

FEDGC: AN ACCURATE AND EFFICIENT FEDERATED LEARNING UNDER GRADIENT CONSTRAINT FOR HETEROGENEOUS DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

Federated Learning (FL) is an important paradigm in large-scale distributed machine learning, which enables multiple clients to jointly learn a unified global model without transmitting their local data to a central server. FL has attracted growing attentions in many real-world applications, such as multi-center cardiovascular disease diagnosis and autonomous driving. Practically, the data across clients are always heterogeneous, i.e., not independently and identically distributed (Non-IID), making the local models suffer from catastrophic forgetting of the initial (or global) model. To mitigate this forgetting issue, existing FL methods may require additional regularization terms or generate pseudo data, resulting to 1) limited accuracy; 2) long training time and slow convergence rate for real-time applications; and 3) high communication cost. In this work, an accurate and efficient *Federated Learning algorithm under Gradient Constraints* (FedGC) is proposed, which provides three advantages: i) High accuracy is achieved by the proposed *Client-Gradient-Constraint based projection method* (CGC) to alleviate the forgetting issue occurred in clients, and the proposed *Server-Gradient-Constraint based projection method* (SGC) to effectively aggregate the gradients of clients; ii) Short training time and fast convergence rate are enabled by the proposed fast *Pseudo-gradient-based mini-batch Gradient Descent* (PGD) method and SGC; iii) Low communication cost is required due to the fast convergence rate and only gradients are necessary to be transmitted between server and clients. In the experiments, four real-world image datasets with three Non-IID types are evaluated, and five popular FL methods are used for comparison. The experimental results demonstrate that our FedGC not only significantly improves the accuracy and convergence rate on Non-IID data, but also drastically decreases the training time. Compared to the state-of-art FedReg, our FedGC improves the accuracy by up to 14.28% and speeds up the local training time by 15.5 times while decreasing 23% of the communication cost.

1 INTRODUCTION

Federated Learning (FL) enables multiple participations / clients to collaboratively train a global model while keeping the training data local due to various concerns such as data privacy and real-time processing. FL has attracted growing attention in many real-world applications, such as multi-center cardiovascular disease diagnosis Linardos et al. (2022), Homomorphic Encryption-based healthcare system Zhang et al. (2022), FL-based real-time autonomous driving Zhang et al. (2021a); Nguyen et al. (2022), FL-based privacy-preserving vehicular navigation Kong et al. (2021), FL-based automatic trajectory prediction Majcherczyk et al. (2021); Wang et al. (2022). However, in practice, the data across clients are always *heterogeneous*, i.e., *not independently and identically distributed* (Non-IID) (Sattler et al., 2020; Zhang et al., 2021b), which hinders the optimization convergence and generalization performance of FL in real-word applications. At each communication round, a client firstly receives the aggregated knowledge of all clients from the server and then locally trains its model using its own data. If the data are Non-IID across clients, the local optimum of each client can be far from the others after local training **and the initial model parameters received from server will be overridden**. Hence, the clients will forget the initially received knowledge from

the server, i.e., the clients suffer from the *catastrophic forgetting* of the learned knowledge from other clients Shoham et al. (2019); Xu et al. (2022). **In other words, there is a drastic performance drop (or loss increase) of model on global data after local training (as detailed in Appendix A.10).**

Recently, several approaches have been proposed to mitigate the catastrophic forgetting in FL, e.g., Federated Curvature (FedCurv) (Shoham et al., 2019) and FedReg Xu et al. (2022). FedCurv utilizes the continual learning method Elastic Weight Consolidation (EWC) (Kirkpatrick et al., 2017) to penalize the clients for changing the most informative parameters. The Fisher information matrix is used in EWC to determine which parameters are informative. However, EWC is not effective for mitigating the catastrophic forgetting Xu et al. (2022) in FL, and FedCurv needs to transmit the Fisher matrix between the server and clients besides model parameters. That significantly increases the communication cost (2.5 times than the baseline FedAvg Xu et al. (2022)). In addition, the calculation of Fisher matrix drastically increases the local training time. FedReg (Xu et al., 2022) is the most recently proposed FL method inspired by the continual learning method Gradient Episodic Memory (GEM) Lopez-Paz & Ranzato (2017). GEM alleviates the catastrophic forgetting by avoiding the increase of loss at previous tasks. However, it requires an episodic memory to contain the representative samples from all previous tasks, which hinders it from being suitable for FL due to data privacy concerns Xu et al. (2022). To resolve this, each client in FedReg firstly generates pseudo data by encoding the knowledge of previous training data learned by the global model, and then regularizes its model parameters by avoiding the increase of loss on pseudo data after local training. Although it uses generated pseudo data to protect data privacy and alleviate the forgetting issue in FL, the data generation will increase a lot of computational and storage costs for clients, especially when clients have large-scale data. In addition, the generation of pseudo data and parameter regularization also significantly increase the local training time. Therefore, these methods are not friendly enough to many real-time applications that concern communication & computational costs.

In this work, we propose an accurate and efficient *Federated Learning algorithm under Gradient Constraints* (FedGC) to improve the performance of FL on Non-IID data and reduce the local training time. **At client**, a fast *Pseudo-gradient-based mini-batch Gradient Descent* (PGD) algorithm is proposed to reduce the local training time while accelerating the convergence rates of FL. The pseudo gradient of a local model is obtained by calculating its gradients over few mini-batches data using gradient descent algorithm. In addition, to mitigate catastrophic forgetting, we propose an effective *Client-Gradient-Constraint based projection method* (CGC). Different from GEM requiring memorized data from other clients and FedReg generating pseudo data at clients, our CGC only utilizes the server gradient (i.e., the aggregated gradient from all clients) to restrict the projected gradient to satisfy the constraint: the angle between these two gradients is less than 90° , in order to enable the local model retains more knowledge received from server. Meanwhile, the projected gradient is also forced to be as close as possible to the pseudo gradient, that enables the local model to learn new knowledge from local data. **At server**, we propose a *Server-Gradient-Constraint based projection method* (SGC) to achieve an optimal server gradient which involves the information of clients participated in aggregation while accelerating the convergence rate by restricting the angles between the server gradient and gradients of participating clients to be less than 90° . Moreover, our FedGC only transmits the gradients between the server and clients. In other words, our FedGC greatly saves communication costs. The contributions are summarized as follows,

- i) High accuracy of our FedGC on Non-IID data is achieved by the proposed CGC to mitigate the catastrophic forgetting occurred in clients and the proposed SGC to effectively aggregate the gradients of clients;
- ii) Short training time and fast convergence rate in our FedGC are enabled by the proposed fast PGD method and SGC;
- iii) Low communication cost is required in our FedGC due to the fast convergence rate and only gradients to be transmitted between server and clients;
- iv) Extensive experimental results illustrate that our FedGC not only improves the performance of FL on Non-IID data with a fast convergence rate but also significantly reduces local training time.

2 RELATED WORKS

Federated learning is an important paradigm in large-scale distributed machine learning. It enables multiple clients to jointly learn a unified global model without transmitting their local data to a central server (McMahan et al., 2017; Bhagoji et al., 2019; Yang et al., 2019). FedAvg (McMahan et al., 2017) is the most popular FL algorithm. In FedAvg, clients first locally train models on local data, and then their model updates (e.g., parameters) are transmitted over the network to a central server, where the updates are aggregated. However, data in many real-world applications are always Non-IID, which degrades the performance of FedAvg and slows down the convergence rate (Li et al. (2020; 2018); Xu et al. (2022)). Many approaches have been proposed to improve the performance and accelerate the convergence rate of FedAvg. Except FedCurv and FedReg, two more popular FL methods resolving the Non-IID are introduced here, including FedProx (Li et al., 2018) and Stochastic Controlled Averaging algorithm (SCAFFOLD) (Karimireddy et al. (2020)). Similar to FedCurv, FedProx (Li et al., 2018) tackles the heterogeneity in FL by adding a regularization term (i.e., the Proximal term) in local objective function. The proximal term represents the l_2 distance between the local model parameters and the initial (global) model parameters. During local training, FedProx minimizes the local loss while restricting the local updates to be close to the initial (global) model. FedProx takes the same communication costs as FedAvg because it does not transmit additional information besides model parameters. However, compared to FedCurv, it lacks of flexibility of model parameters and its stiffness comes at the expense of accuracy. SCAFFOLD (Karimireddy et al., 2020) introduces the control variates c and c_i to correct the client-drift on Non-IID data in local training. The control variates c and c_i aim to guide the local model to update based on the average gradient of participating clients in the previous communication round. However, the variate c_i has the same dimension with the gradient and is transmitted between clients and server, which doubles the communication costs compared with FedAvg. **Moreover, the average gradient in the previous communication round may not satisfy its assumptions that $c_j \approx g_j(y_i)$ and $c \approx \frac{1}{N} \sum_j g_j(y_i)$, especially when deep learning models on image datasets perform at clients (Li et al. (2021)).**

Notably, all the above FL methods rely on stochastic gradient descent (SGD) to train models at clients by performing multiple epochs on full local data. That significantly increases the local training time of clients. In addition, the clients owning small training data need to wait a long time before the clients with large-scale data complete local training. That is not computationally efficient in practice and increases the latency between clients. **Minibatch SGD (Woodworth et al., 2020) is recently proposed to perform local training of clients on several mini-batch data at the same model, and then it calculates the mean gradient by averaging the gradients.**

3 METHOD

Given K clients, each client k has a local dataset \mathcal{D}_k . $\mathbb{D}_k = \{\mathbf{x}_k, \mathbf{y}_k\} \subseteq \mathcal{D}_k$ represents a few mini-batches data including n_k samples. Let T be the number of communication rounds, B be the number of mini-batches for local training. At each communication round, a subset of clients $\mathcal{K} \subseteq [K]$ are sampled uniformly like FedAvg. In the t -th communication round, θ^t represents the parameters of the global model at server and θ_k^t is the parameters of the local model at client k . g_k^t represents the gradient of client k on its local data after local training, g^t is the server gradient that is obtained by aggregating local gradients of clients. Notably, the gradient mentioned in this work is the negative gradient for the convenience of calculation.

3.1 LOCAL TRAINING WITH PGD AND CGC

At the t -th communication round, clients would firstly receive g^{t-1} of the previous round, and then synchronize the parameters of local models with g^{t-1} to ensure clients have the same initially parameters. For client $k \in \mathcal{K}$ with its mini-batches data \mathbb{D}_k , the objective function is given by,

$$\min_{\theta_k^t} L(f(\mathbf{x}_k; \theta_k^t), \mathbf{y}_k) \quad \text{s.t.} \quad \left\langle \frac{\partial L(f(\mathbf{x}_k; \theta_k^t), \mathbf{y}_k)}{\partial \theta_k^t}, g^{t-1} \right\rangle \geq 0 \quad (1)$$

where $\langle \cdot \rangle$ represents the inner product operation. $f(\theta_k^t, \mathbf{x}_k)$ represents the prediction of local model at client k with parameters θ_k^t and input \mathbf{x}_k , $L(f(\theta_k^t, \mathbf{x}_k), \mathbf{y}_k)$ is the loss function of client k on its local data \mathbb{D}_k and

$$L(f(\mathbf{x}_k; \boldsymbol{\theta}_k^t), \mathbf{y}_k) = \frac{1}{|\mathbb{D}_k|} \sum_{(\mathbf{x}_{k,i}, \mathbf{y}_{k,i}) \in \mathbb{D}_k} L(f(\mathbf{x}_{k,i}; \boldsymbol{\theta}_k^t), \mathbf{y}_{k,i}) \quad (2)$$

The constraint in problem (1) indicates that the angle between the current gradient $\frac{\partial L(f(\mathbf{x}_k; \boldsymbol{\theta}_k^t), \mathbf{y}_k)}{\partial \boldsymbol{\theta}_k^t}$ and the server gradient \mathbf{g}^{t-1} is less than 90° . In this way, during local training, the update direction of local model will be not only learned from local data but also restricted by the server gradient. The server gradient obtained by our SGC involves the update direction of clients participated in aggregation (as detailed in Section 3.2). Hence, through solving problem (1), the local model will convergence to the global optimum.

However, (1) is an optimization problem with inequality constraints, which cannot be directly solved by SGD. Hence, we divided the solution into two steps, as shown in Figure 1.

Step 1: Pseudo-gradient-based mini-batch Gradient Descent (PGD)

At t -th communication round, for client $k \in \mathcal{K}$ with its local data \mathbb{D}_k , the loss function Eq. (3) is first minimized,

$$\min_{\boldsymbol{\theta}_{k,B}^t} L(f(\mathbf{x}_k; \boldsymbol{\theta}_{k,B}^t), \mathbf{y}_k) \quad (3)$$

where $\boldsymbol{\theta}_{k,B}^t$ is the model parameters after B mini-batches.

Mini-batch update: Instead of performing multiple epochs on full local data in each communication round like in FedAvg, we prefer to use few mini-batches to train local model to save local training time and avoid the latency between clients. **Most importantly, when data are Non-IID, too many iterations in local training may cause local models to be biased towards their local data and enlarge the difference between the local model and the global model. This is not conducive to the convergence of global model. The experiments of FedAvg using mini-batch update in Appendix A.7 can verify this statement.**

Pseudo-gradient: During this local training, $\mathbf{g}_{k,1}^t, \mathbf{g}_{k,2}^t, \dots, \mathbf{g}_{k,B}^t$ in Figure 1 (b) are obtained. **Different from the mean gradient $\bar{\mathbf{g}}_k^t (= \frac{1}{B} \sum_{b=1}^B \mathbf{g}_{k,b}^t)$ used in Minibatch SGD that averages multiple gradients for several mini-batches data at the same point (e.g., $\boldsymbol{\theta}_{i,0}^t$), we prefer to calculate a pseudo gradient $\tilde{\mathbf{g}}_k^t$ to represent the final gradient after this local training, as shown below.**

$$\tilde{\mathbf{g}}_k^t = \frac{(\boldsymbol{\theta}_{k,B}^t - \boldsymbol{\theta}_{k,0}^t)}{\eta} \quad (4)$$

where $\boldsymbol{\theta}_{k,0}^t$ is the initialized model parameters and η is the learning rate.

Compared to the mean gradient $\bar{\mathbf{g}}_k^t$ that has similar modulus with each gradient $\mathbf{g}_{k,l}^t, 1 \leq l \leq B$, the pseudo gradient $\tilde{\mathbf{g}}_k^t$ would have larger modulus than $\mathbf{g}_{k,l}^t$, as shown in Figure 1(b). In this way, $\tilde{\mathbf{g}}_k^t$ will promote a large update for local model even with only few mini-batches in local training. This will accelerate the convergence rate of FL.

Step 2: Client-Gradient-Constraint based projection method (CGC)

Since the data cross clients are Non-IID, the angle between the pseudo gradient $\tilde{\mathbf{g}}_k^t$ and server gradient \mathbf{g}^{t-1} may be more than 90° , as shown in Figure 1(b). That means the update direction of local model deviates from the global optimum, i.e., the catastrophic forgetting occurs at clients. To mitigate this forgetting, the *Client-Gradient-Constraint based projection method* (CGC) is proposed in this subsection.

At t -th communication round, for client $k \in \mathcal{K}$ with its local data \mathbb{D}_k , the optimization problem of our CGC is given by (5). In problem (5), the projected gradient \mathbf{g}_k^t should be as close as possible

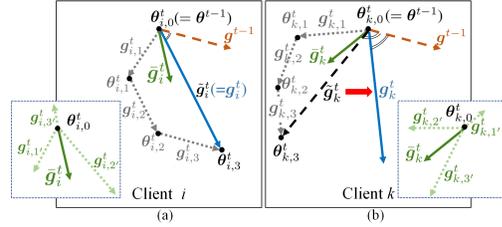


Figure 1: Local training with PGD and CGC performed at two clients of the t -th round. Clients receive the server gradient \mathbf{g}^{t-1} before local training. At step 1, local models perform PGD on three mini-batches, and compute the pseudo gradients (i.e., $\tilde{\mathbf{g}}_i^t$ and $\tilde{\mathbf{g}}_k^t$) by (4). At step 2, clients obtain the projected gradients by performing CGC on the pseudo gradients. Since the angle between $\tilde{\mathbf{g}}_i^t$ and \mathbf{g}^{t-1} is less than 90° , the projected gradient of client i is $\tilde{\mathbf{g}}_i^t$ itself through CGC. In contrast, the angle between $\tilde{\mathbf{g}}_k^t$ and \mathbf{g}^{t-1} is more than 90° , the projected gradient \mathbf{g}_k^t is obtained by performing CGC (red arrow) on $\tilde{\mathbf{g}}_k^t$ to satisfy constraints of (5). $\bar{\mathbf{g}}_i^t$ and $\bar{\mathbf{g}}_k^t$ (green arrows) are the mean gradients used in Mini-batch SGD. The projected gradients usually have larger modulus than the mean gradients.

to the pseudo gradient $\tilde{\mathbf{g}}_k^t$ (in squared L2 norm) while being at an acute angle to the server gradient \mathbf{g}^{t-1} .

$$\min_{\mathbf{g}_k^t} \frac{1}{2} \|\tilde{\mathbf{g}}_k^t - \mathbf{g}_k^t\|^2 \quad \text{s.t.} \quad \langle \mathbf{g}_k^t, \mathbf{g}^{t-1} \rangle - C \geq 0 \quad (5)$$

where $C = 1e - 3$ is a small positive constant to avoid \mathbf{g}_k^t and \mathbf{g}^{t-1} being orthogonal (i.e., to ensure $\langle \mathbf{g}_k^t, \mathbf{g}^{t-1} \rangle > 0$) after projection.

Through solving problem (5), the projected gradient \mathbf{g}_k^t will enable the model retain more knowledge received from server while learning new knowledge from local data, i.e., the model can balance preserving the knowledge received from the server and being adaptive to local data.

By simplifying problem (5), we obtain the primal of CGC Quadratic Program (QP) with inequality constraints:

$$\min_{\mathbf{u}} \frac{1}{2} \mathbf{u}^T \mathbf{u} - \mathbf{h}^T \mathbf{u} + \frac{1}{2} \mathbf{h}^T \mathbf{h} \quad \text{s.t.} \quad C - \mathbf{z}^T \mathbf{u} \leq 0 \quad (6)$$

where $\frac{1}{2} \mathbf{h}^T \mathbf{h}$ is a constant term and can be discarded, $\mathbf{u} = \mathbf{g}_k^t \in \mathbb{R}^p$, $\mathbf{h} = \tilde{\mathbf{g}}_k^t \in \mathbb{R}^p$ and $\mathbf{z} = \mathbf{g}^{t-1} \in \mathbb{R}^p$, and p indicates the number of parameters of local model.

Problem (6) is a QP on p variables and could be measured in millions. Hence, to solve problem (6) efficiently, we convert the primal problem into dual problem, and obtain the dual of the CGC QP:

$$\min_v \frac{1}{2} v^2 \mathbf{z}^T \mathbf{z} + v (\mathbf{h}^T \mathbf{z} - C) \quad \text{s.t.} \quad v \geq 0 \quad (7)$$

where $v \in \mathbb{R}$ is a Lagrange multiplier, and problem (7) is a QP on 1 $\ll p$ variable. To solve problem (7), the python library *quadprog*¹ is used and then the optimum v^* is obtained. The details of calculating dual of CGC QP are provided in Appendix A.1.

Finally, the optimal solution to problem (6) is calculated by $\mathbf{u}^* = \mathbf{h} + v^* \mathbf{z}$, i.e., $\mathbf{g}_k^t = \tilde{\mathbf{g}}_k^t + v^* \mathbf{g}^{t-1}$ after our CGC. In addition, since the whole concatenated gradient has very large dimension, we iteratively perform the gradient projection layer by layer (i.e., layer-wise manner) to reduce the memory overhead. Theoretical analysis of gradient projection is detailed in Appendix A.3

3.2 SERVER AGGREGATION WITH SGC

At t -th communication round, the local gradients (i.e., \mathbf{g}_k^t) of clients are then send to server for aggregation after local training. At server, the aggregated gradient \mathbf{g}^t (i.e., server gradient) is then send back to clients and the parameters $\boldsymbol{\theta}^t$ of global model is calculated by $\boldsymbol{\theta}^t = \boldsymbol{\theta}^{t-1} + \eta \mathbf{g}^t$. For aggregation, most FL methods simply use the weighted average of local gradients at server. When the data across clients are Non-IID, the weighted-average may be only effective for few clients. That is because the server gradient may point to the opposite directions of some local gradients (i.e., the angle between them are more than 90°), that slow down the convergence rate of FL.

To effectively aggregate the local gradients at server and accelerate the convergence rate, we proposed a *Server-Gradient-Constraint based projection method* (SGC). Through our SGC, the projected gradient can point to the positive directions of most local gradients. The optimization problem of SGC is given by,

$$\min_{\mathbf{g}^t} \frac{1}{2} \|\mathbf{g}^t - \bar{\mathbf{g}}^t\|^2 \quad \text{s.t.} \quad \langle \mathbf{g}^t, \mathbf{g}_k^t \rangle - C \geq 0, \forall k \in \mathcal{K} \quad (8)$$

where $\bar{\mathbf{g}}^t = \sum_{k \in \mathcal{K}} \frac{n_k}{N} \mathbf{g}_k^t$ is the weighted average of the local gradients, and $N = \sum_{k \in \mathcal{K}} n_k$. \mathbf{g}_k^t is the local gradient of client k . The constraint in problem (8) is to restrict the angle between the projected gradient \mathbf{g}^t and local gradients at participating clients are less than 90° .

Similarly, we obtain the primal of SGC QP with inequality constrains by simplifying Eq. (8)

$$\min_{\mathbf{z}} \frac{1}{2} \mathbf{z}^T \mathbf{z} - \mathbf{g}^T \mathbf{z} + \frac{1}{2} \mathbf{g}^T \mathbf{g} \quad \text{s.t.} \quad C - \mathbf{G} \mathbf{z} \leq 0 \quad (9)$$

where $\mathbf{z} = \mathbf{g}^t \in \mathbb{R}^p$, $\mathbf{g} = \bar{\mathbf{g}}^t \in \mathbb{R}^p$, and $\mathbf{G} = (\dots, \mathbf{g}_k^t, \dots)^T \in \mathbb{R}^{|\mathcal{K}| \times p}$, $k \in \mathcal{K}$. $\frac{1}{2} \mathbf{g}^T \mathbf{g}$ is a constant term and can be ignored. Problem (9) is a QP on p variables. To solve problem (8) efficiently, we also convert the primal problem (9) into dual problem and obtain the dual of the SGC QP:

$$\min_{\boldsymbol{\lambda}} (\mathbf{G} \mathbf{g} - C)^T \boldsymbol{\lambda} + \frac{1}{2} \boldsymbol{\lambda}^T \mathbf{G} \mathbf{G}^T \boldsymbol{\lambda} \quad \text{s.t.} \quad \boldsymbol{\lambda} \geq 0 \quad (10)$$

where $\boldsymbol{\lambda} \in \mathbb{R}^{|\mathcal{K}|}$ is the Lagrange multiplier, and the problem (10) is a QP on $|\mathcal{K}| \ll p$ variables. The solution $\boldsymbol{\lambda}^*$ of problem (10) is also obtained by the *quadprog* library. The details of calculating dual of SGC QP are provided in Appendix A.2.

¹<https://github.com/quadprog/quadprog>

Finally, the optimal solution to problem (9) is calculated by $z^* = \mathbf{g} + \mathbf{G}^T \lambda^*$, i.e., $\mathbf{g}^t = \bar{\mathbf{g}}^t + \mathbf{G}^T \lambda^*$ after our SGC. When the problem (9) is unsolvable, we simply use $\bar{\mathbf{g}}^t$ to be \mathbf{g}^t .

The pseudo codes of our method are provided in Algorithm 1, and the convergence analysis is detailed in Appendix A.4.

4 EXPERIMENTS

We conduct extensive experiments to compare our FedGC with several popular approaches, including FedAvg, FedProx, FedCurv, SCAFFOLD and FedReg, on real image datasets. The data preparation and experimental details are described below. The performances are evaluated in three aspects: 1) Overall test accuracy; 2) Convergence rate and training time; 3) Communication costs.

4.1 DATASETS

The experiments are conducted on three real image datasets, including Handwritten-Digits, CIFAR-10 (Krizhevsky, 2009) and CIFAR-100 (Krizhevsky (2009)). The data preparation of each dataset is described below. More datasets details are provided in Appendix A.6

Handwritten-Digits It contains four common handwritten-digits datasets, including MNIST (LeCun et al. (1998)), MNIST-M (Zhao et al. (2022)), USPS (Hull, 1994) and SVHN (Netzer et al. (2011)). Each dataset contains 10 classes. To input these images into deep models sharing the same network architecture, all images in these four datasets are pre-processed by reshaping the size to (32, 32, 3), including the training sets and test sets. The training data are split into 4 clients (named HWDigits-4) and 40 clients (named HWDigits-40) respectively. In HWDigits-4, each client owns one handwritten-digits dataset. In HWDigits-40, each client has images belonging to only one class in a dataset under the one-class setting. HWDigits-4 suffers from the attribute skew of Non-IID issue, in which the attributes (i.e., data features) across clients are different. In HWDigits-40, the data across clients may have different labels and attributes, and hence it suffers from both attribute and label skew.

CIFAR-10 and CIFAR-100 Under the one-class setting, the training sets in CIFAR-10 and CIFAR-100 are split into 10 (named CIFAR10-10) and 100 (named CIFAR100-100) clients, respectively, i.e., each client owns only samples of one class. All clients share the test sets in original CIFAR-10 and CIFAR-100 respectively. CIFAR10-10 and CIFAR100-100 both suffer from the label skew under the one-class setting.

4.2 EXPERIMENTAL SETTING

We implement all methods in PyTorch-1.9.0 (Paszke et al., 2019) on a Ubuntu server with two Intel Xen CPUs and 8 NVIDIA RTX 3090 GPUs. SGD with weight decay 0 and momentum 0 is used to train local models at clients. For HWDigits-4 and HWDigits-40, a CNN similar to (McMahan et al., 2017) is adopted in experiments. It contains two 5×5 convolution layers followed by 2×2 Max-Pooling and two fully connected layers (the first layer with 512 units and the second layer with 100 units) with ReLU activation. The communication round $T = 600$ and the number of local epochs E in other compared methods is set to 1. The number of local iterations B (i.e., the number of mini-

Algorithm 1: FedGC

Input: K, T, B , datasets $\mathcal{D} = \cup_{k \in [K]} \mathcal{D}_k, \eta$.
Output: the parameters θ^T of the global model.

- 1 Initialize: server θ^0 ; clients $\theta_k^0 \leftarrow \theta^0$, for $k \in [K]$;
- 2 **for** $t = 1$ to T **do**
- 3 **for** $k \in [K]$ *in parallel* **do**
- 4 # Synchronize parameters
- 5 **if** $t > 1$ **then**
- 6 $\mathbf{g}_k^t \leftarrow \mathbf{g}_k^{t-1}$;
- 7 $\theta_k^t \leftarrow \theta_k^{t-1} + \eta \mathbf{g}_k^t$;
- 8 **end**
- 9 Randomly sample clients $\mathcal{K} \subseteq [K]$;
- 10 **for** $k \in \mathcal{K}$ *in parallel* **do**
- 11 **for** $b = 1$ to B **do**
- 12 Randomly sample a mini-batch data
- 13 $\mathcal{B}_{k,b} \subset \mathcal{D}_k$;
- 14 $\theta_{k,b}^t \leftarrow \text{SGD}(\theta_{k,b-1}^t, \eta)$ on $\mathcal{B}_{k,b}$;
- 15 **end**
- 16 $\tilde{\mathbf{g}}_k^t = \frac{\theta_{k,B}^t - \theta_{k,0}^t}{\eta}$; # Pseudo
- 17 $\mathbf{g}_k^t \leftarrow \text{gradient}$
- 18 $\mathbf{g}_k^t \leftarrow \text{CGC}(\tilde{\mathbf{g}}_k^t, \mathbf{g}_k^{t-1})$ in (6);
- 19 $\theta_k^t \leftarrow \theta_{k,0}^t$;
- 20 **end**
- 21 $\bar{\mathbf{g}}^t \leftarrow \sum_{k \in \mathcal{K}} \frac{n_k}{N} \mathbf{g}_k^t$; $\mathbf{G} \leftarrow (\dots, \mathbf{g}_k^t, \dots)^T$,
- 22 $\mathbf{g}^t \leftarrow \text{SGC}(\bar{\mathbf{g}}^t, \mathbf{G})$ in (9);
- 23 $\theta^t \leftarrow \theta^{t-1} + \eta \mathbf{g}^t$;

^aIf $t = 1$, $\theta_{k,0}^t = \theta_k^{t-1}$, otherwise $\theta_{k,0}^t = \theta_k^t$

batches for local training) in FedGC is 50. The client fraction is set to 1.0 and the batch size is 100. The learning rates η of all methods in experiments are tuned in the range of $\{0.01, 0.05, 0.1, 1.0\}$, and the optimal η of our FedGC is 0.1 and the ones for all other compared methods are 0.01. The learning rate η in FedGC is set to 0.1 and the learning rates in other methods are 0.01. For CIFAR10-10 and CIFAR100-100, a ResNet-9 network with the Fixup initialization (Xu et al., 2022) is trained from scratch. The communication rounds T is set to 500 for CIFAR10-10 and 1000 for CIFAR100-100. $E = 1$ and $B = 1$. The optimal η for our FedGC is 0.1 and the ones for all other compared methods are 0.05.

4.3 OVERALL TEST ACCURACY

Table 1: Comparison results (%) on four datasets (\uparrow). ‘w/o’ and ‘w’ denote the server aggregation in FedGC without and with SGC, respectively.

Type of Non-IID	Attribute skew	Label skew		Attribute&Label skew
Method	HWDigits-4	CIFAR10-10	CIFAR100-100	HWDigits-40
FedAvg	93.07	55.34	40.59	62.11
FedProx	93.69	55.13	32.74	62.01
FedCurv	91.78	55.61	9.48 ¹	62.55
SCAFFOLD	92.29	34.40	2.62 ¹	45.32
FedReg	92.32	57.70	53.02	63.05
Our FedGC(w/o)	94.15	59.64	55.73	74.63
Our FedGC(w)	95.27	60.92	56.56	77.33
#samples per client	7,291~73,257	5,000	500	542~13,861

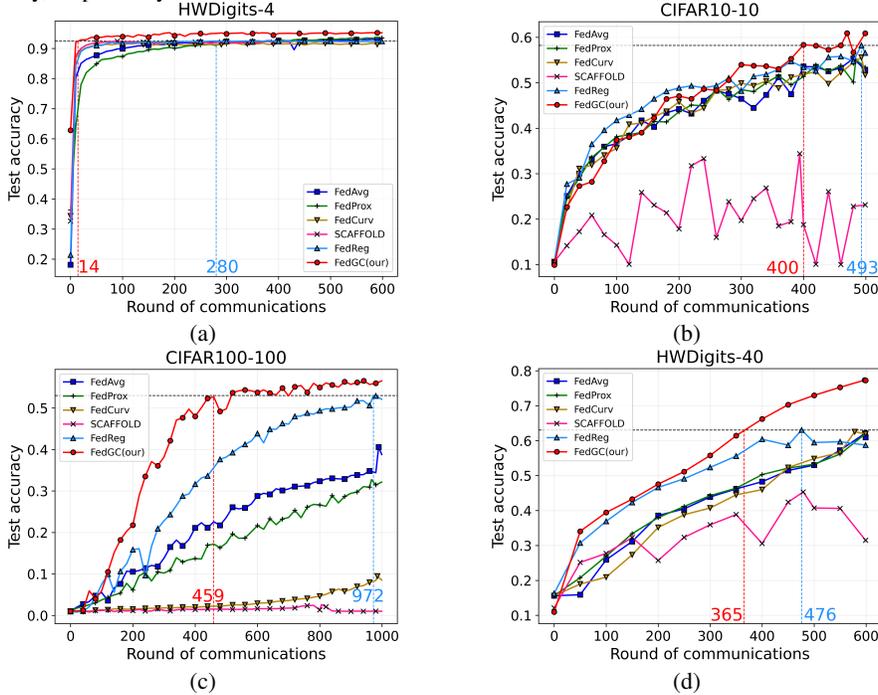
¹ It represents model fails to convergence after 1000 communication rounds.

Table 1 shows the overall test accuracies of all compared methods on HWDigits-4, HWDigits-40, CIFAR10-10 and CIFAR100-100. In Table 1, our FedGC (denotes FedGC(w) if not specified) achieves the highest test accuracies than all other compared methods. Specifically, for HWDigits-40, our FedGC outperforms the SOTA FedReg by 14.28% ($=77.33-63.05$) of accuracy, which clearly shows the superiority of our FedGC in tackling the Non-IID data. FedCurv is designed based on the assumption that the deep neural networks are over-parameterized enough, so that it has a good probability of finding an optimal solution to task B in the neighborhood of previously learned task solution. However, for CIFAR100-100, there are 100 clients and each client contains only samples of one class. Hence, the assumption of over-parameterization cannot be satisfied and it only gets 9.48% of accuracy. FedProx has strong constraints on model parameters, so that it lacks of flexibility, which hinders the learning of new knowledge from the local data. Hence, on CIFAR100-100, it also gets worse performance than the baseline FedAvg. **SCAFFOLD gets the worst performances on most of the compared datasets, because the average gradient in the previous communication round may not satisfy its assumptions, especially when deep learning models on image datasets perform at clients Li et al. (2021).** FedReg is the most recently proposed method, which alleviates the catastrophic forgetting in FL by generating pseudo data. The pseudo data is generated to guarantee that the loss of local model on them is less than that of initial (global) model on them. In this way, the catastrophic forgetting issues can be alleviated. However, when the number of samples between clients varies significantly (e.g., HWDigits-4 and HWDigits-40), the FedReg will be biased to the majority clients (i.e., the clients with many samples). That is because the global model is obtained by weighted averaging the model parameters in FedReg and the global model will be biased to majority clients. In local training, the client with few samples (i.e., minority clients) will forget the learned knowledge from local data after regularizing the model parameters with pseudo data in FedReg. In other words, the FedReg has further magnified the bias of global model to majority clients.

Benefited from the proposed CGC, our FedGC mitigates the catastrophic forgetting of FL by constraining the local gradient at an acute angle to the server gradient and simultaneously minimizing the loss of local model on its local data. Moreover, benefited from the proposed SGC, the projected server gradient in our FedGC can effectively aggregate the knowledge from clients. The gradient projected by both CGC and SGC can reduce the bias toward majority clients. As shown in Table 1, FedGC(w) further improves the performances of FedGC(w/o). Particularly, the SGC can further accelerate the convergence rate of FL (as detailed in Section 4.4). Hence, FedGC(w) achieves the highest accuracy on all compared datasets, and the improvements are up to 14.28% (HWDigits-40) compared to the state-of-the-art (SOTA) FedReg. In addition, our FedGC neither transmits extra

data across the server and clients, nor needs extra storage costs to keep the generated data. The computational costs of CGC and SGC are also extremely small. Therefore, FedGC is very friendly for edge devices with limited computational resources.

Figure 2: Test accuracy vs. communication rounds on four datasets. The black dotted line denotes the highest accuracy of FedReg. The red and blue dotted lines represent the round numbers of FedGC and FedReg reaching this accuracy, respectively.



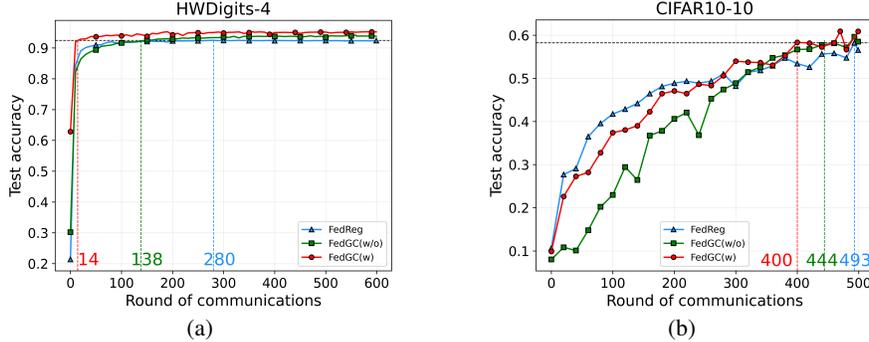
4.4 CONVERGENCE RATE AND TRAINING TIME

Figure 2 illustrates that our FedGC can quickly achieve higher accuracies than other methods on all compared datasets until the end of communication. It indicates that our method can improve both the performance and the convergence rate of FL on Non-IID data. The high convergence rate is mainly benefited from our proposed PGD and SGC. FedReg is also proposed for accelerating the convergence rate of FL by alleviating the forgetting issue, but it is more biased to majority clients, that may slow down its convergence rate on HWDigits-4 and HWDigits-40. For HWDigits-40, FedReg achieves its highest accuracy (63.05%) at round 476 as shown in Figure 2(d) while our FedGC reaches this accuracy at round 365. Meanwhile, our FedGC improves the accuracy by 14.28% while decreasing 23% ($\approx 1 - \frac{365 \times |\theta|}{476 \times |\theta|}$) of communication cost² compared to FedReg. For CIFAR100-100, our FedGC also quickly reaches 53.02% of accuracy (i.e., the accuracy of FedReg), which empirically verifies that the projected gradients in both server and clients are helpful for FL. Among compared methods, **SCAFFOLD is very unstable during training (as shown in Figures 2 (b) and (d))**. In addition, Figure 3 illustrates that our SGC can further improve the convergence rate of FL.

Training time is another effective way to measure the practicality and efficiency of FL methods. The average training time of each client per communication round are summarized in Table 2, where FedAvg is considered as the baseline. Due to the calculation of proximal term in objective function during local training, FedProx increases training time by about 13s³ at each dataset per round. Similarly, SCAFFOLD also requires more time to calculate and transmit the control variate between clients and server. FedCurv dramatically increases the training time due to the time consuming calculation of Fisher information matrix in every mini-batch. FedReg needs to generate pseudo data during local training, so that it also significantly increases the training time. On the contrary, in our FedGC, the gradient projection methods in both clients and server are very efficient due to few

²The calculation is provided in Section 4.5 with FedReg as baseline. $|\theta|$ denotes the amount of parameters.
³ $\approx \frac{(67.99+65.15+64.65+82.41)-(58.63+44+44.74+78.93)}{4}$

Figure 3: Test accuracy vs. communication rounds on HWDigits-4 and CIFAR10-10. The black dotted line denotes the highest accuracy of FedReg. The blue, green and red dotted lines represent the round numbers of FedReg, FedGC(w/o) and FedGC(w) reaching this accuracy, respectively.



variables in the dual problems. Moreover, all compared methods train local model by performing multiple epochs on full local data, while our FedGC trains local model on few mini-batches in each communication round. In this way, it will greatly decrease the training time per round, especially for the clients with many samples (e.g., HWDigits-40, HWDigits-40 and CIFAR100-100). Particularly, our FedGC speeds up the training time by 15.5 ($\approx \frac{274.72}{17.77}$) times compared to SOTA FedReg on HWDigits-40.

Table 2: The training time per communication round and the communication costs (Comm. cost) (\downarrow).

Method	HWDigits-4		CIFAR10-10		CIFAR100-100		HWDigits-40	
	Time(s)	Comm. cost	Time(s)	Comm. cost	Time(s)	Comm. cost	Time(s)	Comm. cost
FedAvg	58.63	$557 \times \theta $	44.00	$480 \times \theta $	44.74	$989 \times \theta $	78.93	$598 \times \theta $
FedProx	67.99	$439 \times \theta $	65.15	N.A.	64.65	N.A.	82.41	N.A.
FedCurv	970.42	N.A. ¹	1055.01	$490 \times 2.5 \theta $	1061.29	N.A.	933.11	$578 \times 2.5 \theta $
SCAFFOLD	81.21	N.A.	63.99	N.A.	72.54	N.A.	115.84	N.A.
FedReg	260.97	N.A.	257.69	$431 \times \theta $	260.72	$531 \times \theta $	274.72	$476 \times \theta $
FedGC (our)	6.33	$19 \times \theta$	18.82	$327 \times \theta$	39.71	$287 \times \theta$	17.77	$357 \times \theta$

¹ N.A. represents the method has not achieved the target accuracy within maximum communication rounds.

4.5 COMMUNICATION COSTS

The comparison of communication costs is shown in Table 2. We calculate the communication cost by the number of communication round when the method achieves the target accuracy (e.g., the highest accuracy of baseline FedAvg) multiplied by the amount of transmitted data (e.g., the amount of parameters $|\theta|$). $2.5|\theta|$ represents that FedCurv needs to transmit the Fisher matrix between the server and clients besides model parameters. In Table 2, our FedGC requires the lowest communication costs on four datasets when achieving the target accuracy. Specifically, on HWDigits-4, our FedGC can decrease 96.6% ($\approx 1 - \frac{19 \times |\theta|}{557 \times |\theta|}$) of the communication cost compared to FedAvg due to the fewest number of communication rounds. Compared to SOTA FedReg, the communication cost can be decreased by 45% ($\approx 1 - \frac{287 \times |\theta|}{531 \times |\theta|}$) on CIFAR100-100.

5 CONCLUSIONS

This paper proposes an accurate and efficient Federated Learning algorithm under Gradient Constraints (FedGC) to address the challenging FL scenarios on Non-IID data: Low accuracy, and long local training time. Our FedGC includes several novel modules to tackle the issues of Non-IID data: 1) CGC (alleviating the forgetting issue for higher accuracy) and PGD (for significantly shorter training time and fast convergence rate) at client; 2) SGC (for effective aggregation of local gradients) at server. Moreover, FedGC enables low communication cost due to fast convergence rate and only transmissions of gradients between server and clients. Extensive experiments are conducted on four datasets comparing several SOTA FL methods. The experimental results show that FedGC can significantly improve the performance and convergence rate of FL on Non-IID data with low communication and storage costs. The local training time can be drastically reduced. Our FedGC improves the accuracy by up to 14.28% and speeds up the local training time by 15.5 times compared to the SOTA FedReg, while decreasing 23% of the communication cost. In a nutshell, FedGC has great potential for many real-world applications that concerns performance, real-time, communication & computational costs, and privacy preservation.

REFERENCES

- Arjun Nitin Bhagoji, Supriyo Chakraborty, Prateek Mittal, and Seraphin Calo. Analyzing federated learning through an adversarial lens. In *International Conference on Machine Learning*, pp. 634–643. PMLR, 2019.
- J.J. Hull. A database for handwritten text recognition research. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(5):550–554, 1994. doi: 10.1109/34.291440.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pp. 5132–5143. PMLR, 2020.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- Qinglei Kong, Feng Yin, Rongxing Lu, Beibei Li, Xiaohong Wang, Shuguang Cui, and Ping Zhang. Privacy-preserving aggregation for federated learning-based navigation in vehicular fog. *IEEE Transactions on Industrial Informatics*, 17(12):8453–8463, 2021.
- A Krizhevsky. Learning multiple layers of features from tiny images. *Master’s thesis, University of Tront*, 2009.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10713–10722, 2021.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*, 2018.
- Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. In *International Conference on Learning Representations*, 2020.
- Akis Linardos, Kaisar Kushibar, Sean Walsh, Polyxeni Gkontra, and Karim Lekadir. Federated learning for multi-center imaging diagnostics: a simulation study in cardiovascular disease. *Scientific Reports*, 12(1):1–12, 2022.
- David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30, 2017.
- Nathalie Majcherczyk, Nishan Srishankar, and Carlo Pinciroli. Flow-fl: Data-driven federated learning for spatio-temporal predictions in multi-robot systems. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 8836–8842, 2021. doi: 10.1109/ICRA48506.2021.9560791.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- Anh Nguyen, Tuong Do, Minh Tran, Binh X. Nguyen, Chien Duong, Tu Phan, Erman Tjiputra, and Quang D. Tran. Deep federated learning for autonomous driving. In *2022 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1824–1830, 2022. doi: 10.1109/IV51971.2022.9827020.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: an imperative style, high-performance deep learning library. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pp. 8026–8037, 2019.

Felix Sattler, Simon Wiedemann, Klaus-Robert Müller, and Wojciech Samek. Robust and communication-efficient federated learning from non-i.i.d. data. *IEEE Transactions on Neural Networks and Learning Systems*, 31(9):3400–3413, 2020.

Neta Shoham, Tomer Avidor, Aviv Keren, Nadav Israel, Daniel Benditkis, Liron Mor-Yosef, and Itai Zeitak. Overcoming forgetting in federated learning on non-iid data. *arXiv preprint arXiv:1910.07796*, 2019.

Chunnan Wang, Xiang Chen, Junzhe Wang, and Hongzhi Wang. Atpfl: Automatic trajectory prediction model design under federated learning framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6563–6572, 2022.

Blake E Woodworth, Kumar Kshitij Patel, and Nati Srebro. Minibatch vs local sgd for heterogeneous distributed learning. *Advances in Neural Information Processing Systems*, 33:6281–6292, 2020.

Chencheng Xu, Zhiwei Hong, Minlie Huang, and Tao Jiang. Acceleration of federated learning with alleviated forgetting in local training. In *International Conference on Learning Representations*, 2022.

Qiang Yang, Yang Liu, Yong Cheng, Yan Kang, Tianjian Chen, and Han Yu. Federated learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 13(3):1–207, 2019.

Hongyi Zhang, Jan Bosch, and Helena Holmström Olsson. Real-time end-to-end federated learning: An automotive case study. In *2021 IEEE 45th Annual Computers, Software, and Applications Conference (COMPSAC)*, pp. 459–468. IEEE, 2021a.

Li Zhang, Jianbo Xu, Pandi Vijayakumar, Pradip Kumar Sharma, and Uttam Ghosh. Homomorphic encryption-based privacy-preserving federated learning in iot-enabled healthcare system. *IEEE Transactions on Network Science and Engineering*, pp. 1–17, 2022. doi: 10.1109/TNSE.2022.3185327.

Lin Zhang, Yong Luo, Yan Bai, Bo Du, and Ling-Yu Duan. Federated learning for non-iid data via unified feature learning and optimization objective alignment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4420–4428, October 2021b.

Sicheng Zhao, Xiangyu Yue, Shanghang Zhang, Bo Li, Han Zhao, Bichen Wu, Ravi Krishna, Joseph E. Gonzalez, Alberto L. Sangiovanni-Vincentelli, Sanjit A. Seshia, and Kurt Keutzer. A review of single-source deep unsupervised visual domain adaptation. *IEEE Transactions on Neural Networks and Learning Systems*, 33(2):473–493, 2022.

Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

A APPENDIX

A.1 THE DUAL OF CGC QP

Here we provide the procedure of how to convert the primal of CGC QP into the dual problem. The primal problem is:

$$\begin{aligned} \min_{\mathbf{u}} \quad & \frac{1}{2} \mathbf{u}^T \mathbf{u} - \mathbf{h}^T \mathbf{u} + \frac{1}{2} \mathbf{h}^T \mathbf{h} \\ \text{s.t.} \quad & \mathbf{C} - \mathbf{z}^T \mathbf{u} \leq 0 \end{aligned} \tag{11}$$

where $\mathbf{u} \in \mathbb{R}^p$, $\mathbf{h} \in \mathbb{R}^p$ and $\mathbf{z} \in \mathbb{R}^p$, and p indicates the number of parameters of local model.

To obtain the dual of (11), we firstly conduct the Lagrange function $L(\mathbf{u}, v)$ with a Lagrange multiplier $v \in \mathbb{R}$ as follow

$$\begin{aligned}
L(\mathbf{u}, v) &= \frac{1}{2}\mathbf{u}^T\mathbf{u} - \mathbf{h}^T\mathbf{u} + \frac{1}{2}\mathbf{h}^T\mathbf{h} + v(C - \mathbf{z}^T\mathbf{u}) \\
&= \frac{1}{2}\mathbf{u}^T\mathbf{u} - \mathbf{h}^T\mathbf{u} + \frac{1}{2}\mathbf{h}^T\mathbf{h} + vC - v\mathbf{z}^T\mathbf{u} \\
&= \frac{1}{2}\mathbf{u}^T\mathbf{u} - (\mathbf{h}^T + v\mathbf{z}^T)\mathbf{u} + \frac{1}{2}\mathbf{h}^T\mathbf{h} + vC
\end{aligned} \tag{12}$$

Because (12) is a quadratic convex function, its minimal value can be obtained by

$$\nabla_{\mathbf{u}}L(\mathbf{u}, v) = \mathbf{u} - (\mathbf{h} + v\mathbf{z}) = 0 \tag{13}$$

$$\implies \mathbf{u} = \mathbf{h} + v\mathbf{z} \tag{14}$$

Hence, the minimal value of (12) is

$$\begin{aligned}
L(\mathbf{h} + v\mathbf{z}, v) &= \frac{1}{2}(\mathbf{h} + v\mathbf{z})^T(\mathbf{h} + v\mathbf{z}) - \mathbf{h}^T(\mathbf{h} + v\mathbf{z}) + \frac{1}{2}\mathbf{h}^T\mathbf{h} + v(C - \mathbf{z}^T(\mathbf{h} + v\mathbf{z})) \\
&= \frac{1}{2}(\mathbf{h}^T + v\mathbf{z}^T)(\mathbf{h} + v\mathbf{z}) - \mathbf{h}^T(\mathbf{h} + v\mathbf{z}) + \frac{1}{2}\mathbf{h}^T\mathbf{h} + vC - v\mathbf{z}^T(\mathbf{h} + v\mathbf{z}) \\
&= -\frac{1}{2}v^2\mathbf{z}^T\mathbf{z} - v\mathbf{h}^T\mathbf{z} + vC
\end{aligned} \tag{15}$$

The dual function $G(v)$ is

$$G(v) = \inf_{\mathbf{u}} L(\mathbf{u}, v) \tag{16}$$

and the dual problem is

$$\begin{aligned}
&\max_v G(v) \\
&\text{s.t. } v \geq 0
\end{aligned} \tag{17}$$

Therefore, the dual of CGC QP is

$$\begin{aligned}
&\min_v \frac{1}{2}v^2\mathbf{z}^T\mathbf{z} + v(\mathbf{h}^T\mathbf{z} - C) \\
&\text{s.t. } v \geq 0
\end{aligned} \tag{18}$$

Problem (18) can be solved by *quadprog* library and the optimum v^* is obtained. The standard form of QP problem in *quadprog* is

$$\begin{aligned}
&\min_{\mathbf{x}} \frac{1}{2}\mathbf{x}^T\mathbf{G}\mathbf{x} - \mathbf{a}^T\mathbf{x} \\
&\text{s.t. } \mathbf{C}^T\mathbf{x} \leq \mathbf{b}
\end{aligned} \tag{19}$$

At last, the optimal solution of the primal CGC OP is calculated by $\mathbf{u}^* = \mathbf{h} + v^*\mathbf{z}$.

A.2 THE DUAL OF SGC QP

Here we provide the procedure of how to convert the primal of SGC QP into the dual problem. The primal problem is:

$$\begin{aligned}
&\min_{\mathbf{z}} \frac{1}{2}\mathbf{z}^T\mathbf{z} - \mathbf{g}^T\mathbf{z} + \frac{1}{2}\mathbf{g}^T\mathbf{g} \\
&\text{s.t. } \mathbf{C} - \mathbf{G}\mathbf{z} \leq 0
\end{aligned} \tag{20}$$

where $\mathbf{z} \in \mathbb{R}^p$, $\mathbf{g} \in \mathbb{R}^p$, and $\mathbf{G} \in \mathbb{R}^{|\mathcal{K}| \times p}$. Problem (20) is a QP on p variables.

Its Lagrange function with a Lagrange multiplier $\boldsymbol{\lambda} \in \mathbb{R}^{|\mathcal{K}|}$ is

$$\begin{aligned}
L(\mathbf{z}, \boldsymbol{\lambda}) &= \frac{1}{2} \mathbf{z}^T \mathbf{z} - \mathbf{g}^T \mathbf{z} + \frac{1}{2} \mathbf{g}^T \mathbf{g} + \boldsymbol{\lambda}^T (\mathbf{C} - \mathbf{G} \mathbf{z}) \\
&= \frac{1}{2} \mathbf{z}^T \mathbf{z} - (\mathbf{g}^T + \boldsymbol{\lambda}^T \mathbf{G}) \mathbf{z} + \frac{1}{2} \mathbf{g}^T \mathbf{g} + \boldsymbol{\lambda}^T \mathbf{C}
\end{aligned} \tag{21}$$

Since $L(\mathbf{z}, \boldsymbol{\lambda})$ is a quadratic convex function, its minimal value can be obtained by

$$\begin{aligned}
\nabla_{\mathbf{z}} L(\mathbf{z}, \boldsymbol{\lambda}) &= \mathbf{z} - (\mathbf{g}^T + \boldsymbol{\lambda}^T \mathbf{G}) = 0 \\
\implies \mathbf{z} &= \mathbf{g} + \mathbf{G}^T \boldsymbol{\lambda}
\end{aligned} \tag{22}$$

So

$$\begin{aligned}
L(\mathbf{g} + \mathbf{G}^T \boldsymbol{\lambda}, \boldsymbol{\lambda}) &= \frac{1}{2} (\mathbf{g} + \mathbf{G}^T \boldsymbol{\lambda})^T (\mathbf{g} + \mathbf{G}^T \boldsymbol{\lambda}) - \mathbf{g}^T (\mathbf{g} + \mathbf{G}^T \boldsymbol{\lambda}) + \frac{1}{2} \mathbf{g}^T \mathbf{g} + \boldsymbol{\lambda}^T \mathbf{C} - \boldsymbol{\lambda}^T \mathbf{G} (\mathbf{g} + \mathbf{G}^T \boldsymbol{\lambda}) \\
&= -\frac{1}{2} \boldsymbol{\lambda}^T \mathbf{G} \mathbf{G}^T \boldsymbol{\lambda} - \frac{1}{2} \mathbf{g}^T \mathbf{G}^T \boldsymbol{\lambda} - \frac{1}{2} \boldsymbol{\lambda}^T \mathbf{G} \mathbf{g} + \boldsymbol{\lambda}^T \mathbf{C} \\
&= -\frac{1}{2} \boldsymbol{\lambda}^T \mathbf{G} \mathbf{G}^T \boldsymbol{\lambda} - \mathbf{g}^T \mathbf{G}^T \boldsymbol{\lambda} + \boldsymbol{\lambda}^T \mathbf{C} \\
&= -(\mathbf{G} \mathbf{g} - \mathbf{C})^T \boldsymbol{\lambda} - \frac{1}{2} \boldsymbol{\lambda}^T \mathbf{G} \mathbf{G}^T \boldsymbol{\lambda}
\end{aligned} \tag{23}$$

Therefore, the dual of SGC QP is

$$\begin{aligned}
\min_{\boldsymbol{\lambda}} \quad & (\mathbf{G} \mathbf{g} - \mathbf{C})^T \boldsymbol{\lambda} + \frac{1}{2} \boldsymbol{\lambda}^T \mathbf{G} \mathbf{G}^T \boldsymbol{\lambda} \\
\text{s.t.} \quad & \boldsymbol{\lambda} \geq 0
\end{aligned} \tag{24}$$

The optimum $\boldsymbol{\lambda}^*$ can be calculated by library *quadprog* and the optimal solution of SGC QP is $\mathbf{z}^* = \mathbf{g} + \mathbf{G}^T \boldsymbol{\lambda}^*$.

A.3 THEORETICAL ANALYSIS OF GRADIENT PROJECTION

Assumption 1 *In each local update, small optimization steps happen and thus we can assume that the function F is locally linear (i.e., convex function).*

To enable the local model retains more knowledge received from server after local update, the loss of local model on global data \mathcal{D} should reduce after local update, i.e.,

$$F(\boldsymbol{\theta}_k^t, \mathcal{D}) < F(\boldsymbol{\theta}^{t-1}, \mathcal{D}) \tag{25}$$

where $\boldsymbol{\theta}_k^t$ is the model parameters at client k after local update and $\boldsymbol{\theta}^{t-1}$ is the initial model (i.e., the global model at previous round) parameters.

Since the global data \mathcal{D} cannot be achieved in FL, we use gradient projection to achieve this objective. For simplifying, \mathcal{D} is omitted in the following equations.

From convexity we know that $\nabla F(\boldsymbol{\theta}^{t-1})^T (\boldsymbol{\theta}_k^t - \boldsymbol{\theta}^{t-1}) \geq 0$ implies $F(\boldsymbol{\theta}_k^t) \geq F(\boldsymbol{\theta}^{t-1})$, so the search direction in local training must satisfy

$$\nabla F(\boldsymbol{\theta}^{t-1})^T (\boldsymbol{\theta}_k^t - \boldsymbol{\theta}^{t-1}) < 0 \tag{26}$$

Thus, it must make an acute angle with negative gradient, i.e.,

$$(\mathbf{g}^{t-1})^T \mathbf{g}_k^t > 0 \Leftrightarrow \langle \mathbf{g}^{t-1}, \mathbf{g}_k^t \rangle > 0 \tag{27}$$

where $\mathbf{g}^{t-1} = -\nabla F_k(\boldsymbol{\theta}^{t-1})$ and $\mathbf{g}_k^t = \boldsymbol{\theta}_k^t - \boldsymbol{\theta}^{t-1}$.

A.4 CONVERGENCE ANALYSIS

Here we give a simple proof that our FedGC has a faster convergence rate than FedAvg. We analyze the convergence of our FedGC by finding an upper bound ξ of $\mathbb{E}[F(\boldsymbol{\theta}^T)] - F^*$, i.e., $\mathbb{E}[F(\boldsymbol{\theta}^T)] - F^* \leq \xi$, where $F(\boldsymbol{\theta}^T)$ represents the final global model with parameter $\boldsymbol{\theta}^T$ in FL and F^* is the optimal model for all clients' data (i.e., the upper bound of global model in FL).

From A.3, our FedGC can guarantee that the loss of local model on global data reduces after local update (i.e., $F(\boldsymbol{\theta}_k^t) < F(\boldsymbol{\theta}^{t-1})$), while FedAvg cannot. In other words, our FedGC can guarantee that the loss of global model gradually reduces (i.e., $F(\boldsymbol{\theta}^t) < F(\boldsymbol{\theta}^{t-1})$), while FedAvg cannot. In detail,

$$F(\boldsymbol{\theta}^t) = F\left(\sum_{k=1}^K p_k \boldsymbol{\theta}_k^t\right) \quad (28)$$

where the aggregation is averaged for simplifying, and $\sum_{k=1}^K p_k = 1$.

Let Assumption 1 holds, we can obtain

$$F(\boldsymbol{\theta}^t) = F\left(\sum_{k=1}^K p_k \boldsymbol{\theta}_k^t\right) = \sum_{k=1}^K p_k F(\boldsymbol{\theta}_k^t) < \sum_{k=1}^K p_k F(\boldsymbol{\theta}^{t-1}) = F(\boldsymbol{\theta}^{t-1}) \quad (29)$$

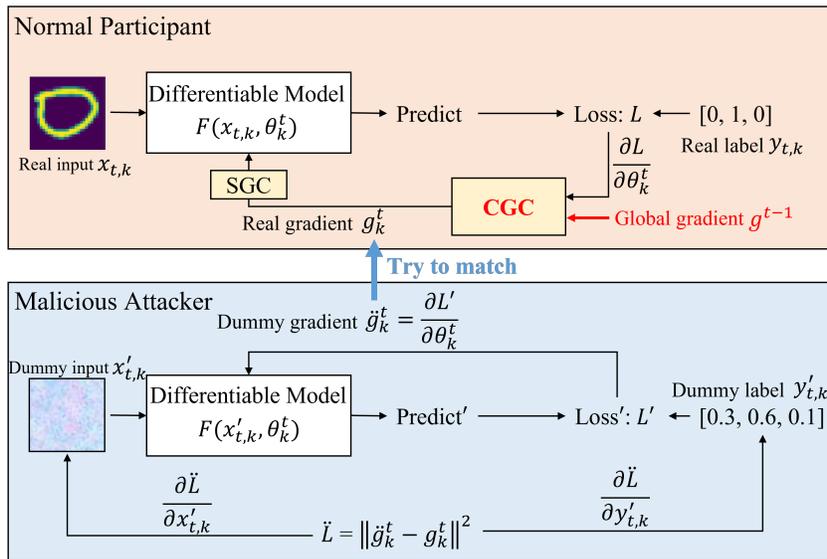
Thus, after T communication rounds (i.e., one fixed T), the difference between $\mathbb{E}[F(\boldsymbol{\theta}^T)]$ and F^* of our FedGC should be smaller than that of FedAvg. Thus, we can conclude that $\xi^{GC} < \xi^{Avg}$.

A.5 PRIVACY-PRESERVING ANALYSIS

Recently, Deep Leakage from Gradients (DLG) Zhu et al. (2019) attack has attracted growing attentions, which can completely steal real data from gradients. However, under our FedGC, the real data cannot be reversely obtained. The details are shown below.

In DLG, a pair of ‘‘dummy’’ input and label are first randomly generated to perform the usual forward and backward operations. In order to obtain the real data reversely, the dummy gradients from the dummy data are firstly derived, then DLG optimizes the dummy inputs and labels by minimizing the Euclidean distance between the dummy gradients and the real gradients. **However, matching the gradients cannot make the dummy data close to the real data when our FedGC performs.**

Figure 4: DLG in our FedGC.



In Figure 4, at t -th communication round, an honest client k samples a minibatch $(x_{t,k}, y_{t,k})$ from its own data and an evil client randomly initializes a dummy input $x'_{t,k}$ with dummy label $y'_{t,k}$. The

objective of the evil client is

$$x'_{t,k}, y'_{t,k} = \arg \min_{x'_{t,k}, y'_{t,k}} \|\ddot{\mathbf{g}}_k^t - \mathbf{g}_k^t\|^2 \quad (30)$$

where

$$\ddot{\mathbf{g}}_k^t = \frac{\partial L' (F(x'_{t,k}, \theta_k^t), y'_{t,k})}{\partial \theta_k^t} \quad (31)$$

$$\mathbf{g}_k^t = CGC \left(\frac{\partial L (F(x_{t,k}, \theta_k^t), y_{t,k})}{\partial \theta_k^t}, \mathbf{g}^{t-1} \right) \quad (32)$$

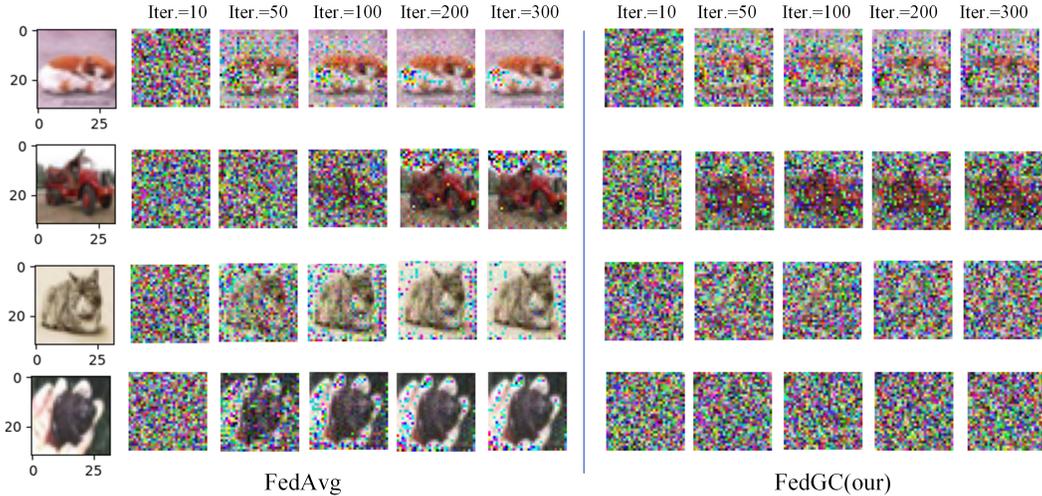
From Eq.(32), we can observe that in our FedGC, \mathbf{g}_k^t is not only determined by $\frac{\partial L}{\partial \theta_k^t}$ but also dependent on the global gradient \mathbf{g}^{t-1} by our CGC, and the CGC projection cannot be reversely derived. Thus, we can obtain

$$(x'_{t,k}, y'_{t,k}) \neq (x_{t,k}, y_{t,k}) \quad (33)$$

Therefore, we can conclude that real data cannot be recovered with gradient inversion attacks (e.g., DLG) in our FedGC and our FedGC can ensure preserving of privacy.

In addition, we conduct the gradient inversion attacks experiment to compare FedAvg and our FedGC with DLG on CIFAR10-10 dataset. The backbone is ResNet-9. We set the number of iterations in DLG to 300. As shown in Figure 5, the quality of the images recovered from our FedGC is significantly worse than those recovered from FedAvg, exhibiting a better privacy protection capability of our FedGC.

Figure 5: Images recovered from updated gradients of FedAvg and FedGC on CIFAR10-10.



A.6 MORE DATASETS DETAILS

Handwritten-Digits MNIST has a training sets of 60,000 examples and a test set of 10,000 examples. The images in MNIST are grayscale with image size 28×28 . MNIST-M⁴ contains 59,001 training and 90,001 test RGB images with size 28×28 . USPS contains 9,298 grayscale images with image size 16×16 . There are 73,257 RGB images in training set of SVHN⁵ and 26,032 RGB images in test set with image size 32×32 .

CIFAR-10 and **CIFAR-100** CIFAR-10 contains 60,000 32×32 color images of 10 classes, including 5,000 images in training set and 1,000 images in test set per class. CIFAR-100 consists of 60,000

⁴<http://yaroslav.ganin.net/>

⁵<http://ufldl.stanford.edu/housenumbers/>

32×32 color images of 100 classes. In CIFAR-100, there are 500 images in training set and 100 images in test set per class.

A.7 EXPERIMENT OF FEDAVG ON HWDIGITS-4 UNDER MINI-BATCH UPDATE

In this experiment, the number of mini-batches B for local training in FedAvg is in the range of $\{1, 10, 30, 50\}$, and the number of communication rounds T is 600 for all experiments.

Table 3: Experiment of FedAvg on HWDigits-4 under mini-batch update.

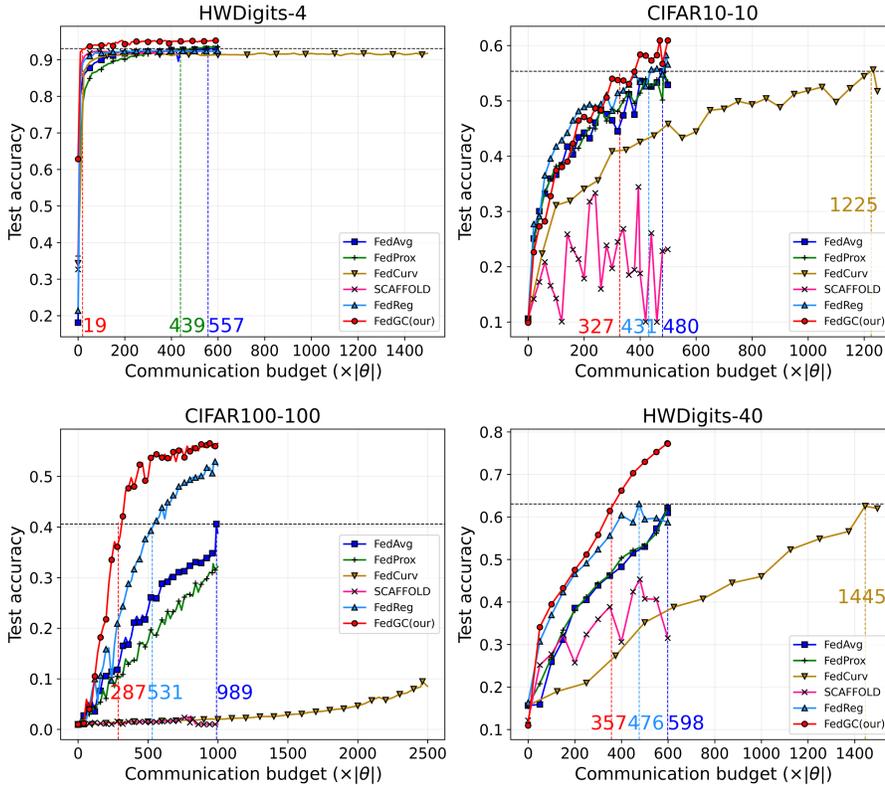
Setting	$B = 1$	$B = 10$	$B = 20$	$B = 50$	$E = 1$ (i.e., $B = 73 \sim 732$)
FedAvg	85.22	92.75	93.43	93.59	93.07

In Table 3, we can observe that FedAvg can achieve better performance when $B = 50$ while requiring less iterations in local training stage compared to $E = 1$. The experimental results are consistent with the statement in Section 3.1.

A.8 LEARNING CURVES

The learning curves are drawn in Figure 6, where the x-axis is scaled by communication-budget (i.e., in multiples of $|\theta|$). The figures show that our FedGC can reach the target accuracy (i.e., the highest accuracy of baseline FedAvg) with the smallest communication budget.

Figure 6: Learning curves of compared methods in terms of communication budget.



A.9 ABLATION EXPERIMENTS

We conduct the ablation experiment to disentangle CGC & SGC on all four compared datasets.

The results in Table 4 illustrate that the performances of our FedGC will degrade without CGC, in particular for the label skew datasets. That is because the issue of catastrophic forgetting seriously

Table 4: Ablation experiments.

Method	HWDigits-4	CIFAR10-10	CIFAR100-100	HWDigits-40
FedAvg	93.07	55.34	40.59	62.11
FedReg	92.32	57.70	53.02	63.05
FedGC (w/o CGC)	93.38	54.09	41.64	74.51
FedGC (w/o SGC)	94.15	59.64	55.73	74.63
FedGC	95.37	60.92	56.56	77.33

influences the performance of FL, and our CGC can mitigate this forgetting issue. Therefore, in our FedGC, CGC is very important to improve the performance of FL on Non-IID data.

A.10 AVERAGE LOSSES OF EACH CLIENT ON GLOBAL DATA BEFORE AND AFTER LOCAL TRAINING IN FEDAVG

Here we give an example to showcase that the forgetting issue is important. For FedAvg on CIFAR10-10, we calculate L^{t-1} and L_k^t to show the average losses of each client on global data before and after local training respectively.

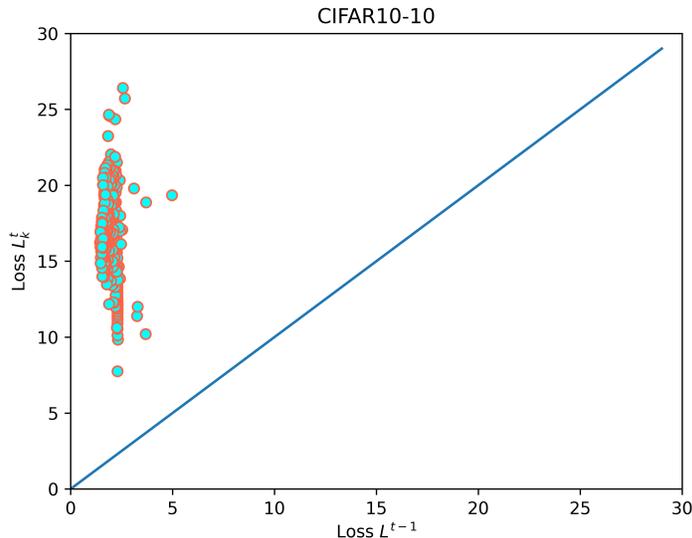
$$L^{t-1} = \frac{1}{|K|} \sum_{k=1}^K L(\theta_k^{t-1}, \mathcal{D}) \quad (34)$$

where K is the number of clients, θ_k^{t-1} is the initial model before local training, \mathcal{D} is the global data.

$$L_k^t = \frac{1}{|K|} \sum_{k=1}^K \frac{1}{|K|} \sum_{i=1}^K L(\theta_i^t, \mathcal{D}_k) \quad (35)$$

where θ_i^t is the local model after local training and \mathcal{D}_k is the local data.

Figure 7: Average losses of each client on global data before (L^{t-1}) and after (L_k^t) local training in FedAvg, and each point indicates one communication round.



The Figure 7 shows that in FedAvg, the average losses L_k^t on global data for all communication rounds are significantly larger than L^{t-1} after local training. This is because the clients forget the knowledge learned from global model after local training in FedAvg.