity of a client's "average train loss value" (which is not shared with the server) to the addition or removal of a sample in the client's local dataset is relatively small, as this value is resulted from an "averaging over the whole dataset". However, even if the "average train loss value" of a client changes slightly as a result of addition or removal of a sample (a record-level change), the client's selected model remains the same, especially in the third stage of RC-DPFL that one of the M existing cluster models results in a smaller loss by a large margin. Therefore, in the third stage, the sensitivity of the local model selection of a client to a record-level change is effectively zero (in contrast to the sensitivity of the average train loss, which was "small"). This is true especially for our algorithm, which transitions to a loss-based strategy after some training progress in the first E_c rounds. Our claim of "effectively zero record-level privacy leakage of local model selections" is an immediate consequence of the their "zero sensitivity to record-level changes" discussed above. Please remember that the Definition 3.1 considered in this work is concerned with "record-level" changes.

I GRADIENT ACCUMULATION

When training large models with DP-SGD, increasing the batch size results in memory exploding during training or finetuning. This might happen even when we are not using DP training. On the other hand, using a small batch size results in larger stochastic noise in batch gradients. Also, in the case of DP training, using a small batch size results in fast increment of DP noise (as explained in Lemma 4.1 in details). Therefore, if the memory budget of devices allow, we prefer to avoid using small batch sizes. But what if there is a limited memory budget? A solution for virtually increasing batch size is "gradient accumulation", which is very useful when the available physical memory of GPU is insufficient to accommodate the desired batch size. In gradient accumulation, gradients are computed for smaller batch sizes and summed over multiple batches, instead of updating model parameters after computing each batch gradient. When the accumulated gradients reach the target logical batch size, the model weights are updated with the accumulated batch gradients. The page in https://opacus.ai/api/batch_memory_manager.html shows the implementation of gradient accumulation for DP training.