

# Convergence Visualizer of Decentralized Federated Distillation with Reduced Communication Costs

1<sup>st</sup> Akihito Taya

*Institute of Industrial Science,  
The University of Tokyo  
Tokyo, JAPAN  
taya-a@iis.u-tokyo.ac.jp*

2<sup>nd</sup> Yuuki Nishiyama

*Center for Spatial Information Science,  
The University of Tokyo  
Chiba, JAPAN  
yuukin@iis.u-tokyo.ac.jp*

3<sup>rd</sup> Kaoru Sezaki

*Center for Spatial Information Science,  
The University of Tokyo  
Chiba, JAPAN  
sezaki@iis.u-tokyo.ac.jp*

**Abstract**—Federated learning (FL) achieves collaborative learning without the need for data sharing, thus preventing privacy leakage. To extend FL into a fully decentralized algorithm, researchers have applied distributed optimization algorithms to FL by considering machine learning (ML) tasks as parameter optimization problems. Conversely, the consensus-based multi-hop federated distillation (CMFD) proposed in the authors' previous work makes neural network (NN) models get close with others in a function space rather than in a parameter space. Hence, this study solves two unresolved challenges of CMFD: (1) communication cost reduction and (2) visualization of model convergence. Based on a proposed dynamic communication cost reduction method (DCCR), the amount of data transferred in a network is reduced; however, with a slight degradation in the prediction accuracy. In addition, a technique for visualizing the distance between the NN models in a function space is also proposed. The technique applies a dimensionality reduction technique by approximating infinite-dimensional functions as numerical vectors to visualize the trajectory of how the models change by the distributed learning algorithm.

**Index Terms**—federated learning, multi-hop networks, consensus algorithm, decentralized machine learning, visualization

## I. INTRODUCTION

Owing to its communication-efficiency and privacy-preserving properties, federated learning (FL) is a promising framework for the collaborative learning of Internet of Things (IoT) sensor devices, wearable sensors, smartphones, smart homes, and connected vehicles [1]–[3]. McMahan *et al.* [1] originally proposed FL as a distributed machine learning (ML) algorithm to prevent the leakage of user privacy included in local data samples by avoiding the upload of user data to a central server. The emergence of FL has garnered considerable attention, and many researchers have extended FL to address the remaining challenges, e.g., non-independent and identically distributed (IID) datasets, as well as further reduction of communication costs.

One problem is that typical FL algorithms require a central server to manage the learning process and aggregate the collected parameters to update a global neural network

(NN) model. In centralized FL (FL), the server capacity for computation and communication can be a bottleneck, and there is a concern regarding the Big Tech monopoly of well-trained models. However, decentralized FL (FL) [4]–[6] is an important extension for solving such problems. In DFL, client devices communicate directly with others to share the model-update information without a central server. Consequently, DFL is suitable for IoT applications in which sensor devices are connected via wireless sensor networks (WSN).

We proposed consensus-based multi-hop federated distillation (CMFD) [7] as a DFL using knowledge distillation (KD). The concept of CMFD is to solve convex functional optimization rather than non-convex parameter optimization. Although NNs are non-convex, ML tasks can be convex problems in terms of functional optimization. Therefore, CMFD attempts to solve the optimization problem in a distributed manner by directly aggregating functions rather than the NN parameters. This concept can be realized by decentralizing algorithms of federated distillation (FD). FD is a variation of FL, in which devices share the output values of their local models, called distilled values, rather than the parameters.

Although FD requires lower communication costs than parameter-aggregation-based FL, there is a potential to reduce the amount of transferred data. In addition, the communication cost can be reduced if the shared distilled values are limited; however, the prediction performance is degraded. To overcome the accuracy-communication tradeoff, this paper proposes a dynamic communication cost reduction method (DCCR) as an extension of CMFD. DCCR enables devices to share a few distilled values with others at the beginning of training and gradually increase the amount of shared data to improve prediction accuracy.

A key feature of the DCCR is communication-less synchronization of the distributed random selection. Decentralized FD (FD) must share input values with other devices to aggregate the distilled values. Hence, the devices should synchronize the information indicating which samples are used to calculate the distilled values. To suppress bothersome communication, the DCCR leverages a pseudo-random number generator (PRNG)

to obtain randomly selected synchronized data samples in a distributed manner.

Another challenge in DFD is validating model convergence properties. In contrast to CFL, where devices easily synchronize their models to the global one that is managed by a server, local models may not converge to the same one in finite time with a decentralized algorithm, even though the theoretical analysis guarantees convergence assuming infinite time. Moreover, when KD is applied, our interest lies in the convergence of prediction functions rather than that of the NN parameters. Because the convergence of models cannot be quantified directly, a novel framework for validating their convergence should be developed.

To validate the convergence properties of the trained models, we develop a visualization framework that visualizes how the prediction models trained by the local devices approach others in a function space. Although dimensionality-reduction techniques have been widely adopted for high-dimensional data visualization, their targets are high-dimensional data, which implies that they cannot be applied to NN models because they have infinite dimensions as functions. Therefore, we approximate these functions as finite-dimensional vectors such that the distances between the two functions remain unchanged. Subsequently, a dimensionality-reduction technique is applied to the finite-dimensional vectors to obtain a low-dimensional projection.

The contributions of this paper are summarized as follows:

- We propose a DCCR as an extension of CMFD. Devices first utilize a few values for distillation to suppress the communication costs and gradually increase the amount of shared information to improve the prediction accuracy toward the end of the training. In addition, we leverage a PRNG to share information of random sampling with zero-cost communication.
- We provide a scheme for visualizing how local models approach each other in a function space. This scheme leverages a dimensionality-reduction technique to project local models onto a 2D/3D space, demonstrating that the local models get close to each other in the function space, even with the proposed communication-efficient method.

## II. RELATED WORKS

### A. Decentralized and distillation-based federated learning

In the context of the IoT, where sensor devices are connected via WSN, DFL has advantages over centralized algorithms. Although CFL is the predominant approach in FL, the communication and computational capacity of a central server can become challenging when the number of participating devices increases. Hence, DFL has been developed to address this limitation [4], [5], [8], [9]. The DFL assumes that the client devices are connected via multi-hop networks and communicate directly with each other to share information.

By extending FedAvg [1], DFL algorithms that share the parameters of the NN models among devices have been proposed [5], [8], [9]. However, these schemes require numerous

communication resources because of their numerous parameters. Conversely, DFD algorithms send the output values of local prediction models, referred to as distilled values, to adjacent devices [10], [11]. Our previous study [7] was also categorized as DFD. DFD reduces the communication cost of FL because the number of distilled values is typically less than the number of parameters of the deep neural network (DNN); however, DFD requires a number of distilled values to improve the training performance. Therefore, a communication-efficient method for the DFDs should be developed.

There are several methods for reducing the communication costs of FL [12]–[15]. For instance, model pruning [12], [13] removes redundant structures and parameters from NN models. Although this technique reduces the communication cost of sharing model gradients, it is not suitable for FD where model gradients are not shared. Another approach to suppress network traffic is client selection, in which only a few clients upload their model updates to the server [14], [15]. However, this cannot be applied to decentralized algorithms. In addition, quantization is a basic approach for reducing communication costs. In FD, distilled values, rather than model parameters, are quantized [16]. Hence, we focus on reducing the number of distilled values utilized in each communication round, which can be easily combined with compression methods.

### B. Visualization of neural networks

Uddin *et al.* [17] visualized the progress of CFL by plotting the mutual information between the ground truth information and output logits produced by every global model in each communication round. Although this technique can visualize the relationship between the ground truth and the global model at the central server, the relationship between more than two devices cannot be visualized in decentralized scenarios.

Dimensionality-reduction techniques are suitable approaches for visualization. They have been developed to project high-dimensional data onto a lower-dimensional space while preserving the original data structure and relationships within the data. Principal component analysis (PCA) is a traditional technique that utilizes singular value decomposition (SVD) [18] to linearly transform high-dimensional data into a lower-dimensional space. In addition, uniform manifold approximation and projection (UMAP) [19] and t-distributed stochastic neighbor embedding (t-SNE) [20] are nonlinear dimensionality-reduction techniques used for high-dimensional-data visualization. UMAP approximates the high-dimensional manifold structure and optimizes lower-dimensional embedding by minimizing the cross-entropy between the high- and low-dimensional representations. We adopt UMAP because it tends to output more separated clusters than t-SNE.

## III. PROBLEM STATEMENT AND THE CONCEPT OF CMFD

We consider a system where some devices with sensors are connected via a multi-hop network. Each device  $i$  trains its local model  $f_i$  utilizing a local dataset  $\mathcal{D}_i$ . To avoid privacy leakage, the dataset  $\mathcal{D}_i$  is not shared with other devices In

contrast to a typical DFL, where the devices share information of parameter updates, the devices share distilled values, i.e., the outputs of  $f_i$ . When calculating the distilled values, all devices utilize the same input  $\mathcal{D}_s$ , which is assumed to be a public dataset or an artificially generated dataset [21].

The key concept proposed in our previous work [7] was to operate a consensus-based optimization algorithm in a function space to avoid non-convexity problems. Although training NNs is typically defined as the optimization of parameters, parameter optimization encounters challenges because of non-convexity in the DFL context. Therefore, prediction functions rather than parameters should be optimized as follows:

$$\underset{f \in L_\mu^2}{\text{minimize}} \quad L_\mu(f) = \sum_{i \in \mathcal{U}} L_{\mu_i}(f), \quad (1)$$

where  $\mu$ ,  $\mu_i$ , and  $L_\mu^2$  represent the global and local probability measure of the input space, and the  $L^2$  space with respect to the probability measure  $\mu$ , respectively. We denote  $L_\mu$  and  $L_{\mu_i}$  as the global and local loss functions, respectively. Notably, typical loss functions (e.g., Kullback-Leibler (KL) divergence [22] and mean squared error (MSE)) satisfy the following inequities for all  $f_1, f_2$ , and  $k \in [0, 1]$ :

$$L_\mu(tf_1 + (1-t)f_2) \leq tL_\mu(f_1) + (1-t)L_\mu(f_2), \quad (2)$$

$$L_{\mu_i}(tf_1 + (1-t)f_2) \leq tL_{\mu_i}(f_1) + (1-t)L_{\mu_i}(f_2). \quad (3)$$

This implies that (1) is a convex functional optimization problem, although optimizing the parameters  $w_i$  of the NN is a non-convex problem. Hence, CMFD focuses on updating the prediction models  $f$  directly in the direction of reducing  $L_\mu(f)$  when aggregating the local update information.

In CMFD, the devices are allowed to have different parameter sets at the end provided that the trained models are the same because functions are treated as identical in  $L_\mu^2$  iff the distance between the functions is equal to zero as follows:

$$d_{L_\mu^2}(f_1, f_2) := \sqrt{\int_{\mathcal{X}} |d_{\mathcal{Y}}(f_1(\mathbf{x}), f_2(\mathbf{x}))|^2 d\mu} = 0, \quad (4)$$

where  $\mathcal{X} \subseteq \mathbb{R}^X$ ,  $\mathcal{Y} \subseteq \mathbb{R}^Y$ ,  $d_{L_\mu^2}$  and  $d_{\mathcal{Y}}$  denote the input and output space, and the distance of functions and output values, respectively. It should be noted that functions with different parameters can be the same because of their non-convexity.

In CMFD, the devices iterate the following two steps in the  $t$ th communication round:

- 1) Update the local NN parameters  $w_i^t$  by stochastic gradient descent (SGD) utilizing  $\mathcal{D}_i$ .
- 2) Update  $w_i^t$  by KD using the received distilled values  $y_{j,x}^t$  from adjacent devices  $j \in \mathcal{N}_i$ .

These steps emulate consensus-based optimization in the function space  $L_\mu^2$ , and therefore the NN models  $f_i$  are optimized appropriately because of the convexity of (1).

#### IV. DYNAMIC COMMUNICATION COST REDUCTION

Compared with FL, FD reduces communication costs because the dimensions of the distilled values are fewer than the parameters of DNN. However, there are demands to further

---

#### Algorithm 1 Pseudo code of DCCR

---

**Require:** Local and shared train data:  $\mathcal{D}_i, \mathcal{D}_s$ ,  
Learning and sharing rate:  $\eta_t, \varepsilon$

```

1: while not converged do
2:   for all device  $i = 1, \dots, n$  do
3:     for minibatch  $\tilde{\mathcal{D}}$  in  $\mathcal{D}_i$  do
4:        $\hat{w}_i^t \leftarrow w_i^t - \eta_t \nabla_w \text{Loss}(f(\cdot; w_i^t), \tilde{\mathcal{D}})$ 
5:     end for
6:     Select subset  $\tilde{\mathcal{D}}^t$  of  $\mathcal{D}_s$  using the shared PRNG.
7:     for all  $x \in \tilde{\mathcal{D}}^t$  do
8:        $y_{i,x}^t \leftarrow f(x; \hat{w}_i^t)$ 
9:     end for
10:    send  $y_{i,x}^t$  to neighbor devices
11:   end for
12:   for all device  $i = 1, \dots, n$  do
13:     receive  $y_{j,x}^t$  from neighbor devices
14:     for minibatch  $\tilde{\mathcal{D}}$  in  $\tilde{\mathcal{D}}^t$  do
15:        $w_i^{t+1} \leftarrow \hat{w}_i^t - \varepsilon n_i \nabla_w c(\hat{w}_i^t)$ 
16:     end for
17:   end for
18: end while

```

---

reduce the amount of data transferred in WSN because of limited communication resources. Communication costs can be reduced by limiting the dataset  $\mathcal{D}_s$  utilized to calculate the distilled values  $y_{i,x}^t$ , which, unfortunately, results in performance degradation. To overcome the tradeoff between communication costs and prediction accuracy, we extend CMFD by utilizing DCCR in which the devices limit the distilled values to be shared at the beginning of training to reduce communication costs, whereas the number of the distilled values used for training is gradually increased to improve the prediction accuracy.

The pseudo-code for the proposed algorithm is presented in Alg. 1. Each device performs local SGD to update the local parameters  $w_i^t$  on lines 3–5. After the local SGD phase, the device randomly selects  $n_d^t$  data samples from  $\mathcal{D}_s$  to make a temporary dataset  $\tilde{\mathcal{D}}^t$  on Line 6.  $\tilde{\mathcal{D}}^t$  is selected such that the same samples are included across devices in a distributed manner, as explained in the last part of this section. Using the selected samples  $\tilde{\mathcal{D}}^t$ , each device calculates the output values of its local model  $f_i := f(\cdot; \hat{w}_i^t)$  in lines 7–9. The output values, which are referred to as the distilled values  $y_{i,x}^t$ , are sent to adjacent devices (line 10). After receiving the distilled values, each device updates its local model (line 15).

In contrast to parameter-aggregation-based FL algorithms, which calculate the weighted average of the parameters  $\hat{w}_i^t$ , FD performs SGD for model aggregation to make the models closer in a function space. The concept of model aggregation is illustrated in Fig. 1. The devices share the distilled values  $y_{i,x}^t$ , and each device calculates the average of the received distilled values, which becomes the new target value. To make

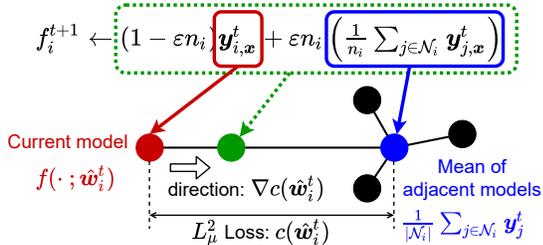


Fig. 1: Concept of distillation-based function aggregation. The local model approaches the other models in a function space by SGD, treating the mean of the adjacent output values as a target value.

the local models closer to other models in the function space, the devices update their local models using the gradient of the distance between their local models and the obtained target values, that is, the devices perform SGD according to the following criteria:

$$c(\hat{\mathbf{w}}_i^t) := \sum_{\mathbf{x} \in \tilde{\mathcal{D}}^t} \left\| f(\mathbf{x}; \hat{\mathbf{w}}_i^t) - \frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} \mathbf{y}_{j,\mathbf{x}}^t \right\|^2, \quad (5)$$

where  $|\cdot|$  denotes the cardinality of a set. The second term in the norm is the target value in the  $t$ th communication round.

When selecting  $\tilde{\mathcal{D}}^t$  from  $\mathcal{D}_s$ , the DCCR dynamically changes the number of data samples  $n_d^t$ , which is determined as follows:

$$n_d^t := d \left( \left\lfloor \frac{t}{\tau} \right\rfloor + 1 \right), \quad (6)$$

where  $\lfloor \cdot \rfloor$  denotes the floor function. Coefficients  $d$  and  $\tau$  control parameters for determining the rate of increase of  $n_d^t$ . When  $\tau$  is equal to  $Td/|\mathcal{D}_s|$ , all data samples in  $\mathcal{D}_s$  are utilized in the last  $\tau$  round of training. Because  $n_d^t$  increases with communication round  $t$ , DCCR suppresses the traffic load at the beginning of training while updating the NN models carefully utilizing sufficient distilled values at the end.

#### Zero-cost information sharing of random samples:

Selecting  $\tilde{\mathcal{D}}^t$  in a distributed manner requires a strategy to share the selected samples to address the accuracy-communication tradeoff. When all devices utilize all the samples, the sample-selection information does not need to be shared. In DCCR, however, if each device selects different data samples in a distributed manner, the cost function  $c(\hat{\mathbf{w}}_i^t)$  cannot be calculated. To suppress the costs of sharing sample-selection information, the DCCR leverages a PRNG that is included in the standard libraries of daily used programming languages.

Although the PRNG approximates the properties of random sequences, we utilize its deterministic features in DCCR. The PRNG generates the same series from the same initial key, which is called the seed. Therefore, if the devices share the information of the initial key of the PRNG before beginning

training, the devices can obtain the same randomly selected data samples  $\tilde{\mathcal{D}}^t$  in a distributed manner without communicating with the others at each iteration.

## V. VISUALIZATION OF FUNCTION CONVERGENCE

As described in Sec. III, our algorithm focuses on updating the NN models to converge in the function space rather than in the parameter space. This section provides a scheme that maps functions from  $L_\mu^2$  to positions in 2D or 3D space to visualize the proximity of the NN models to others in the function space. Using this scheme, we can observe the dynamics of the model updates in the training phase.

Dimensionality-reduction techniques are typically adopted to project high-dimensional data onto a low-dimensional space. However, the prediction models  $f_i$  have infinite dimensions, indicating that dimensionality-reduction techniques cannot be applied directly. Considering that we are interested in the distance between the models in the convergence analysis, it is sufficient that a visualized image reflects the distance between the prediction models. Therefore, the proposed visualization scheme evaluates the empirically approximated distances of the prediction models.

The ideal distance between the two models is defined by (4), which cannot be calculated when using a NN with numerous parameters. Hence, we use the following empirical representation:

$$d_{L_\mu^2}(f_i, f_j) \approx \sqrt{\sum_{\mathbf{x} \sim \mu} |d_{\mathcal{Y}}(f_i(\mathbf{x}), f_j(\mathbf{x}))|^2}. \quad (7)$$

When  $\mathcal{Y}$  is a subset of  $\mathbb{R}^Y$ , the right-hand side (RHS) of (7) is equivalent to the distance between the flattened output vectors corresponding to specific input values. In the context of FD, the shared dataset  $\mathcal{D}_s$  can be used as the shared input values. Whereas the distributions  $\mu_i$  of local datasets  $\mathcal{D}_i$  differ among devices under a non-IID condition,  $\mathcal{D}_s$  is expected to have an IID by sampling appropriately from the public dataset.

Let  $\mathbf{y}_{i,\mathcal{D}_s}^t$  denote the vectorized output values of device  $i$  defined as follows:

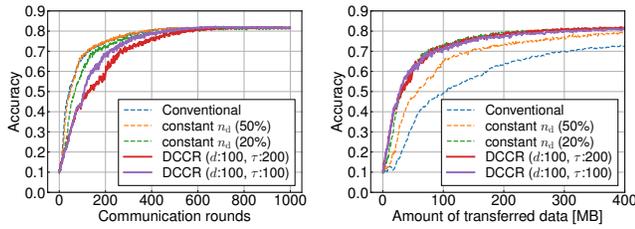
$$\mathbf{y}_{i,\mathcal{D}_s}^t := [f_i^t(\mathbf{x}_1)^\top \dots f_i^t(\mathbf{x}_{|\mathcal{D}_s|})^\top]^\top \in \mathbb{R}^{Y|\mathcal{D}_s|}, \quad (8)$$

where  $A^\top$  denotes the transpose of matrix  $A$ . As mentioned previously, the output vectors  $\mathbf{y}_{i,\mathcal{D}_s}^t$  and  $\mathbf{y}_{j,\mathcal{D}_s}^t$  can be utilized rather than the models  $f_i$  and  $f_j$  when projecting models onto a low-dimensional space with dimensionality-reduction techniques because the distance of the output vector coincides with the empirical approximation of the models (7). UMAP [19] is a distance-based dimension reduction tool that can be utilized to visualize high-dimensional data. By treating (8) as a feature vector, its dimensions can be reduced to two or three using the UMAP. Consequently, positions of the dimension-reduced vectors can then be plotted in a 2D or 3D space.

## VI. EVALUATION

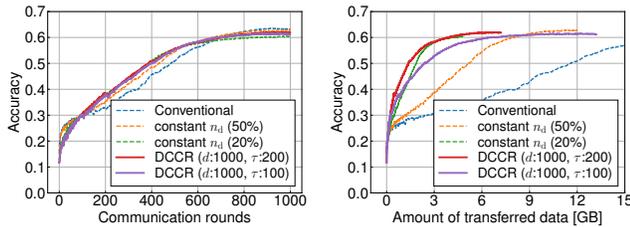
### A. Prediction accuracy with DCCR

In this section, we discuss the evaluation results of the DCCR using fashion MNIST (F-MNIST) [23] and CI-



(a) Prediction accuracy vs. communication rounds. (b) Prediction accuracy vs. communication costs.

Fig. 2: Prediction accuracy as functions of communication rounds and communication costs when learning F-MNIST. Limiting the data samples in distillation reduced the communication costs without degrading the prediction accuracy.



(a) Prediction accuracy vs. communication rounds. (b) Prediction accuracy vs. communication costs.

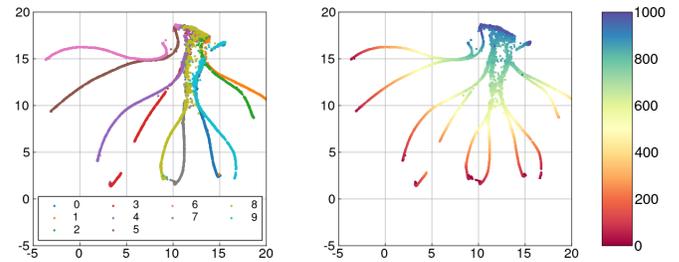
Fig. 3: Prediction accuracy as functions of communication rounds and the amount of transferred data when learning CIFAR10. Compared with the conventional method, the DCCR significantly reduced the total amount of transferred data.

FAR10 [24]. The simulation parameters are listed in Table I. We evaluated the ring lattice topology in which each device communicated with six adjacent devices. Device  $i$  possesses a non-IID local dataset including two categories labeled  $i$  and  $((i + 1) \bmod 10)$ , assuming a situation where nearby devices have similar data distributions, where  $\bmod$  denotes the modulo operator. We utilized the optimization software Optuna [25] to optimize the learning rate  $\eta$  and sharing rate  $\varepsilon$  that appear in Alg. 1.

The results obtained for the F-MNIST and CIFAR10 datasets are shown in Figs. 2 and 3, respectively. The left and right graphs in each figure show the prediction accuracy as a function of communication rounds and the amount of transferred data in the whole network, respectively. Figs. 2 and 3 show the accuracy of CMFD using different sizes  $n_d^t$  of the temporary dataset  $\tilde{\mathcal{D}}^t$ . The dashed lines, labeled as ‘‘Conven-

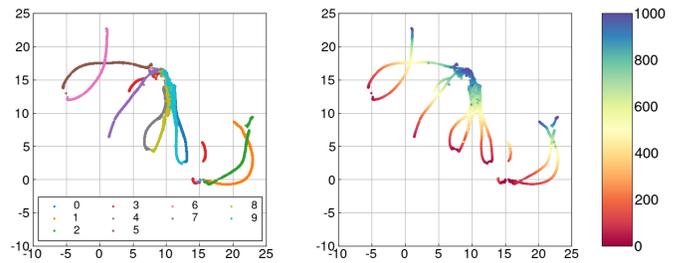
TABLE I: Simulation parameters

Parameters	Values
Num. devices	10
Network topology	Ring lattice
Num. training data per device $ \mathcal{D}_i $	1000 (F-MNIST) 4000 (CIFAR10)
Num. share data $ \mathcal{D}_s $	1000 (F-MNIST) 10000 (CIFAR10)



(a) Color code: device ID (b) Color code: communication rounds

Fig. 4: Visualized learning process of F-MNIST using DCCR with  $d = 100, \tau = 200$ . Plotted points represent the two-dimensional projection of the local models during training. In the left image, different colors represent different device models. In the right image, points are colored by the communication rounds to track the progress of learning over time. All local models converge to near positions at the end of the training.



(a) Color code: device ID (b) Color code: communication rounds

Fig. 5: Visualized learning process of F-MNIST when 20% samples were used for distillation in each communication round. Positions and colors represent the same information as Fig. 4. Some models could not converge to the same position as the others.

ditional,’’ represent the accuracy of the original CMFD that uses all data samples in  $\mathcal{D}_s$  for distillation. Additionally, the figures illustrates the performance when either 20% or 50% of samples are randomly selected from  $\mathcal{D}_s$  during each communication round; these are labeled as ‘‘constant  $n_d^t$ .’’ The figures also present the performance of DCCR with two parameter settings. Parameters  $d$  and  $\tau$  control the increase rate on  $n_d^t$  based on (6). DCCR with  $\tau = 200$  uses less samples and is more communication efficient than that with  $\tau = 100$ .

Although different strategies of  $n_d^t$  achieved similar accuracy in the end, they differed in the amount of data required to reach that accuracy. As shown in Fig. 2(a), CMFD without DCCR labeled as ‘‘Conventional’’ achieved a higher accuracy than that with DCCR at the beginning of the training because it utilizes all samples of  $\mathcal{D}_s$  for distillation, whereas distillation data are limited when using DCCR. However, DCCR has advantages in terms of communication cost. It is shown that DCCR with  $d = 100$  and  $\tau = 100$  requires approximately

84 MB data to achieve an accuracy of 70% in Fig. 2(b), whereas the conventional scheme requires 300 MB of data to achieve the same accuracy. In addition, similar characteristics can be observed in Fig. 3 showing the performance of training CIFAR10. Although a slight difference can be observed in the accuracy with a constant  $n_d^t$  and DCCR, the local models did not converge with a constant  $n_d^t$  as discussed in the following subsection.

### B. Analysis with visualized function trajectories

The visualized images of the convergence of the trained models  $f_i$  when learning F-MNIST are presented in Figs. 4 and 5. Fig. 4 shows the trajectory of the trained models when DCCR ( $n_d^t = 100, \tau = 200$ ) is applied. Fig. 5 also shows the trajectory with a constant  $n_d^t$ , where 20% samples were randomly selected in each round. The colors of the plots in Figs. 4(a) and 4(b) represent the device ID and number of communication rounds, respectively. The positions of the plotted points represent the UMAP-based two-dimensional projection of the local models during training. Since UMAP is a distance-based dimensionality-reduction technique, models that output similar values are projected to nearby positions. When a line is drawn connecting the plots of each device's model, the trajectory illustrates how the model changes gradually through gradient-descent-based training. At the beginning of the training, local models are plotted at distant positions, but they move closer to each other as training progresses.

Comparing the two schemes—DCCR shown in Fig. 4 and a scheme with constant  $n_d^t$  shown in Fig. 5—DCCR achieved superior convergence. Although both schemes required similar amounts of transferred data and attained comparable accuracy, DCCR outperformed in terms of convergence. The underlying reason is that in the distillation steps with a small-sized  $n_d^t$ , the models could not get close to each other sufficiently. In contrast, DCCR enabled the devices to utilize all samples in  $\mathcal{D}_s$  for distillation at the end of the training. Consequently, DCCR achieved better convergence compared with a scheme using constant  $n_d^t$ .

Figs. 6 and 7 show the trajectories of the trained models when learning the CIFAR10. When we projected the prediction models onto a 2D space, we resampled 1000 samples from the shared data  $\mathcal{D}_s$  to reduce the computation time without degrading the visualization precision. Furthermore, the characteristics similar to those of F-MNIST can be observed comparing Figs. 6 and 7.

## VII. CONCLUSION

We proposed a DCCR that suppresses the transferred data in the DFD. With DCCR, the devices first share a small number of distilled values with adjacent devices and then gradually increase the number of shared values to improve the convergence performance. We also developed a visualization scheme using dimensionality-reduction techniques. The results showed the trajectory of how NN models approach others in a function space.

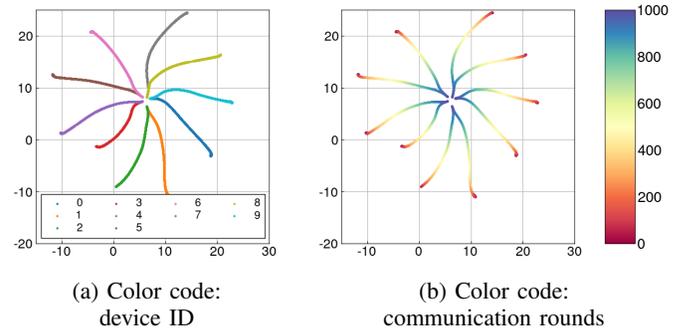


Fig. 6: Visualized learning process of CIFAR10 using DCCR with  $d = 100, \tau = 100$ . Positions and colors represent the same information as Fig. 4. All local models successfully converged to the same positions at the end of the training.

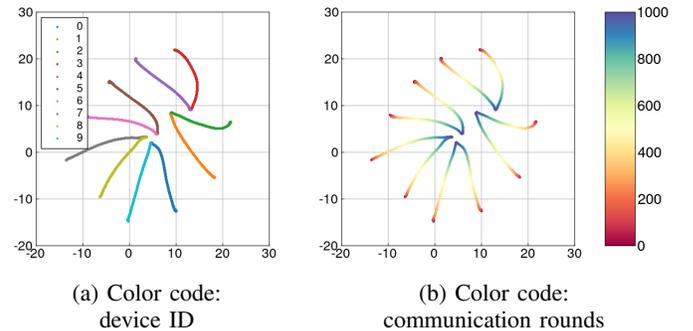


Fig. 7: Visualized learning process of CIFAR10 when 50% of samples were used for distillation in each communication round. Positions and colors represent the same information as Fig. 4. The models could not converge to the same position.

## REFERENCES

- [1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Proc. 20th Int. Conf. Artificial Intelligence and Statistics (AISTATS)*, (Fort Lauderdale, FL, USA), pp. 1273–1282, Apr. 2017.
- [2] P. Kairouz *et al.*, "Advances and open problems in federated learning," *Foundations and Trends® in Machine Learning*, vol. 14, pp. 1–210, June 2021.
- [3] H. Uddin, M. Gibson, G. A. Safdar, T. Kalsoom, N. Ramzan, M. Ur-Rehman, and M. A. Imran, "IoT for 5G/B5G applications in smart homes, smart cities, wearables and connected cars," in *2019 IEEE 24th International Workshop on Computer Aided Modeling and Design of Communication Links and Networks (CAMAD)*, pp. 1–5, Sept. 2019.
- [4] E. T. M. Beltrán, M. Q. Pérez, P. M. S. Sánchez, S. L. Bernal, G. Bovet, M. G. Pérez, G. M. Pérez, and A. H. Celdrán, "Decentralized federated learning: Fundamentals, state-of-the-art, frameworks, trends, and challenges," *arXiv preprint cs.LG, arXiv:2211.08413*, Nov. 2022.
- [5] S. Savazzi, M. Nicoli, and V. Rampa, "Federated learning with cooperating devices: A consensus approach for massive IoT networks," *IEEE Internet Things J.*, vol. 7, pp. 4641–4654, May 2020.
- [6] A. Elgabli, J. Park, A. S. Bedi, M. Bennis, and V. Aggarwal, "GADMM: Fast and communication efficient framework for distributed machine learning," *Journal of Machine Learning Research*, vol. 21, pp. 1–39, Mar. 2020.
- [7] A. Taya, T. Nishio, M. Morikura, and K. Yamamoto, "Decentralized and model-free federated learning: Consensus-based distillation in function space," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 8, pp. 799–814, Sept. 2022.

- [8] A. Lalitha, O. C. Kilinc, T. Javidi, and F. Koushanfar, "Peer-to-peer federated learning on graphs," *arXiv preprint arXiv:1901.11173*, Jan. 2019.
- [9] K. Niwa, N. Harada, G. Zhang, and W. B. Kleijn, "Edge-consensus learning: Deep learning on P2P networks with nonhomogeneous data," in *Proc. 26th ACM SIGKDD Int. Conf. Knowledge Discovery & Data Mining*, (Virtual Conference), pp. 668–678, Aug. 2020.
- [10] R. Anil, G. Pereyra, A. Passos, R. Ormandi, G. E. Dahl, and G. E. Hinton, "Large scale distributed neural network training through online distillation," *arXiv preprint arXiv:1804.03235*, Apr. 2018.
- [11] Y. Zhang, T. Xiang, T. M. Hospedales, and H. Lu, "Deep mutual learning," in *Proc. 2018 IEEE/CVF Conf. Computer Vision and Pattern Recognition (CVPR)*, pp. 4320–4328, IEEE, 2018.
- [12] Y. Jiang, S. Wang, V. Valls, B. J. Ko, W.-H. Lee, K. K. Leung, and L. Tassiulas, "Model pruning enables efficient federated learning on edge devices," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–13, 2022.
- [13] H. Zhang, J. Liu, J. Jia, Y. Zhou, H. Dai, and D. Dou, "FedDUAP: Federated learning with dynamic update and adaptive pruning using shared data on the server," *arXiv preprint arXiv:2204.11536*, 2022.
- [14] W. Xu, W. Fang, Y. Ding, M. Zou, and N. Xiong, "Accelerating federated learning for IoT in big data analytics with pruning, quantization and selective updating," *IEEE Access*, vol. 9, pp. 38457–38466, 2021.
- [15] T. Nishio and R. Yonetani, "Client selection for federated learning with heterogeneous resources in mobile edge," in *Proc. 2019 IEEE international conference on communications (ICC)*, (Shanghai, China), pp. 1–7, May 2019.
- [16] F. Sattler, A. Marban, R. Rischke, and W. Samek, "Communication-efficient federated distillation," *arXiv preprint arXiv:2012.00632*, 2020.
- [17] M. P. Uddin, Y. Xiang, X. Lu, J. Yearwood, and L. Gao, "Mutual information driven federated learning," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 7, pp. 1526–1538, 2021.
- [18] C. M. Bishop and N. M. Nasrabadi, *Pattern recognition and machine learning*, vol. 4. Springer, 2006.
- [19] L. McInnes, J. Healy, and J. Melville, "UMAP: Uniform manifold approximation and projection for dimension reduction," *arXiv preprint arXiv:1802.03426*, 2018.
- [20] L. van der Maaten and G. Hinton, "Visualizing data using t-SNE," *Journal of Machine Learning Research*, vol. 9, no. 11, pp. 2579–2605, 2008.
- [21] E. Jeong, S. Oh, H. Kim, J. Park, M. Bennis, and S.-L. Kim, "Communication-efficient on-device machine learning: Federated distillation and augmentation under non-IID private data," *arXiv preprint arXiv:1811.11479*, Nov. 2018.
- [22] T. Van Erven and P. Harremos, "Rényi divergence and Kullback-Leibler divergence," *IEEE Trans. Inf. Theory*, vol. 60, pp. 3797–3820, July 2014.
- [23] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [24] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," *Technical report, University of Toronto*, 2009.
- [25] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A next-generation hyperparameter optimization framework," in *Proc. of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '19, (Anchorage, AK, USA), pp. 2623–2631, Association for Computing Machinery, July 2019.