

**Notice:** This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

arXiv:2110.05178v1 [cs.LG] 11 Oct 2021

# Gradual Federated Learning with Simulated Annealing

Luong Trung Nguyen, Junhan Kim, and Byonghyo Shim

Information System Laboratory

Department of Electrical and Computer Engineering, Seoul National University

Email: {ltnguyen,junhankim,bshim}@islab.snu.ac.kr

## Abstract

Federated averaging (FedAvg) is a popular federated learning (FL) technique that updates the global model by averaging local models and then transmits the updated global model to devices for their local model update. One main limitation of FedAvg is that the average-based global model is not necessarily better than local models in the early stage of the training process so that FedAvg might diverge in realistic scenarios, especially when the data is non-identically distributed across devices and the number of data samples varies significantly from device to device. In this paper, we propose a new FL technique based on simulated annealing. The key idea of the proposed technique, henceforth referred to as *simulated annealing-based FL* (SAFL), is to allow a device to choose its local model when the global model is immature. Specifically, by exploiting the simulated annealing strategy, we make each device choose its local model with high probability in early iterations when the global model is immature. From extensive numerical experiments using various benchmark datasets, we demonstrate that SAFL outperforms the conventional FedAvg technique in terms of the convergence speed and the classification accuracy.

## I. INTRODUCTION

Federated learning (FL) is an emerging distributed learning technique where hundreds or thousands of devices jointly train a common machine learning (ML) model without exchanging

A part of this paper was presented at the International Conference on Acoustics, Speech, and Signal Processing (ICASSP), 2021 [1].

This work was supported by Samsung Research Funding & Incubation Center for Future Technology of Samsung Electronics under Project Number SRFC-IT1901-17.

their local dataset with the centralized server or other devices [2]–[6]. A wide range of FL applications include human face recognition, next-word prediction, resource allocation, device tracking, basestation association, cyberattack detection, to name just a few [7]–[10]. In the FL-based approach, a learning task is performed in an iterative fashion, mainly following by three steps (see Fig. 1). First, a server sets up a common ML model and then broadcasts the model to the user devices. Second, user devices train the model locally and individually using their own local datasets. Third, the server evaluates the model by aggregating the locally trained parameters sent by the devices.

The central challenge of FL is to improve the learning capability of user devices without sharing their own datasets with other devices. In fact, due to various reasons such as user privacy and limited resources (e.g., computing hardware, battery power, network capacity, bandwidth), data generated in one device cannot be transmitted to the server or other devices. One well-known approach to deal with this issue is federated averaging (FedAvg) [2]. In this technique, instead of transmitting data, each device transmits locally trained parameters (e.g., gradients or updated model parameters) to the server. The server updates the global model by averaging the local parameters and then sends the updated model back to devices for the local model update.

While FedAvg is effective in solving nonconvex problem, it has been shown that FedAvg and its variants might diverge in realistic scenarios where the data is non-identically distributed across devices (e.g., data of different languages in the next-word prediction application) and/or the number of data samples significantly varies from device to device [3], [20]. One important reason for the divergence of FedAvg is that the average-based global model is not necessarily better than locally trained models so that just relying on the global model might degrade the entire learning process [11]–[13]. To illustrate this, we consider a simple FL task whose goal is to minimize the cost function given by

$$J(\mathbf{w}; \mathcal{D}) = \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} (y_i - \mathbf{x}_i^T \mathbf{w})^2 + \|\mathbf{w}\|_1, \quad (1)$$

where  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_i$  is the training dataset and  $\|\mathbf{w}\|_1 = |w_1| + |w_2|$  is the  $\ell_1$ -norm of  $\mathbf{w}$ . For simplicity, we consider two devices with the local datasets  $\mathcal{D}_1 = \{([\frac{1}{4} \ 0]^T, -1)\}$  and  $\mathcal{D}_2 = \{([0 \ \frac{3}{2}]^T, 1)\}$ . One can easily check that the parameters  $\mathbf{w}$  minimizing the cost function  $J(\mathbf{w}; \mathcal{D})$  with respect to  $\mathcal{D}_1$  and  $\mathcal{D}_2$  are  $\mathbf{w}^{(1)} = [0 \ 0]^T$  and  $\mathbf{w}^{(2)} = [0 \ \frac{4}{9}]^T$ , respectively (see Appendix A). Using  $\mathbf{w}^{(1)}$  and  $\mathbf{w}^{(2)}$ , we obtain the average-based model  $\bar{\mathbf{w}}$  evaluated at the server:  $\bar{\mathbf{w}} =$

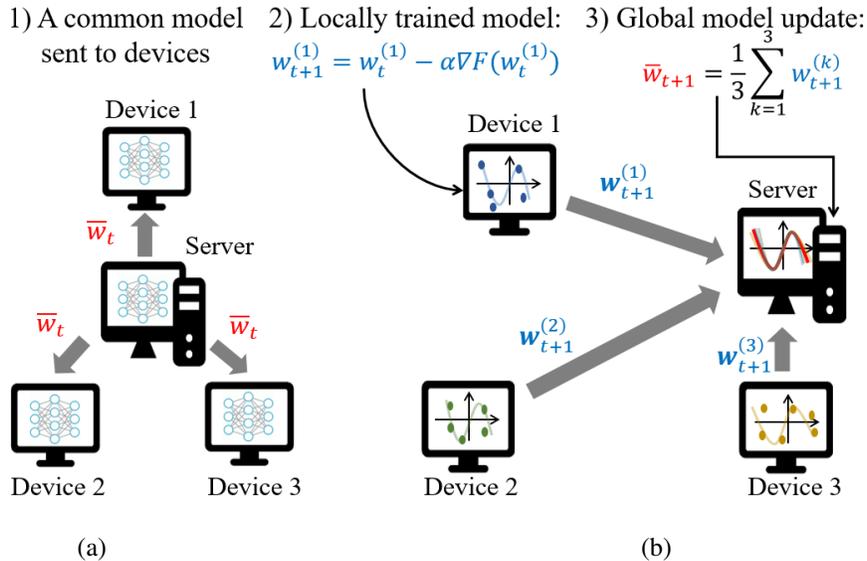


Fig. 1. Federated learning operation: (a) centralized server sends a common model to user devices and (b) each device locally trains the model using its own local data and then upload the trained network parameters to the server to globally update the model.

$\frac{1}{2}(\mathbf{w}^{(1)} + \mathbf{w}^{(2)}) = [0 \ \frac{2}{9}]^T$ . Since the optimum weight over  $\mathcal{D}_1 \cup \mathcal{D}_2$  is  $\mathbf{w}_* = [0 \ \frac{4}{9}]^T$ , we have

$$\|\mathbf{w}_* - \bar{\mathbf{w}}\|_2 = \frac{2}{9} > 0 = \|\mathbf{w}_* - \mathbf{w}^{(2)}\|_2,$$

which implies that the average value  $\bar{\mathbf{w}}$  is worse than the locally generated value  $\mathbf{w}^{(2)}$ . In this scenario, clearly, it would be better for the second device to use its own solution  $\mathbf{w}^{(2)}$  instead of the server feedback  $\bar{\mathbf{w}}$ . Simply put, the moral of the story is that collaboration might do more harm than good, especially when things are not ready.

Our intent in this paper is to put forth a simple yet effective FL strategy overcoming the problem we mentioned. Key idea of the proposed approach, referred to as the simulated annealing-based FL (SAFL), is that we encourage each device to stay with its locally trained model instead of relying on the collaborative learning model in the early stage of the learning process. When the collaborative model becomes mature and reliable after the reasonable number of iterations, we use the server-generated model to update the device. This idea can be well explained using the simulated annealing (SA) strategy. In the SA strategy, the solution space is searched by imposing perturbations on the estimates of parameters [14]–[18]. In the early stage (a.k.a., *heating stage*), the SA algorithm decides to move the system to a new (presumably perturbation) state with high

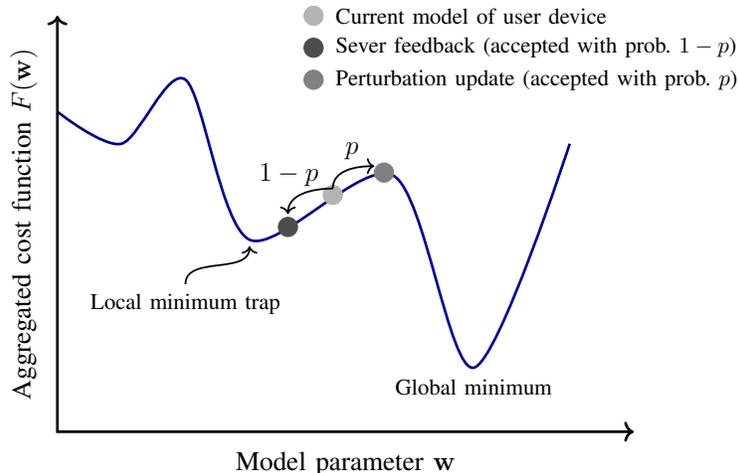


Fig. 2. The SA-based update strategy of the proposed SAFL. Here, the local model selection probability  $p$  is set to exponentially decay with respect to the training iteration.

probability, even though the new state might not be better than the current state, to avoid the chance of trapping in the local optima. In the later stage (a.k.a., *cooling stage*), the SA algorithm reduces the exploration of the perturbation space.

Inspired by the SA strategy, the proposed SAFL updates the local model of each user probabilistically. To be specific, SAFL decides whether the device keeps its own locally updated model with some modification (i.e., perturbation update) or uses the global evaluation model provided by the server (i.e., server feedback) (see Fig. 2). In the early iterations where the global model is immature, we give a favor to the locally updated model by setting the local model selection probability high. As the number of iterations increases, we gradually reduce this probability so that the device relies more on the server feedback, which helps to avoid the overfitting to the local dataset.

The main contributions of this paper are summarized as follows:

- We propose a new FL technique called SAFL (Section II) inspired by the SA technique. From extensive numerical experiments on various datasets including MNIST, Fashion-MNIST, CIFAR-10, and Google speech commands, we demonstrate that the proposed SAFL technique is effective and in fact outperforms the conventional FedAvg technique by a large margin in terms of accuracy and convergence speed (Section V). Specifically, in the MNIST dataset, SAFL converges two times faster than FedAvg and also achieves more than 50%

improvement in the classification accuracy.

- We analyze the performance of the proposed SAFL technique (Section III). Specifically, we show that under some suitable conditions, the mean squares error (MSE) of SAFL satisfies

$$E[\|\widehat{\mathbf{w}} - \mathbf{w}_*\|_2^2] \leq \xi$$

after  $\mathcal{O}(\frac{1}{\xi})$  iterations where  $\widehat{\mathbf{w}}$  is the evaluated model parameters and  $\mathbf{w}_*$  is the optimal model parameters (see Theorem III.1).

- We extend SAFL to the scenario where the performance of the average-based global model is degraded due to non-i.i.d. data and data imbalance among devices (Section IV). Our key idea is to detect biased local updates by measuring the performance gap between the global and local models. Specifically, if the performance gap is large, then we consider the local update as a biased update and do not upload it to the server. In doing so, we can exclude the biased local update in the update of the global model and prevent the performance degradation of the global model. From the numerical results, we demonstrate that the extended SAFL is effective in handling the non-i.i.d. data and reducing the number of local updates uploaded to the server (see Section V).

We briefly summarize notations used in this paper. For a vector  $\mathbf{a} \in \mathbb{R}^n$ ,  $\text{Diag}(\mathbf{a}) \in \mathbb{R}^{n \times n}$  is the diagonal matrix formed by  $\mathbf{a}$ .  $\|\mathbf{a}\|_2$  stands for the spectral norm (i.e., the largest singular value) of  $\mathbf{a}$ . The inner product of two vectors  $\mathbf{a}$  and  $\mathbf{b}$  is defined as  $\langle \mathbf{a}, \mathbf{b} \rangle = \mathbf{a}^T \mathbf{b}$ .  $\mathbf{A} \odot \mathbf{B}$  is the Hadamard product (or element-wise multiplication) of two matrices  $\mathbf{A}$  and  $\mathbf{B}$ . Given a function  $f : \mathbf{X} \in \mathbb{R}^{n_1 \times n_2} \rightarrow f(\mathbf{X}) \in \mathbb{R}$ ,  $\nabla_{\mathbf{X}} f(\mathbf{X})$  is the Euclidean gradient of  $f(\mathbf{X})$  with respect to  $\mathbf{X}$ , i.e.,  $[\nabla_{\mathbf{X}} f(\mathbf{X})]_{ij} = \frac{\partial f(\mathbf{X})}{\partial y_{ij}}$ .  $\mathbf{1} = [1 \ 1 \ \dots \ 1]^T$  is all-ones vector.

## II. PROPOSED SAFL ALGORITHM

We consider a communication system consisting of one central server and  $n$  user devices. The server generates a global model with parameters  $\mathbf{w}$  and then transmits the generated model to  $s$  selected devices ( $1 \leq s \leq n$ ). Each selected device has its own dataset  $\mathcal{D}^{(k)} = \{(\mathbf{x}_i^{(k)}, y_i^{(k)})\}_{i=1}^{m_k}$  to train the local model, where  $\mathbf{x}_i^{(k)} \in \mathbb{R}^q$  is an input data sample (e.g., image),  $y_i^{(k)}$  is the class label of  $\mathbf{x}_i^{(k)}$ , and  $m_k$  is the number of data samples in the  $k$ -th device. We consider the standard FL setting where devices cannot exchange their own datasets with other devices or the central server. In each iteration, FedAvg updates the model parameters  $\mathbf{w}$  (e.g., weights and biases)

by taking the following steps. First, using its own dataset  $\mathcal{D}^{(k)}$ , each user device updates the model parameters locally to minimize the loss function  $F(\mathbf{w}; \mathcal{D}^{(k)})$ .<sup>1</sup> For example, the update expression of the model parameters  $\mathbf{z}_t^{(k)}$  at the  $k$ -th device is

$$\mathbf{z}_t^{(k)} = \mathbf{w}_{t-1}^{(k)} - \alpha \left. \frac{\partial F(\mathbf{w}; \mathcal{D}^{(k)})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}_{t-1}^{(k)}}, \quad (2)$$

where  $\alpha$  is the learning rate and  $\mathbf{w}_{t-1}^{(k)}$  is the local model parameters after  $t-1$  iterations. Second, the server aggregates the local updates  $\mathbf{z}_t^{(k)}$  to evaluate the global model parameters. The update expression of the global evaluation model is

$$\bar{\mathbf{z}}_t = \sum_{k=1}^n \eta_k \mathbf{z}_t^{(k)}, \quad (3)$$

where  $\eta_k$  is the coefficient satisfying  $\sum_k \eta_k = 1$ .<sup>2</sup> Note that when  $s < n$ , we simply set  $\eta_k = 0$  for non-selected devices. Finally, the server transmits the globally updated parameters  $\bar{\mathbf{z}}_t$  to the selected devices to update the local models. That is, the local model parameters  $\mathbf{w}_t^{(k)}$  is updated as

$$\mathbf{w}_t^{(k)} = \begin{cases} \bar{\mathbf{z}}_t & \text{if the device receives } \bar{\mathbf{z}}_t \\ \mathbf{z}_t^{(k)} & \text{otherwise} \end{cases}. \quad (4)$$

One potential drawback of the conventional FedAvg technique is that an entire FL process can be degraded by applying the hard-decision rule in (4). This is because the global evaluation model  $\bar{\mathbf{z}}_t$  is not necessarily better than locally updated parameters  $\mathbf{z}_t^{(k)}$  in many practical scenarios. For example, in the next word prediction application, a language model is trained to predict which word comes next when the initial text fragment is given. In heterogeneous scenarios, users with different countries might use their own mother languages with different grammar and word combination rules (e.g., a subject-verb-object (SVO) rule is used in English, while a subject-object-verb (SOV) rule is used in Korean). Since the next word prediction task is performed with different language rules, the average-based model might perform much worse than the locally trained language model of a local device.

As another example, one can consider the face and object recognition problem where a classification model is trained to identify the user's face ID. The local dataset collected from

<sup>1</sup>For example, if the mean squared error (MSE) is employed as a loss function, then  $F(\mathbf{w}; \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}} \frac{1}{2} (\mathbf{x}^T \mathbf{w} - y)^2$ .

<sup>2</sup>We consider the generic setting of  $\eta_k$  which is an arbitrary value defined by user. A typical setting of  $\eta_k$  is  $\eta_k = \frac{m_k}{\sum_k m_k}$  [2].

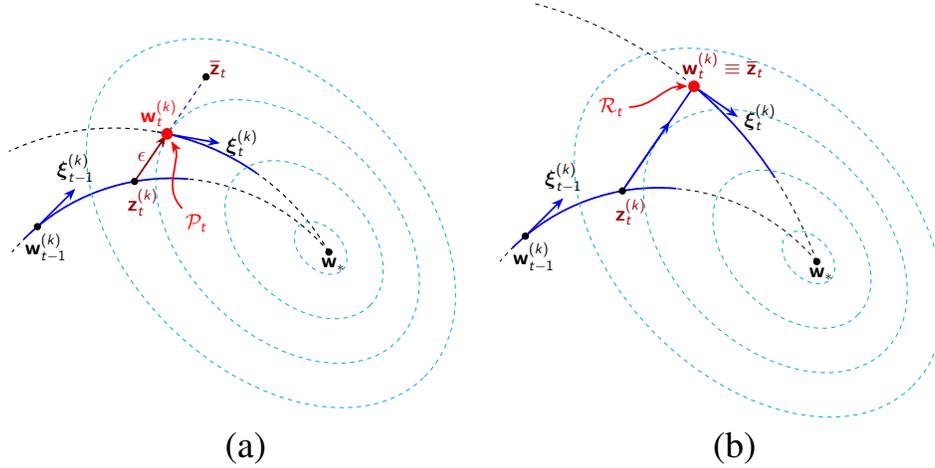


Fig. 3. In each iteration of the proposed SAFL technique, each device (a) moves to the perturbation state  $\mathcal{P}_t^{(k)}$  with probability  $p$  or (b) stays in the normal state  $\mathcal{R}_t^{(k)}$  with probability  $1 - p$ . In this figure,  $\xi_t^{(k)}$  is a descent direction.

user's personal images is often non-i.i.d. distributed across devices. Since the global model is aggregated by averaging the locally trained models, it may overfit to the local data. In this case, if the device uses the average-based model exclusively, the device might also suffer the overfitting problem, even when the good training dataset is available. Indeed, it has been shown that FedAvg can diverge in such non-i.i.d. scenario [3], [20].

Inspired by this observation, we first define a weighted sum model that incorporates  $\mathbf{z}_t^{(k)}$  and  $\bar{\mathbf{z}}_t$ . The corresponding local update model is expressed as

$$\mathbf{w}_t^{(k)} = \begin{cases} \epsilon \bar{\mathbf{z}}_t + (1 - \epsilon) \mathbf{z}_t^{(k)} & \text{if the device receives } \bar{\mathbf{z}}_t \\ \mathbf{z}_t^{(k)} & \text{otherwise} \end{cases}. \quad (5)$$

where  $\epsilon$  is the regularization parameter used to control the contribution of the global evaluation model in (5). For example, by setting  $\epsilon = 1$ , the update expression (5) is returned to the conventional FL case. Whereas, by setting  $\epsilon = 0$ , the device ignores the server feedback  $\bar{\mathbf{z}}_t$  and continues to use the locally trained model  $\mathbf{z}_t^{(k)}$ .

For the model selection, we consider a strategy inspired by the SA algorithm. In the proposed SAFL, we define the normal state  $\mathcal{R}_t$  and the perturbation state  $\mathcal{P}_t$  as  $\bar{\mathbf{z}}_t$  and  $\epsilon \bar{\mathbf{z}}_t + (1 - \epsilon) \mathbf{z}_t^{(k)}$ , respectively (see Fig. 3).  $\mathcal{P}_t$  is accepted with probability  $p = \exp(-\frac{t}{L})$  where  $L$  is a positive constant (a.k.a., the maximum temperature of SA [15]), while  $\mathcal{R}_t$  is with probability  $1 - p$ . To

TABLE I  
THE PROPOSED SAFL ALGORITHM

---

**Algorithm 1:** Proposed SAFL

---

**Input:**  $T$ : max iteration  
 $E$ : max local epoch  
 $L$ : control parameter  
 $\{\eta_k\}_k$ : weight coefficients  
 $\{\mathbf{w}_0^{(k)}\}_k$ : parameter initialization of the devices  
 $s$ : number of selected devices each round  
 $t = 1$ : initial iteration

---

**While**  $t < T$  and a stopping criterion is not met **do:**

**For** the server **do:**

**If** the server receives  $\mathbf{z}_t^{(k)}$  from the devices **then do:**

$\bar{\mathbf{z}}_t = \sum_{k=1}^n \eta_k \mathbf{z}_t^{(k)}$

Select a random set of devices  $S_t$  satisfying  $|S_t| = s$

Send  $\bar{\mathbf{z}}_t$  to  $S_t$

**End If**

**End For**

**For** device  $k \in S_t$  **in parallel do:**

**For**  $e = 1$  to  $E$  **do:**

**For** all example  $A_t^{(i_k)}$ ,  $i_k \in \{1, 2, \dots, m_k\}$  **do:**

$\mathbf{z}_t^{(k)} = \mathbf{w}_{t-1}^{(k)} - \alpha \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)})$

**If** the device receives  $\bar{\mathbf{z}}_t$  from the server **then do:**

Generate  $\mathbf{u}_t^{(k)}$  using (6)

$\mathbf{w}_t^{(k)} = \mathbf{u}_t^{(k)} \odot \bar{\mathbf{z}}_t + (\mathbf{1} - \mathbf{u}_t^{(k)}) \odot \mathbf{z}_t^{(k)}$

**Else do:**

$\mathbf{w}_t^{(k)} = \mathbf{z}_t^{(k)}$

**End If**

$t = t + 1$

**End For**

**End For**

Send  $\mathbf{z}_t^{(k)}$  to the server

**End For**

**End While**

---

**Output:**  $\hat{\mathbf{w}}_t = \sum_k \mathbf{w}_t^{(k)}$

---

be specific, let  $\mathbf{u}_t^{(k)}$  be the random vector whose  $j$ -th element  $u_j$  satisfies

$$u_j = \begin{cases} \epsilon & \text{with probability } p = \exp\left(-\frac{t}{L}\right), \\ 1 & \text{with probability } 1 - p, \end{cases} \quad (6)$$

then the local update expression (5) can be reformulated as (see Fig. 3)

$$\mathbf{w}_t^{(k)} = \begin{cases} \mathbf{u}_t^{(k)} \odot \bar{\mathbf{z}}_t + (\mathbf{1} - \mathbf{u}_t^{(k)}) \odot \mathbf{z}_t^{(k)} & \text{if receives } \bar{\mathbf{z}}_t \\ \mathbf{z}_t^{(k)} & \text{otherwise} \end{cases}. \quad (7)$$

Note that the model selection probability  $p$  decays exponentially with the number of iteration. In early iterations (i.e.,  $p$  is close to one), each device relies on its locally trained model and thus the local model would be trained mainly by the local dataset. In later iterations (i.e.,  $p$  is close to zero), the device uses the global evaluation model which is presumably more robust to the overfitting problem than the locally trained model.

We note that the server update procedure of SAFL is essentially the same as the conventional FedAvg so that various fusion models can be easily integrated to SAFL [33], [34], [37]–[41]. For example, if we integrate the inverse distance aggregation (IDA) fusion model [34] and SAFL, the coefficient  $\eta_k$  is expressed as [37]

$$\eta_k = \frac{\|\bar{\mathbf{z}}_t - \mathbf{z}_t^{(k)}\|_2^{-1}}{\sum_{k=1}^n \|\bar{\mathbf{z}}_t - \mathbf{z}_t^{(k)}\|_2^{-1}}. \quad (8)$$

In Algorithm I, we summarize the proposed SAFL algorithm.

### III. CONVERGENCE ANALYSIS OF SAFL

In this section, we analyze the convergence behavior of the proposed SAFL. For simplicity, we consider the scenario where each participating device updates its local model using the stochastic gradient descent (SGD) [19]. Let  $\delta_t$  be a user-predefined value satisfying

$$\delta_t = \begin{cases} 1 & \text{if the } k\text{-th device receives the server feedback } \bar{\mathbf{z}}_t, \\ 0 & \text{else,} \end{cases} \quad (9)$$

Then, the update expressions (2) and (7) can be reformulated as

$$\mathbf{z}_t^{(k)} = \mathbf{w}_{t-1}^{(k)} - \alpha \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)}), \quad (10)$$

$$\mathbf{w}_t^{(k)} = \delta_t \mathbf{u}_t^{(k)} \odot \bar{\mathbf{z}}_t + (\mathbf{1} - \delta_t \mathbf{u}_t^{(k)}) \odot \mathbf{z}_t^{(k)}. \quad (11)$$

where the input data  $A_t^{(i_k)} = (\mathbf{x}_{i_k}, y_{i_k}) \in \mathcal{D}^{(k)}$  is sampled identically and independently at each iteration. Here, we put no assumption on the data distribution so that our analysis results can be

applied for both i.i.d. and non-i.i.d. scenarios. Also note that  $F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)})$  is the cost function with respect to the data sample  $A_t^{(i_k)}$  and  $F_k(\mathbf{w}_{t-1}^{(k)})$  is the empirical risk function defined as

$$F_k(\mathbf{w}_{t-1}^{(k)}) = \frac{1}{|\mathcal{D}^{(k)}|} \sum_{i_k} F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)}). \quad (12)$$

Before proceeding, we summarize the assumptions used in our analysis:

**A1**  $F_k(\mathbf{w})$  is non-negative:  $F_k(\mathbf{w}) \geq 0$  and  $F_k(\mathbf{w}_*) = 0$

**A2**  $F_k(\mathbf{w})$  is a smooth convex function:  $\lambda \mathbf{I} \succeq \nabla^2 F_k(\mathbf{w}) \succeq \mu \mathbf{I}$  for  $\lambda \geq \mu \geq 0$ .

**A3** The stochastic gradient  $\nabla F_k(\mathbf{w}_t^{(k)}; A_t^{(i_k)})$  has a bounded variance:

$$\text{tr}(\text{Var}(\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)} | \mathbf{w}_{t-1}^{(k)})) \leq \sigma_k^2. \quad (13)$$

It is worth mentioning that these assumptions are used in various machine learning problems, such as linear regression, Tikhonov regularization, logistic regression, and support vector machine (SVM) [10]–[12].

Without loss of generality, we focus on the minimization problem<sup>3</sup> of the empirical risk. Hence, **A1** ensures that the objective function is to be minimized to zero. When the objective function has a nonzero lower bound, say,  $F_k(\mathbf{w}) \geq F_0$  for some constant  $F_0$ , we simply define a new objective function  $\tilde{F}_k(\mathbf{w}) = F_k(\mathbf{w}) - F_0$  and easily extend the analysis results to  $\tilde{F}_k(\mathbf{w})$ . Assumption **A2** is popularly used to guarantee a linear convergence rate of many gradient descent-based machine learning techniques [19]. Equivalently, **A2** can be expressed as [19]

**A2a**  $\nabla F_k(\mathbf{w})$  is  $\lambda$ -Lipschitz continuous:

$$\|\nabla F_k(\mathbf{w}_2) - \nabla F_k(\mathbf{w}_1)\|_2 \leq \lambda \|\mathbf{w}_2 - \mathbf{w}_1\|_2, \forall \mathbf{w}_1, \mathbf{w}_2. \quad (14)$$

**A2b**  $F_k(\mathbf{w})$  is a  $\mu$ -strongly convex function:

$$\begin{aligned} F_k(\mathbf{w}_2) &\geq F_k(\mathbf{w}_1) + \nabla F_k(\mathbf{w}_1)^T (\mathbf{w}_2 - \mathbf{w}_1) \\ &\quad + \frac{\mu}{2} \|\mathbf{w}_2 - \mathbf{w}_1\|_2^2, \forall \mathbf{w}_1, \mathbf{w}_2 \end{aligned} \quad (15)$$

for  $\lambda \geq \mu \geq 0$ .

Intuitively, **A2** ensures that there exists a quadratic lower bounds on the growth of the objective function. In our analysis, we use **A2**, together with Taylor's expansion, to build a universal upper

<sup>3</sup>The maximization problem can be converted into a minimization problem with the same solution by multiplying the objective function by  $-1$ .

bound on the MSE of the local updates  $\mathbf{w}_t^{(k)}$ . Assumption **A3** is referred to as bounded variance condition in the literature [19], which is widely used in the SGD convergence analysis [20]–[22].

In our main theorem, under **A1**, **A2**, and **A3**, we show that the proposed SAFL converges linearly<sup>4</sup> to an accurate solution.

**Theorem III.1.** *Under A1, A2, and A3, the MSE error bound of SAFL satisfies*

$$E[\|\widehat{\mathbf{w}}_t - \mathbf{w}_*\|_2^2] \leq (1 - \alpha\mu)^{2t}\zeta + \frac{\alpha}{\mu} \sum_{k=1}^n \eta_k^2 \sigma_k^2 \frac{1 - (1 - \alpha\mu)^{2q}}{1 - e^{-c}(1 - \alpha\mu)^{2q}}, \quad (16)$$

where  $\widehat{\mathbf{w}}_t = \sum_{k=1}^n \eta_k \mathbf{w}_t^{(k)}$ ,  $\zeta = \max_k E[\|\mathbf{w}_0^{(k)} - \mathbf{w}_*\|_2^2]$ ,  $q$  is the largest number of local iterations and  $c = (1 - p(1 - \epsilon^2))(\frac{1 - \alpha(2\lambda - \mu)}{1 - \alpha\mu})^2$  for some  $p$  and  $\epsilon$ , provided that  $\alpha < \frac{1}{2\lambda - \mu}$ .

**Remark III.2.** *The right-hand side of (16) consists of two terms: 1) the first term  $(1 - \alpha\mu)^{2t}\zeta$  converges linearly to zero with the iteration  $t$  and 2) the second term is a function of the learning rate  $\alpha$  and can be reduced with a small  $\alpha$ . In fact, when  $\alpha = \mathcal{O}(\frac{1}{t})$ , we can further show that SAFL converges sublinearly to the optimal solution.*

**Corollary III.3.** *Under the same conditions of Theorem III.1, if  $\alpha_t = \frac{\alpha_0}{t+1}$  for some  $\alpha_0$  satisfying  $\frac{2 - \sqrt{2}}{\mu} < \alpha_0 < \frac{2 + \sqrt{2}}{\mu}$ , then the MSE bound of SAFL satisfies*

$$E[\|\widehat{\mathbf{w}}_t - \mathbf{w}_*\|_2^2] \leq \frac{c}{t+1}, \quad (17)$$

where  $c = \max\{\frac{2\alpha_0^2(\max_k \sigma_k^2)}{2 - (2 - \mu\alpha_0)^2}, \max_k E[\|\mathbf{w}_0^{(k)} - \mathbf{w}_*\|_2^2]\}$ .

*Proof.* See Appendix E. □

*One can see that the MSE of the proposed SAFL scales in the order of  $\mathcal{O}(\frac{1}{t})$ . This MSE bound matches to the latest results of federated optimization bound [20]–[22].*

**Remark III.4.** *In Theorem III.1, the impact of the network size  $n$  on the MSE bound is captured by the factor  $\sum_{k=1}^n \eta_k^2 \sigma_k^2$ . In particular, when the local dataset has the same size (i.e.,  $\eta_k = \frac{1}{n}$ ), we*

<sup>4</sup>A sequence  $\{u_t\}_{t=1}^\infty$  is said to converge linearly to  $u_*$  if there exists a number  $\lambda \in (0, 1)$  such that  $\lim_{t \rightarrow \infty} \frac{|u_{t+1} - u_*|}{|u_t - u_*|} = \lambda$ . Also, if  $\lambda = 1$ , then the sequence is said to converge sublinearly to  $u_*$ .

have

$$\begin{aligned} E[\|\widehat{\mathbf{w}}_t - \mathbf{w}_*\|_2^2] &\leq (1 - \alpha\mu)^{2t}\zeta + \frac{\alpha\bar{\sigma}^2}{n\mu} \frac{1 - (1 - \alpha\mu)^{2q}}{1 - e^{-c}(1 - \alpha\mu)^{2q}} \\ &\stackrel{(a)}{\leq} (1 - \alpha\mu)^{2t}\zeta + \frac{\alpha\bar{\sigma}^2}{n\mu}, \end{aligned}$$

where  $\bar{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n \sigma_k^2$  and (a) is because  $\frac{1 - (1 - \alpha\mu)^{2q}}{1 - e^{-c}(1 - \alpha\mu)^{2q}} \leq 1$ . When  $t$  is large and  $\alpha$  is fixed, this MSE bound decays and converges to  $\frac{\alpha\bar{\sigma}^2}{n\mu}$ , which means that the quality of the SAFL solution improves with the number of participating devices  $n$ .

**Remark III.5.** While we use the convexity assumption **A2** to facilitate our analysis, our main result can be readily extended to the case where  $F_k(\mathbf{w})$  is not necessarily a strong convex function. For example, we consider the non-negative function  $F_k(\mathbf{w})$  satisfying

**A4**  $F_k(\mathbf{w})$  is  $\mu$ -strongly quasi-convex:

$$\langle \nabla F_k(\mathbf{w}), \mathbf{w} - \mathbf{w}_* \rangle \geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}_*\|_2^2. \quad (18)$$

**A5** The stochastic gradient  $\nabla F_k(\mathbf{w}_t^{(k)}; A_t^{(i_k)})$  has a bounded variance:

$$E[\|\nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)})\|_2^2] \leq \sigma_k^2. \quad (19)$$

Note that **A4** is weaker than **A2** since if **A2** holds true, then we have

$$\begin{aligned} \langle \nabla F_k(\mathbf{w}), \mathbf{w} - \mathbf{w}_* \rangle &\stackrel{(a)}{=} F_k(\mathbf{w}) - F_k(\mathbf{w}_*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}_*)^T \nabla^2 F_k(\boldsymbol{\xi})(\mathbf{w} - \mathbf{w}_*) \\ &\stackrel{(b)}{\geq} \frac{1}{2}(\mathbf{w} - \mathbf{w}_*)^T \nabla^2 F_k(\boldsymbol{\xi})(\mathbf{w} - \mathbf{w}_*) \\ &\stackrel{(c)}{\geq} \frac{\mu}{2} \|\mathbf{w}_* - \mathbf{w}\|_2^2, \end{aligned}$$

where  $\boldsymbol{\xi}$  is a point between  $\mathbf{w}$  and  $\mathbf{w}_*$ , (a) is due to Taylor's expansion, (b) is because  $F_k(\mathbf{w}) - F_k(\mathbf{w}_*) \geq 0$ , and (c) is because  $\nabla^2 F_k(\boldsymbol{\xi}) \succeq \mu \mathbf{I}$ . We also note that **A4** does not imply **A2**, meaning that the quasi-strong convexity does not imply the convexity of  $F_k(\mathbf{w})$  [23]. For a complete review of the functional classes satisfying this condition, see [24]. Interestingly, using **A4** instead of **A2**, one can show that the proposed SAFL still has the same convergence rate  $\mathcal{O}(\frac{1}{t})$ .

**Theorem III.6.** Under A1, A4, and A5, if  $\alpha_t = \frac{\alpha_0}{t+1}$  for some  $\alpha_0$  satisfying  $\alpha_0 > \frac{1}{\mu}$ , the MSE bound of SAFL satisfies

$$E\left[\sum_k \eta_k \|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2\right] \leq \frac{c}{t+1}, \quad (20)$$

where  $c = \max\left\{\frac{\alpha_0^2 \sum_k \eta_k \sigma_k^2}{\mu \alpha_0 - 1}, E\left[\sum_k \eta_k \|\mathbf{w}_0^{(k)} - \mathbf{w}_*\|_2^2\right]\right\}$ .

*Proof.* See Appendix F. □

**Remark III.7.** In Theorem III.1, we put no constraint on the number of local iterations  $q$  (a.k.a., the synchronization interval [22]). Therefore, the communication rounds required for  $T$  iterations is  $\mathcal{O}(Tq^{-1})$ , which is comparable to the latest results of existing distributed SGD techniques [22].

We are now ready to prove Theorem III.1.

*Proof of Theorem III.1:* In our proof, we first show that the bound of  $E[\|\widehat{\mathbf{w}}_t - \mathbf{w}_*\|_2^2]$  is expressed in terms of  $E[\|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2]$  and  $\|E[\mathbf{w}_t^{(k)} - \mathbf{w}_*]\|_2^2$ . We then build the upper bounds for each of these. That is,

$$\begin{aligned} E[\|\widehat{\mathbf{w}}_t - \mathbf{w}_*\|_2^2] &= E\left[\left\|\sum_k \eta_k \mathbf{w}_t^{(k)} - \mathbf{w}_*\right\|_2^2\right] \\ &= \sum_k \eta_k^2 E[\|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2] \\ &\quad + \sum_i \sum_{j \neq i} \eta_i \eta_j E[\langle \mathbf{w}_t^{(i)} - \mathbf{w}_*, \mathbf{w}_t^{(j)} - \mathbf{w}_* \rangle] \\ &\stackrel{(a)}{\leq} \sum_k \eta_k^2 E[\|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2] \\ &\quad + \sum_i \sum_{j \neq i} \eta_i \eta_j \|E[\mathbf{w}_t^{(i)} - \mathbf{w}_*]\|_2 \|E[\mathbf{w}_t^{(j)} - \mathbf{w}_*]\|_2 \\ &\stackrel{(b)}{=} \left(\sum_k \eta_k^2\right) \max_k E[\|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2] \\ &\quad + \left(1 - \sum_k \eta_k^2\right) \max_k \|E[\mathbf{w}_t^{(k)} - \mathbf{w}_*]\|_2^2, \end{aligned} \quad (21)$$

where (a) is from the Cauchy-Schwarz inequality and (b) is because  $\sum_i \sum_{j \neq i} \eta_i \eta_j = (\sum_k \eta_k)^2 - \sum_k \eta_k^2 = 1 - \sum_k \eta_k^2$ .

In the following lemmas, we provide the upper bounds of  $E[\|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2]$  and  $\|E[\mathbf{w}_t^{(k)} - \mathbf{w}_*]\|_2^2$ .

**Lemma III.8.** *Under the same conditions of Theorem III.1, we have*

$$E[\|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2] \leq (1 - \alpha\mu)^{2t} e^{-c\lfloor \frac{t}{q} \rfloor} \zeta_0 + \frac{\alpha^2 \sigma_k^2}{1 - (1 - \alpha\mu)^2} \frac{1 - (1 - \alpha\mu)^{2q}}{1 - e^{-c}(1 - \alpha\mu)^{2q}}, \quad (22)$$

where  $\zeta_0 = \max_i E[\|\mathbf{w}_0^{(i)} - \mathbf{w}_*\|_2^2]$ .

*Proof.* See Appendix B. □

**Lemma III.9.** *Under the same conditions of Theorem III.1, we have*

$$\|E[\mathbf{w}_t^{(k)} - \mathbf{w}_*]\|_2 \leq (1 - \alpha\mu)^t \sqrt{\zeta}, \quad (23)$$

where  $\zeta = \max_k E[\|\mathbf{w}_0^{(k)} - \mathbf{w}_*\|_2^2]$ .

*Proof.* See Appendix C. □

Finally, using (21), (22), and (23), we have

$$\begin{aligned} E[\|\widehat{\mathbf{w}}_t - \mathbf{w}_*\|_2^2] &\leq \sum_k \eta_k^2 (1 - \alpha\mu)^{2t} e^{-c\lfloor \frac{t}{q} \rfloor} \zeta + (1 - \sum_k \eta_k^2) (1 - \alpha\mu)^{2t} \zeta \\ &\quad + \sum_k \eta_k^2 \sigma_k^2 \frac{\alpha^2}{1 - (1 - \alpha\mu)^2} \frac{1 - (1 - \alpha\mu)^{2q}}{1 - e^{-c}(1 - \alpha\mu)^{2q}} \\ &= (1 - \sum_k \eta_k^2 (1 - e^{-c\lfloor \frac{t}{q} \rfloor})) (1 - \alpha\mu)^{2t} \zeta \\ &\quad + \sum_k \eta_k^2 \sigma_k^2 \frac{\alpha^2}{1 - (1 - \alpha\mu)^2} \frac{1 - (1 - \alpha\mu)^{2q}}{1 - e^{-c}(1 - \alpha\mu)^{2q}} \\ &\stackrel{(a)}{\leq} (1 - \alpha\mu)^{2t} \zeta + \frac{\sum_k \eta_k^2 \sigma_k^2 \alpha^2}{1 - (1 - \alpha\mu)^2} \frac{1 - (1 - \alpha\mu)^{2q}}{1 - e^{-c}(1 - \alpha\mu)^{2q}} \\ &\stackrel{(b)}{\leq} (1 - \alpha\mu)^{2t} \zeta + \frac{\alpha}{\mu} \sum_k \eta_k^2 \sigma_k^2 \frac{1 - (1 - \alpha\mu)^{2q}}{1 - e^{-c}(1 - \alpha\mu)^{2q}}, \end{aligned}$$

where (a) is because  $1 - \sum_k \eta_k^2 (1 - e^{-c\lfloor \frac{t}{q} \rfloor}) \leq 1$  and (b) is because  $1 - (1 - \alpha\mu)^2 = \alpha\mu(2 - \alpha\mu) \geq \alpha\mu$ , which establishes Theorem III.1. □

#### IV. EXTENDED SAFL FOR THE OVERFITTING PROBLEM

In many practical scenarios, the FL performance can be degraded due to various reasons such as biased user data, training failures, model poisoning attack, and adversarial attacks [7]–[10], [20], [25]. For example, when a user device trains its local model using non-representative data

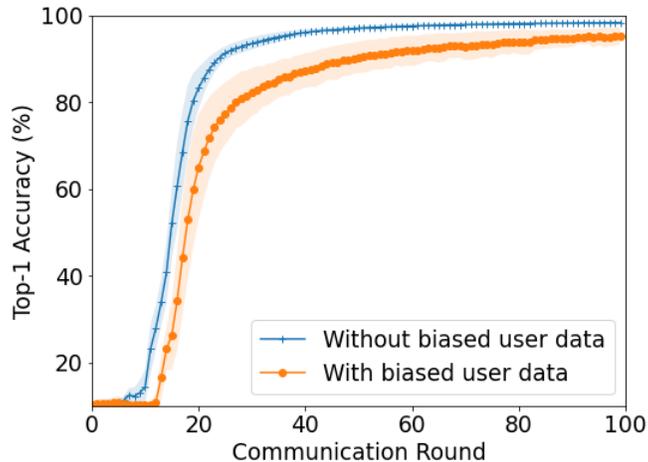


Fig. 4. The performance of the FL network with and without biased user data.

(i.e., certain elements in the dataset are more heavily weighted and represented than others), then the local update of the device might cause a model overfitting problem, resulting in the degradation of the entire FL network performance. To illustrate this behavior, we consider a FL network performing the MNIST classification (see Section V for the detailed setting of the FL network). Depending on the number of digit labels in the local datasets, devices can be classified into two groups: 1) a group with local dataset containing multiple digits (say, 1, 2, 5, 7, and 9) and 2) a group with dataset containing only one digit (say, 2). Due to the data bias, devices in the second group can only learn features of one digit and, as a result, locally trained models might fail to predict other digits (1, 5, 7, and 9). In fact, when the locally trained models of the second group are overfitted to the biased dataset, there would be a performance degradation in the global evaluation model. In our example, if the server uses the local updates of the second group in the update of the global model, then the accuracy of the global model is degraded significantly (see Fig. 4).

In the above example, since the data distributions of devices are known a priori, we can prevent the performance degradation of the global model by excluding local updates of the second group in the update of the global model. In general, however, it is very difficult for the server to exclude those biased local updates since the local datasets are not revealed to the server due to the privacy of the user data. Instead of making the server to exclude biased local updates,

---

**Algorithm 2:** Extended SAFL
 

---

**Input:**  $T$ : max iteration

 $E$ : max local epoch

 $L$ : maximum temperature

 $\{\eta_k\}_k$ : weight coefficients

 $\{\mathbf{w}_0^{(k)}\}_k$ : parameter initialization of the devices

 $s$ : number of selected devices each round

 $t = 1$ : initial iteration

 $\{q_0^{(k)}\}_k = 1$ : initial probability of the local update

---

**While**  $t < T$  and a stopping criterion is not met **do:**
**For** the server **do:**
**If** the server receives  $\mathbf{z}_t^{(k)}$  from the devices **then do:**

$$\bar{\mathbf{z}}_t = \sum_{k=1}^n \eta_k \mathbf{z}_t^{(k)}$$

 Select a random set of devices  $S_t$  satisfying  $|S_t| = s$ 

 Send  $\bar{\mathbf{z}}_t$  to  $S_t$ 
**End If**
**End For**
**For** device  $k \in S_t$  **in parallel do:**
**For**  $e = 1$  to  $E$  **do:**
**For** all example  $A_t^{(i_k)}$ ,  $i_k \in \{1, 2, \dots, m_k\}$  **do:**

$$\mathbf{z}_t^{(k)} = \mathbf{w}_{t-1}^{(k)} - \alpha \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)})$$

**If** the device receives  $\bar{\mathbf{z}}_t$  from the server **then do:**

 Generate  $\mathbf{u}_t^{(k)}$  using (6)

$$\mathbf{w}_t^{(k)} = \mathbf{u}_t^{(k)} \odot \bar{\mathbf{z}}_t + (\mathbf{1} - \mathbf{u}_t^{(k)}) \odot \mathbf{z}_t^{(k)}$$

 Compute  $q_t^{(k)}$  using (25)

**Else do:**

$$\mathbf{w}_t^{(k)} = \mathbf{z}_t^{(k)}$$

$$q_t^{(k)} = q_{t-1}^{(k)}$$

**End If**
 $t = t + 1$ 
**End For**
**End For**

 Send  $\mathbf{z}_t^{(k)}$  to the server with probability  $q_t^{(k)}$ 
**End For**
**End While**


---

**Output:**  $\hat{\mathbf{w}}_t = \sum_k \mathbf{w}_t^{(k)}$ 


---

we modify SAFL such that each device can decide whether to upload its local update to the

server or not. This decision is done by measuring the performance gap between global and local models. Let  $h(\bar{\mathbf{z}}_{t-1})$  and  $h(\mathbf{z}_t^{(k)})$  be the accuracies of the global evaluation model and the local model in the  $k$ -th device, respectively, then the performance gap  $\Delta_t^{(k)}$  between the global and local models is defined as

$$\Delta_t^{(k)} = \frac{|h(\bar{\mathbf{z}}_{t-1}) - h(\mathbf{z}_t^{(k)})|}{h(\bar{\mathbf{z}}_{t-1}) + h(\mathbf{z}_t^{(k)}) + \epsilon}, \quad (24)$$

where  $\epsilon$  is a small constant to avoid division by zero (e.g.,  $\epsilon = 10^{-6}$ ). If  $\Delta_t^{(k)}$  is large, then we consider the local update as a biased update and do not upload the local update to the server. To do so, we set the probability  $q_t^{(k)}$  that the local update  $\mathbf{z}_t^{(k)}$  is uploaded to the server as<sup>5</sup>

$$q_t^{(k)} = \exp\left(-\frac{\Delta_t^{(k)}}{\nu}\right), \quad (25)$$

where  $\nu$  is a regularization parameter. Since the probability  $q_t^{(k)}$  decays exponentially with the performance gap  $\Delta_t^{(k)}$ , if  $\Delta_t^{(k)}$  is large, then it is highly likely that the device does not send its local update  $\mathbf{z}_t^{(k)}$  to the server. By excluding the biased local update  $\mathbf{z}_t^{(k)}$  in the update of the global model, we can prevent the performance degradation of the global model.

While the communication cost of SAFL is the same as that of FedAvg, the extended SAFL can reduce the number of local updates uploaded to the server. Let  $X$  be the total local updates of  $n$  devices after  $T$  communication rounds and let  $X_t^{(k)}$  be the random variable indicating whether the  $k$ -th device sends the local update to the server, i.e.,  $P(X_t^{(k)} = 1) = q_t^{(k)}$  and  $P(X_t^{(k)} = 0) = 1 - q_t^{(k)}$ . Then, we have  $X = \sum_{t=1}^T \sum_{k=1}^n X_t^{(k)}$  and thus

$$E[X] = \sum_{t=1}^T \sum_{k=1}^n \exp\left(-\frac{\Delta_t^{(k)}}{\nu}\right) \leq \sum_{t=1}^T \sum_{k=1}^n 1 = nT, \quad (26)$$

where  $nT$  is the total local updates of FedAvg.

In Algorithm II, we summarize the extended SAFL algorithm.

## V. SIMULATION

In this section, we investigate the empirical performance of the proposed SAFL on various benchmark datasets, which has been popularly used in the FL evaluation. We first summarize the datasets used in our experiments as follows:

<sup>5</sup>The choice of the exponential decay is based on our empirical experiences.

TABLE II  
DEEP NEURAL NETWORKS

LeNet-5		Light VGGNet	
Layer	Filter Stride	Layer	Filter Stride
conv5-6	1	conv3-64	1
avg-pool-2	2	conv3-128	1
conv5-16	1	max-pool-2	2
avg-pool-2	-	conv3-128	1
FC-120	-	max-pool-2	2
FC-84	-	conv3-128	1
FC-10	-	max-pool-2	2
softmax	-	conv3-128	1
-	-	max-pool-2	2
-	-	global-avg-pool	-
-	-	conv1-10	1
-	-	softmax	-
0.04M params		1.76M params	

- *MNIST* [26]: a dataset consisting of 70,000 images of handwritten digits between 0 and 9. All the images are divided into two groups: 60,000 images for the training set and 10,000 images for the test set<sup>6</sup>.
- *Fashion-MNIST* [27]: a dataset containing 70,000 grayscale images of clothing (e.g., sneakers, shirts, shoes, and bags). These images are classified into 10 categories. The dataset is divided into two sets: the training set of 60,000 images and the test set of 10,000 images.
- *CIFAR10* [28]: a dataset of color images popularly used in image classification. It consists of 60,000 images of  $32 \times 32$  pixels from 10 categories: airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck. The dataset is divided into two: training set of 50,000 images and test set of 10,000 images.
- *Google speech commands dataset* [29]: a dataset popularly used in speech recognition tasks. It consists of 65,000 utterances of 30 short words. To pre-process the GSC dataset, we compute the first 13 mel-frequency cepstral coefficients (MFCC) of a speech signal using 80 filterbanks. To be specific, we first perform a 1024-point short-time Fourier transform

<sup>6</sup>The training and the test sets are split by the command `tf.keras.datasets.mnist.load_data()` in Tensorflow.

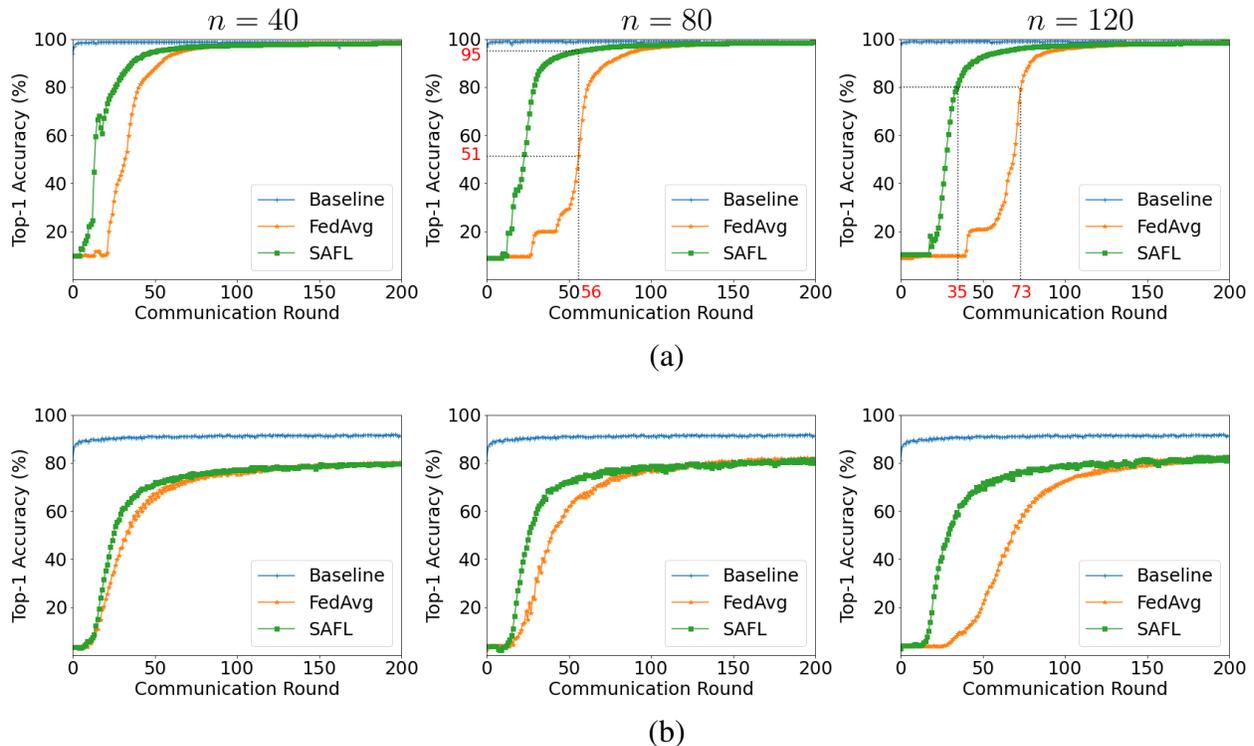


Fig. 5. Test accuracy of SAFL on different datasets: (a) MNIST and (b) GSC.

TABLE III  
TEST ACCURACY IN THE 50-TH COMMUNICATION ROUND.

Dataset	Baseline		AvgFed		SAFL	
	Top-1 Acc.	Top-5 Acc.	Top-1 Acc.	Top-5 Acc.	Top-1 Acc.	Top-5 Acc.
MNIST	99%	99%	30%	90%	94%	99%
FMNIST	92%	99%	58%	98%	75%	99%
CIFAR10	77%	98%	21%	75%	35%	88%
GSC	91%	98%	62%	91%	72%	92%

(STFT) with frames of 64ms and 75% overlap (at 16kHz sampling frequency) and then compute the power spectrum and MFCC.

In our experiments, each training set is partitioned into  $n$  subsets  $(\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_n)$  of user devices, and each of which is the local dataset in each device. We consider the heterogeneous

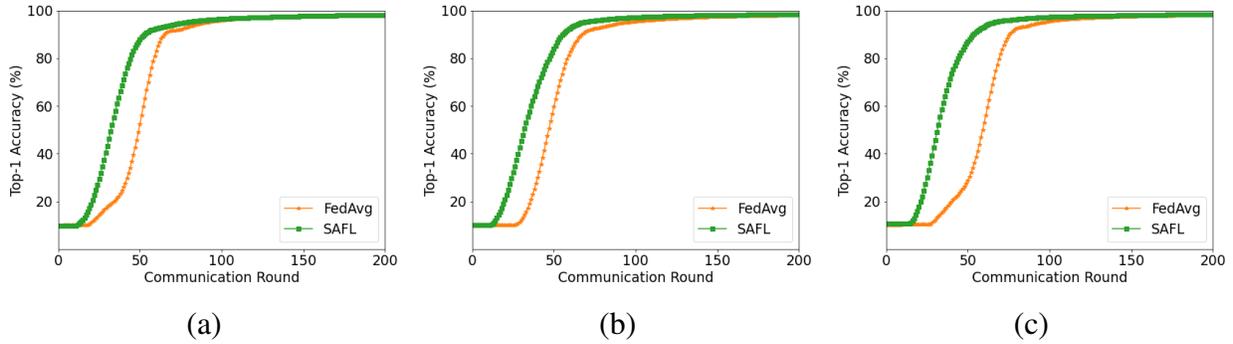


Fig. 6. Performance of SAFL for different fraction  $R$  of selected devices: (a)  $R = 0.3$ , (b)  $R = 0.5$ , and (c)  $R = 0.7$ .

TABLE IV  
EFFECT OF THE PARAMETERS  $\epsilon$  AND  $L$  IN THE PROPOSED SAFL.

$\epsilon$	Baseline		FedAvg		SAFL $L = 10$		SAFL $L = 80$	
	Test Cost	Test Acc.	Test Cost	Test Acc.	Test Cost	Test Acc.	Test Cost	Test Acc.
0.3	0.075	98.8%	0.476	88.2%	0.322	91.2%	<b>0.186</b>	<b>95.3%</b>
0.5	0.075	98.8%	0.476	88.2%	<b>0.206</b>	<b>94.7%</b>	0.204	95.1%
0.7	0.075	98.8%	0.476	88.2%	0.335	91.5%	0.217	94.6%

TABLE V  
SAFL ACCURACY FOR DIFFERENT BATCH SIZES.

Dataset size	FL Technique	Batch size			
		50	100	150	200
Half of dataset	FedAvg	97.59%	97.15%	91.53%	93.03%
	SAFL	<b>97.99%</b>	<b>97.55%</b>	<b>97.25%</b>	<b>94.01%</b>
All of dataset	FedAvg	98.44%	97.42%	97.01%	95.74%
	SAFL	<b>98.47%</b>	<b>98.12%</b>	<b>97.67%</b>	<b>97.08%</b>

scenarios where  $\mathcal{D}_k$  consists of just a few number of class labels (not all the labels) whose sizes are all different ( $m_k = |\mathcal{D}_k|$ ). To be specific, we first set  $m_k = \max(\lfloor x_k \rfloor, 1)$  where

TABLE VI  
TOTAL LOCAL UPDATES UPLOADED TO THE SERVER.

Iter	FedAvg		SAFL		Extended SAFL	
	Total updates	Test Acc.	Total updates	Test Acc.	Total updates	Test Acc.
60	3050	52.22%	3050	93.72%	2023	94.83%
70	3550	78.91%	3550	95.35%	2448	95.87%
80	4050	92.32%	4050	96.22%	2883	96.51%

$x_k \sim N(\bar{m}, \sigma^2)$  is a normal random variable with mean  $\bar{m}$  and variance  $\sigma^2$  and  $\lfloor x_k \rfloor$  is the integer satisfying  $\lfloor x_k \rfloor \leq x_k < \lfloor x_k \rfloor + 1$ . Then we select a number of digits (e.g., at most 7 digits) for each device at random and choose  $m_k$  samples randomly from the training subset containing only the selected digit labels.

For the MNIST classification, we use LeNet-5, a CNN model consisting of two sets of convolutional and pooling layers, followed by two fully-connected layers and the softmax classifier [30]. For the fashion-MNIST, CIFAR10, and GSC classifications, we use the VGGNet, a CNN model using only  $3 \times 3$  convolutional kernels [31]. The parameter settings of the CNN architectures are shown in Table II. As a loss function in the training process, we use the cross-entropy:

$$H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^{10} (y_i \ln \hat{y}_i + (1 - y_i) \ln(1 - \hat{y}_i)), \quad (27)$$

where  $\hat{\mathbf{y}} = [\hat{y}_1 \ \dots \ \hat{y}_{10}]^T$  is the predicted softmax output and  $\mathbf{y} = [y_1 \ \dots \ y_{10}]^T$  is the one-hot vector of the true label. For all experiments, we set the learning rate  $\alpha$  to a fixed constant ( $\alpha = 0.002$ ) and set the number of local epochs to  $E = 3$ . We initialize the local model of each device with a different random seed.

We first evaluate the test accuracy of SAFL for different network size ( $n = 40, 80,$  and  $120$ ). In this experiment, we set the parameters  $\epsilon = 0.3$ ,  $L = 80$ ,  $\bar{m} = 600$ , and  $\sigma^2 = 100$ . In Fig. 5, we plot the test accuracy of SAFL and the conventional FedAvg as a function of the communication round. The baseline is the centralized machine learning using the whole dataset. From the results, we observe that the accuracy of all the FL algorithms improves after a sufficient communication rounds (e.g., 100 rounds) and the performance of all FL algorithms eventually converges to the

accuracy of the centralized learning technique. In particular, the proposed SAFL outperforms the standard FL technique by a large margin. For example, for  $n = 80$ , SAFL achieves the test accuracy of 95% at the 56-th communication round, resulting in an accuracy improvement of more than 50% (see Fig. 5a). We also observe that the proposed SAFL converges faster than the standard FL technique. For example, when  $n = 120$ , SAFL achieves the accuracy of 80% in 35 communication rounds, while the standard FL technique requires more than 70 rounds to achieve the same level of accuracy (see Fig. 5a). Similar results can be observed from the GSC dataset (see Fig. 5b). In Table. III, we show the top-1 and top-5 accuracy of SAFL in the early stage of SA for  $n = 80$ . From these experiments, we observe that SAFL outperforms the conventional approaches, resulting in an 17% improvement of the top-1 accuracy on the FMNIST dataset.

We next examine the impact of the hyperparameters  $\epsilon$  and  $L$  on the performance of SAFL. In this MNIST experiment, we set  $n = 50$  and run simulations for different values  $\epsilon = 0.3, 0.5$ , and  $0.7$ . In Table. IV, we show the test cost and the test accuracy evaluated at the 50-th communication round. The best performance of SAFL is highlighted with bold digits. For example, when  $\epsilon = 0.3$  and  $L = 80$ , SAFL achieves the smallest MSE (i.e.,  $\text{MSE} = 0.186$ ) and the best accuracy (i.e., 95.3%).

We test the performance of SAFL for different training batch sizes ( $B = 50, 100, 150$ , and  $200$ ). For all the MNIST experiments, we set the parameters  $n = 100$ ,  $\epsilon = 0.3$ ,  $L = 80$ ,  $\bar{m} = 600$ ,  $\sigma^2 = 100$ . The MNIST accuracy are tested after  $T = 100$  communication rounds. From the results, we observe that the small and moderate batch size can be used to enhance the accuracy of the FL networks, especially when the data size is reduced by half. For example, the batch size  $B = 50$  gives more than 97.99% SAFL accuracy while the batch size  $B = 200$  results in less than 97.08% accuracy (see Table. V).

We also test the accuracy of SAFL for different fraction of selected devices ( $R = 0.3, 0.5$ , and  $0.7$ ). From the results, we observe that SAFL outperforms FedAvg, resulting in more than 50% improvement of the test accuracy after 50 communication rounds when  $R = 0.7$  (see Fig. 6).

Next we evaluate the test accuracy of the extended SAFL as a function of the total local updates uploaded to the server. Here, we set  $n = 100$  local devices and count the total local updates in different iterations ( $T = 60, 70$ , and  $80$ ). We run 100 trials and compute the mean values (see Table VI). From the results, we observe that SAFL has the same communication cost

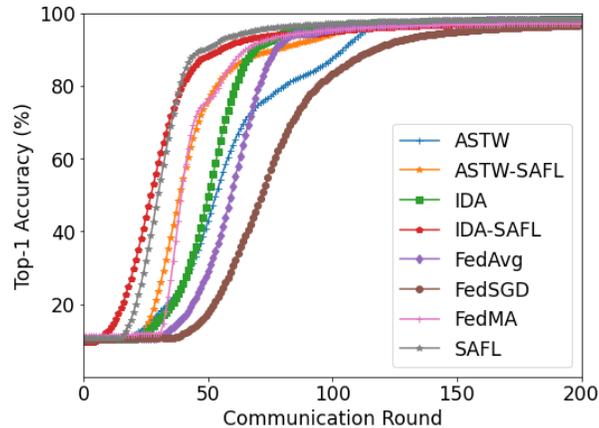


Fig. 7. The learning performance of the FL techniques.

as FedAvg. While the accuracy of the extended SAFL is comparable to the SAFL accuracy, the extended SAFL significantly reduces the number of the local updates uploaded by the devices, resulting in more than 30% reduction of the local updates (see Table VI).

Finally, we compare the performance of SAFL with the state-of-the-art FL techniques including the temporally weighted aggregation asynchronous (ASTW) [33], IDA [34], FedAvg, FedSGD [35], and FedMA [36]. We also test the combined algorithms: ASTW-SAFL and IDA-SAFL which are combined version of SAFL and ASTW/IDA fusion models [37]. From the results, we observe that SAFL outperforms FedSGD and FedAvg by a large margin, resulting in more than 50% improvement of the test accuracy after 50 communication rounds. The performance of SAFL is comparable to that of FedMA. We also observe that the combination of SAFL and state-of-the-art data fusion model can boost up the learning accuracy significantly. For example, IDA-SAFL can achieve more than 80% accuracy after 50 communication rounds, resulting in more than 30% improvement of the test accuracy over the conventional IDA.

## VI. CONCLUSION

In this paper, we proposed a FL technique that greatly improves the accuracy and convergence speed of FL. Motivated by the observation that the average-based global model is not necessarily better than local models, the proposed SAFL technique allows each device to choose its own model instead of the global model in the early stage of FL. From the convergence analysis,

we showed that SAFL sublinearly converges to the optimal solution under suitable conditions. Also, from the numerical experiments based on various benchmark datasets, we demonstrated that SAFL outperforms the conventional FL technique in terms of the convergence speed and the classification accuracy. In this work, we restricted our attention to the single-task learning scenario. Our future work will be directed toward the extension to the multi-tasking scenario [32].

## APPENDIX A

### PROOF OF $\mathbf{w}^{(1)}$ , $\mathbf{w}^{(2)}$ , AND $\mathbf{w}_*$

*Proof.* We first find the solution  $\mathbf{w}^{(1)}$ . Let  $(\mathbf{x}_1, y_1) = (\left[ \frac{1}{4} \ 0 \right]^T, -1)$  and  $\mathbf{w} = \left[ w_1 \ w_2 \right]^T$ . Then, we have

$$\begin{aligned} \mathbf{w}^{(1)} &= \arg \min_{\mathbf{w}} J(\mathbf{w}, \mathcal{D}_1) \\ &= \arg \min_{\mathbf{w}} (y_1 - \mathbf{x}_1^T \mathbf{w})^2 + \|\mathbf{w}\|_1 \\ &= \arg \min_{\mathbf{w}} \left(-1 - \frac{1}{4}w_1\right)^2 + |w_1| + |w_2| \\ &= \begin{bmatrix} \arg \min_{w_1} \left(-1 - \frac{1}{4}w_1\right)^2 + |w_1| \\ \arg \min_{w_2} |w_2| \end{bmatrix} \\ &\stackrel{(a)}{=} \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \end{aligned}$$

where (a) is because  $\left(-1 - \frac{1}{4}w_1\right)^2 + |w_1| = 1 + \frac{1}{16}w_1^2 + \frac{1}{2}w_1 + |w_1| \geq 1 + \frac{1}{16}w_1^2 \geq 1$  and the equality holds if and only if  $w_1 = 0$ . Similarly, we can find out the solutions  $\mathbf{w}^{(2)} = \left[ 0 \ \frac{4}{9} \right]^T$  and  $\mathbf{w}_* = \left[ 0 \ \frac{4}{9} \right]^T$ , which is the desired results. □

## APPENDIX B

### PROOF OF LEMMA III.8

*Proof.* In this proof, we first show a recursive inequality of the MSE and then build the upper bound of the MSE.

Let  $\mathbf{e}_t^{(k)} = \mathbf{w}_t^{(k)} - \mathbf{w}_*$  and  $\Delta \mathbf{z}_t^{(k)} = \mathbf{z}_t^{(k)} - \bar{\mathbf{z}}_t$ , then from (7), and (10), we have

$$\begin{aligned}
\mathbf{e}_t^{(k)} &= \mathbf{w}_t^{(k)} - \mathbf{w}_* \\
&= \delta_t \mathbf{u}_t^{(k)} \odot \bar{\mathbf{z}}_t + (\mathbf{1} - \delta_t \mathbf{u}_t^{(k)}) \odot \mathbf{z}_t^{(k)} - \mathbf{w}_* \\
&= \mathbf{z}_t^{(k)} - \mathbf{w}_* - \delta_t \mathbf{u}_t^{(k)} \odot \Delta \mathbf{z}_t^{(k)} \\
&= (\mathbf{w}_{t-1}^{(k)} - \mathbf{w}_*) - \delta_t \mathbf{u}_t^{(k)} \odot \Delta \mathbf{z}_t^{(k)} - \alpha \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)}) \\
&= \mathbf{e}_{t-1}^{(k)} - \delta_t \mathbf{u}_t^{(k)} \odot \Delta \mathbf{z}_t^{(k)} - \alpha (\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_*)) + \alpha (\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)})).
\end{aligned}$$

Applying Taylor's expansion yields

$$\begin{aligned}
\mathbf{e}_t^{(k)} &= \mathbf{e}_{t-1}^{(k)} - \delta_t \mathbf{u}_t^{(k)} \odot \Delta \mathbf{z}_t^{(k)} - \alpha \nabla^2 F(\boldsymbol{\xi}_{t-1}^{(k)}) \mathbf{e}_{t-1}^{(k)} + \alpha (\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)})) \\
&= (\mathbf{I} - \alpha \nabla^2 F_k(\boldsymbol{\xi}_{t-1}^{(k)})) \mathbf{e}_{t-1}^{(k)} - \delta_t \text{Diag}(\mathbf{u}_t^{(k)}) \Delta \mathbf{z}_t^{(k)} + \alpha (\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)})),
\end{aligned}$$

where  $\boldsymbol{\xi}_{t-1}^{(k)}$  is a point in the line segment of two endpoints  $\mathbf{w}_{t-1}^{(k)}$  and  $\mathbf{w}_*$ .

Taking the conditional variance of  $\mathbf{e}_t^{(k)}$ , we have

$$\begin{aligned}
E[\|\mathbf{e}_t^{(k)}\|_2^2 | \mathbf{w}_{t-1}^{(k)}, \Delta \mathbf{z}_t^{(k)}, \mathbf{u}_t^{(k)}] &= \|E[\mathbf{e}_t^{(k)} | \mathbf{w}_{t-1}^{(k)}, \Delta \mathbf{z}_t^{(k)}, \mathbf{u}_t^{(k)}]\|_2^2 + \text{tr}(\text{Var}(\mathbf{e}_t^{(k)} | \mathbf{w}_{t-1}^{(k)}, \Delta \mathbf{z}_t^{(k)}, \mathbf{u}_t^{(k)})) \\
&= \|\mathbf{A} \mathbf{e}_{t-1}^{(k)} - \delta_t \text{Diag}(\mathbf{u}_t^{(k)}) \Delta \mathbf{z}_t^{(k)}\|_2^2 + \alpha^2 \text{tr}(\text{Var}(\nabla F_k(\mathbf{w}_{t-1}^{(k)}) \\
&\quad - \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)}) | \mathbf{w}_{t-1}^{(k)}, \Delta \mathbf{z}_t^{(k)}, \mathbf{u}_t^{(k)})),
\end{aligned}$$

where  $\mathbf{A} = \mathbf{I} - \alpha \nabla^2 F_k(\boldsymbol{\xi}_{t-1}^{(k)})$ .

By Assumption **A3**, the stochastic gradient has a bounded variance:

$$\text{tr}(\text{Var}(\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(i_k)}) | \mathbf{w}_{t-1}^{(k)})) \leq \sigma_k^2. \quad (28)$$

Thus, we have

$$\begin{aligned}
E[\|\mathbf{e}_t^{(k)}\|_2^2 | \mathbf{w}_{t-1}^{(k)}, \Delta \mathbf{z}_t^{(k)}, \mathbf{u}_t^{(k)}] &\leq \|\mathbf{A} \mathbf{e}_{t-1}^{(k)} - \delta_t \text{Diag}(\mathbf{u}_t^{(k)}) \Delta \mathbf{z}_t^{(k)}\|_2^2 + \alpha^2 \sigma_k^2 \\
&= \|\mathbf{A} \mathbf{e}_{t-1}^{(k)}\|_2^2 + \delta_t \|\text{Diag}(\mathbf{u}_t^{(k)}) \Delta \mathbf{z}_t^{(k)}\|_2^2 + \alpha^2 \sigma_k^2 \\
&\quad - 2\delta_t (\Delta \mathbf{z}_t^{(k)})^T \text{Diag}(\mathbf{u}_t^{(k)}) \mathbf{A} \mathbf{e}_{t-1}^{(k)}.
\end{aligned}$$

Taking expectations of the last inequality, and noting the law of total expectation ( $E[X] =$

$E[E[X|Y]]$ , we have

$$\begin{aligned} E[\|\mathbf{e}_t^{(k)}\|_2^2] &= E[E[\|\mathbf{e}_t^{(k)}\|_2^2 | \mathbf{w}_{t-1}^{(k)}, \Delta \mathbf{z}_t^{(k)}, \mathbf{u}_t^{(k)}]] \\ &\leq E[\|\mathbf{A}\mathbf{e}_{t-1}^{(k)}\|_2^2] + \delta_t E[\|\text{Diag}(\mathbf{u}_t^{(k)})\Delta \mathbf{z}_t^{(k)}\|_2^2] + \alpha^2 \sigma_k^2 \\ &\quad - 2\delta_t E[(\Delta \mathbf{z}_t^{(k)})^T \text{Diag}(\mathbf{u}_t^{(k)})\mathbf{A}\mathbf{e}_{t-1}^{(k)}]. \end{aligned}$$

By Assumption **A2**, the positive definite matrix  $\mathbf{A}$  satisfies  $\|\mathbf{A}\mathbf{e}_{t-1}^{(k)}\|_2^2 \leq (1 - \alpha\mu)^2 \|\mathbf{e}_{t-1}^{(k)}\|_2^2$ . It follows that

$$\begin{aligned} E[\|\mathbf{e}_t^{(k)}\|_2^2] &\leq (1 - \alpha\mu)^2 E[\|\mathbf{e}_{t-1}^{(k)}\|_2^2] + \alpha^2 \sigma_k^2 \\ &\quad + \delta_t E[\|\text{Diag}(\mathbf{u}_t^{(k)})\Delta \mathbf{z}_t^{(k)}\|_2^2] - 2\delta_t E[(\Delta \mathbf{z}_t^{(k)})^T \text{Diag}(\mathbf{u}_t^{(k)})\mathbf{A}\mathbf{e}_{t-1}^{(k)}]. \end{aligned} \quad (29)$$

Using the law of total expectation, one can easily check that

$$\begin{aligned} E[\|\text{Diag}(\mathbf{u}_t^{(k)})\Delta \mathbf{z}_t^{(k)}\|_2^2] &= E[E[\|\text{Diag}(\mathbf{u}_t^{(k)})\Delta \mathbf{z}_t^{(k)}\|_2^2 | \mathbf{w}_{t-1}^{(k)}, \Delta \mathbf{z}_t^{(k)}]] \\ &= (1 - p(1 - \epsilon^2)) E[\|\Delta \mathbf{z}_t^{(k)}\|_2^2] \end{aligned} \quad (30)$$

$$\begin{aligned} E[(\Delta \mathbf{z}_t^{(k)})^T \text{Diag}(\mathbf{u}_t^{(k)})\mathbf{A}\mathbf{e}_{t-1}^{(k)}] &= E[E[(\Delta \mathbf{z}_t^{(k)})^T \text{Diag}(\mathbf{u}_t^{(k)})\mathbf{A}\mathbf{e}_{t-1}^{(k)} | \mathbf{w}_{t-1}^{(k)}, \Delta \mathbf{z}_t^{(k)}]] \\ &= (1 - p(1 - \epsilon)) E[(\Delta \mathbf{z}_t^{(k)})^T \mathbf{A}\mathbf{e}_{t-1}^{(k)}] \end{aligned} \quad (31)$$

From (29), (30), and (31), we have

$$\begin{aligned} E[\|\mathbf{e}_t^{(k)}\|_2^2] &\leq (1 - \alpha\mu)^2 E[\|\mathbf{e}_{t-1}^{(k)}\|_2^2] + \delta_t (1 - p(1 - \epsilon^2)) E[\|\Delta \mathbf{z}_t^{(k)}\|_2^2] \\ &\quad - 2\delta_t (1 - p(1 - \epsilon)) E[(\Delta \mathbf{z}_t^{(k)})^T \mathbf{A}\mathbf{e}_{t-1}^{(k)}] + \alpha^2 \sigma_k^2 \\ &= (1 - \alpha\mu)^2 E[\|\mathbf{e}_{t-1}^{(k)}\|_2^2] - \delta_t E[\Omega] + \alpha^2 \sigma_k^2 \\ &\stackrel{(a)}{\leq} (1 - \alpha\mu)^2 E[\|\mathbf{e}_{t-1}^{(k)}\|_2^2] + \alpha^2 \sigma_k^2 \\ &\quad - \delta_t (1 - p(1 - \epsilon^2)) (1 - \alpha(2\lambda - \mu))^2 E[\|\mathbf{e}_{t-1}^{(k)}\|_2^2] \\ &\stackrel{(b)}{\leq} \tau_t E[\|\mathbf{e}_{t-1}^{(k)}\|_2^2] + \alpha^2 \sigma_k^2, \end{aligned} \quad (32)$$

where  $\Omega = 2(1 - p(1 - \epsilon))(\Delta \mathbf{z}_t^{(k)})^T \mathbf{A}\mathbf{e}_{t-1}^{(k)} - (1 - p(1 - \epsilon^2))\|\Delta \mathbf{z}_t^{(k)}\|_2^2$ ,  $\tau_t = (1 - \alpha\mu)^2 - \delta_t(1 - p(1 - \epsilon^2))(1 - \alpha(2\lambda - \mu))^2$ , and (a) is because  $E[\Omega] \geq (1 - p(1 - \epsilon^2))(1 - \alpha(2\lambda - \mu))^2 E[\|\mathbf{e}_{t-1}^{(k)}\|_2^2]$  (see Appendix D).

Using the recursion relationship (between  $E[\|\mathbf{e}_t^{(k)}\|_2^2]$  and  $E[\|\mathbf{e}_{t-1}^{(k)}\|_2^2]$ ) in (32), and noting that  $\tau_t > 0$  as long as  $\alpha \leq \frac{1}{2\lambda - \mu}$ , we have

$$\begin{aligned}
E[\|\mathbf{e}_t^{(k)}\|_2^2] &\leq \prod_{j=1}^t \tau_j E[\|\mathbf{e}_0^{(k)}\|_2^2] + \alpha^2 \sigma_k^2 \left(1 + \sum_{i=0}^{t-2} \prod_{j=0}^i \tau_{t-j}\right) \\
&= (1 - \alpha\mu)^{2t-2\lfloor \frac{t}{q} \rfloor} \tau^{\lfloor \frac{t}{q} \rfloor} E[\|\mathbf{e}_0^{(k)}\|_2^2] \\
&\quad + \alpha^2 \sigma_k^2 \sum_{i=0}^{q-1} (1 - \alpha\mu)^{2i} \sum_{j=0}^{\lceil \frac{t+1}{q} \rceil - 1} \tau^j (1 - \alpha\mu)^{2j(q-1)} \\
&= (1 - \alpha\mu)^{2t-2\lfloor \frac{t}{q} \rfloor} \tau^{\lfloor \frac{t}{q} \rfloor} E[\|\mathbf{e}_0^{(k)}\|_2^2] \\
&\quad + \alpha^2 \sigma_k^2 \left( \frac{1 - (1 - \alpha\mu)^{2q}}{1 - (1 - \alpha\mu)^2} \right) \frac{1 - (\tau(1 - \alpha\mu)^{2q-2})^{\lceil \frac{t+1}{q} \rceil}}{1 - \tau(1 - \alpha\mu)^{2q-2}} \\
&\stackrel{(a)}{=} (1 - \alpha\mu)^{2t-2\lfloor \frac{t}{q} \rfloor} \tau^{\lfloor \frac{t}{q} \rfloor} E[\|\mathbf{e}_0^{(k)}\|_2^2] \\
&\quad + \frac{\alpha^2 \sigma_k^2}{1 - (1 - \alpha\mu)^2} \frac{1 - (1 - \alpha\mu)^{2q}}{1 - \tau(1 - \alpha\mu)^{2q-2}} \\
&\stackrel{(b)}{\leq} (1 - \alpha\mu)^{2t} e^{-c\lfloor \frac{t}{q} \rfloor} E[\|\mathbf{e}_0^{(k)}\|_2^2] \\
&\quad + \frac{\alpha^2 \sigma_k^2}{1 - (1 - \alpha\mu)^2} \frac{1 - (1 - \alpha\mu)^{2q}}{1 - e^{-c}(1 - \alpha\mu)^{2q}},
\end{aligned}$$

where  $q$  is the number of local iterations,  $\tau = (1 - \alpha\mu)^2 - (1 - p(1 - \epsilon^2))(1 - \alpha(2\lambda - \mu))^2$ ,  $c = (1 - p(1 - \epsilon^2))\left(\frac{1 - \alpha(2\lambda - \mu)}{1 - \alpha\mu}\right)^2$ , (a) is because  $1 - (\tau(1 - \alpha\mu)^{2q-2})^{\lceil \frac{t+1}{q} \rceil} \leq 1$ , and (b) is because  $\tau \leq (1 - \alpha\mu)^2 e^{-c}$ , which is the desired result.  $\square$

## APPENDIX C

### PROOF OF LEMMA III.9

*Proof.* We first show a recursive inequality of  $\|E[\mathbf{w}_t^{(k)} - \mathbf{w}_*]\|_2$  and then build an upper bound of this term.

Let  $\mathbf{e}_t^{(k)} = \mathbf{w}_t^{(k)} - \mathbf{w}_*$  and  $\Delta \mathbf{z}_t^{(k)} = \mathbf{z}_t^{(k)} - \bar{\mathbf{z}}_t$ , then from (7), and (10), we have

$$\begin{aligned}
\|E[\mathbf{A}\mathbf{e}_t^{(k)}]\|_2 &= \|E[\mathbf{A}(\mathbf{w}_t^{(k)} - \mathbf{w}_*)]\|_2 \\
&= \|E[\mathbf{A}(\delta_t \mathbf{u}_t^{(k)} \odot \bar{\mathbf{z}}_t + (\mathbf{1} - \delta_t \mathbf{u}_t^{(k)}) \odot \mathbf{z}_t^{(k)} - \mathbf{w}_*)]\|_2 \\
&\stackrel{(a)}{=} \|\delta_t(\epsilon p + 1 - p)E[\mathbf{A}(\bar{\mathbf{z}}_t - \mathbf{w}_*)] + (1 - \delta_t(\epsilon p + 1 - p))E[\mathbf{A}(\mathbf{z}_t^{(k)} - \mathbf{w}_*)]\|_2
\end{aligned}$$

$$= \|\delta_t(\epsilon p + 1 - p)E[\mathbf{A}(\sum_i \eta_i \mathbf{z}_t^{(i)} - \mathbf{w}_*)] + (1 - \delta_t(\epsilon p + 1 - p))E[\mathbf{A}(\mathbf{z}_t^{(k)} - \mathbf{w}_*)]\|_2, \quad (33)$$

where  $\mathbf{A}$  is a matrix independent of  $\mathbf{u}_t^{(k)}$  and (a) is because  $E[u_j] = \epsilon p + 1 - p$  for all element  $u_j$  of  $\mathbf{u}_t^{(k)}$ .

Applying Jensen's inequality yields

$$\begin{aligned} \|E[\mathbf{A}\mathbf{e}_t^{(k)}]\|_2 &\stackrel{(b)}{\leq} \delta_t(\epsilon p + 1 - p) \sum_i \eta_i \|E[\mathbf{A}(\mathbf{z}_t^{(i)} - \mathbf{w}_*)]\|_2 \\ &\quad + (1 - \delta_t(\epsilon p + 1 - p)) \|E[\mathbf{A}(\mathbf{z}_t^{(k)} - \mathbf{w}_*)]\|_2 \\ &\stackrel{(c)}{\leq} \max_i \|E[\mathbf{A}(\mathbf{z}_t^{(i)} - \mathbf{w}_*)]\|_2, \end{aligned} \quad (34)$$

where the last inequality is because  $\sum_i \eta_i = 1$ .

Also, from the update expression (10), we have

$$\begin{aligned} E[\mathbf{A}(\mathbf{z}_t^{(k)} - \mathbf{w}_*)] &= E[\mathbf{A}(\mathbf{w}_{t-1}^{(k)} - \alpha \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(jk)}) - \mathbf{w}_*)] \\ &= E[\mathbf{A}(\mathbf{e}_{t-1}^{(k)} - \alpha(\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_*)))] \\ &\quad + \alpha E[\mathbf{A}(\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(jk)}))]. \end{aligned} \quad (35)$$

Using Taylor's expansion, we have

$$\begin{aligned} E[\mathbf{A}(\mathbf{z}_t^{(k)} - \mathbf{w}_*)] &= E[\mathbf{A}(\mathbf{e}_{t-1}^{(k)} - \alpha \nabla^2 F(\boldsymbol{\xi}_{t-1}^{(k)})\mathbf{e}_{t-1}^{(k)})] \\ &\quad + \alpha E[\mathbf{A}(\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(jk)}))] \\ &= E[\mathbf{A}\mathbf{G}_{t-1}^{(k)}\mathbf{e}_{t-1}^{(k)}] \\ &\quad + \alpha E[\mathbf{A}(\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(jk)}))], \end{aligned}$$

where  $\mathbf{G}_{t-1}^{(k)} = \mathbf{I} - \alpha \nabla^2 F_k(\boldsymbol{\xi}_{t-1}^{(k)})$  and  $\boldsymbol{\xi}_{t-1}^{(k)}$  is a point in the line segment of two endpoints  $\mathbf{w}_{t-1}^{(k)}$  and  $\mathbf{w}_*$ .

Using the law of total expectation yields

$$\begin{aligned} E[\mathbf{A}(\mathbf{z}_t^{(k)} - \mathbf{w}_*)] &= E[\mathbf{A}\mathbf{G}_{t-1}^{(k)}\mathbf{e}_{t-1}^{(k)}] + \alpha E[\mathbf{A}(\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(jk)}))] \\ &= E[\mathbf{A}\mathbf{G}_{t-1}^{(k)}\mathbf{e}_{t-1}^{(k)}] + \alpha E[\mathbf{A}E[\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(jk)}) | \mathbf{A}, \mathbf{w}_{t-1}^{(k)}]] \\ &= E[\mathbf{A}\mathbf{G}_{t-1}^{(k)}\mathbf{e}_{t-1}^{(k)}], \end{aligned} \quad (36)$$

where the last equality is because the input data  $A_t^{(i_k)}$  is sampled identically and independently in each iteration and  $E[\nabla F_k(\mathbf{w}_{t-1}^{(k)}) - \nabla F_k(\mathbf{w}_{t-1}^{(k)}; A_t^{(j_k)}) | \mathbf{A}, \mathbf{w}_{t-1}^{(k)}] = 0$ .

From (34) and (36), we have

$$\begin{aligned} \|E[\mathbf{A}\mathbf{e}_t^{(k)}]\|_2 &\leq \max_i \|E[\mathbf{A}(\mathbf{z}_t^{(i)} - \mathbf{w}_*)]\|_2 \\ &= \max_i E[\mathbf{A}\mathbf{G}_{t-1}^{(i)}\mathbf{e}_{t-1}^{(i)}]. \end{aligned}$$

Letting  $j = \arg \max_i E[\mathbf{A}\mathbf{G}_{t-1}^{(i)}\mathbf{e}_{t-1}^{(i)}]$ , we have

$$\|E[\mathbf{A}\mathbf{e}_t^{(k)}]\|_2 \leq E[\mathbf{A}\mathbf{G}_{t-1}^{(j)}\mathbf{e}_{t-1}^{(j)}].$$

Applying this inequality yields

$$\begin{aligned} \|E[\mathbf{e}_t^{(k)}]\|_2 &\leq \|E[\mathbf{G}_{t-1}^{(i_{t-1})}\mathbf{e}_{t-1}^{(i_{t-1})}]\|_2 \\ &\leq \|E[\mathbf{G}_{t-1}^{(i_{t-1})}\mathbf{G}_{t-2}^{(i_{t-2})}\mathbf{e}_{t-2}^{(i_{t-2})}]\|_2 \\ &\leq \|E[\mathbf{G}_{t-1}^{(i_{t-1})}\mathbf{G}_{t-2}^{(i_{t-2})}\mathbf{G}_{t-3}^{(i_{t-3})}\dots\mathbf{G}_0^{(i_0)}\mathbf{e}_0^{(i_0)}]\|_2 \\ &\stackrel{(a)}{\leq} E[\|\mathbf{G}_{t-1}^{(i_{t-1})}\dots\mathbf{G}_0^{(i_0)}\mathbf{e}_0^{(i_0)}\|_2] \\ &\stackrel{(b)}{\leq} (1 - \alpha\mu)^t E[\|\mathbf{e}_0^{(i_0)}\|_2] \\ &\stackrel{(c)}{\leq} (1 - \alpha\mu)^t \sqrt{\zeta}, \end{aligned} \tag{37}$$

where  $i_{t-1} = \arg \max_i \|E[\mathbf{G}_{t-1}^{(i)}\mathbf{e}_{t-1}^{(i)}]\|_2$  and  $i_{t-j} = \arg \max_i \|E[\mathbf{G}_{t-1}^{(i_{t-1})}\dots\mathbf{G}_{t-j}^{(i)}\mathbf{e}_{t-j}^{(i)}]\|_2$  ( $j = 2, 3, \dots, t$ ),  $\zeta = \max_i E[\|\mathbf{e}_0^{(i)}\|_2^2]$ , (a) is because of Jensen's inequality ( $\|\gamma\mathbf{x} + (1 - \gamma)\mathbf{y}\|_2 \leq \gamma\|\mathbf{x}\|_2 + (1 - \gamma)\|\mathbf{y}\|_2$ ), (b) is because  $\|\mathbf{G}_{t-1}^{(i)}\| \leq 1 - \alpha\mu$ , and (c) is because  $\zeta \geq E[\|\mathbf{e}_0^{(i_0)}\|_2^2] \geq (E[\|\mathbf{e}_0^{(i_0)}\|_2])^2$ , which is the desired result. □

## APPENDIX D

$$\text{PROOF OF } \Omega \geq (1 - p(1 - \epsilon^2))(1 - \alpha(2\lambda - \mu))^2 \|\mathbf{e}_{t-1}^{(k)}\|_2^2$$

*Proof.* In this proof, we show that  $\Omega \geq (1 - p(1 - \epsilon^2))\|\mathbf{e}_{t-1}^{(k)}\|_2^2$  for some values of  $\epsilon$  and  $p$  as long as  $\alpha < \frac{1}{2\lambda - \mu}$ . In fact, we have

$$\begin{aligned} \Omega &= 2(1 - p(1 - \epsilon))(\Delta\mathbf{z}_t^{(k)})^T \mathbf{A}\mathbf{e}_{t-1}^{(k)} \\ &\quad - (1 - p(1 - \epsilon^2))\|\Delta\mathbf{z}_t^{(k)}\|_2^2 \end{aligned}$$

$$\begin{aligned}
&= \|\mathbf{e}_{t-1}^{(k)}\|_2^2 (2(1-p(1-\epsilon))\rho s \\
&\quad - (1-p(1-\epsilon^2))s^2), \tag{38}
\end{aligned}$$

where  $\rho = \frac{(\Delta \mathbf{z}_t^{(k)})^T \mathbf{A} \mathbf{e}_{t-1}^{(k)}}{\|\Delta \mathbf{z}_t^{(k)}\|_2 \|\mathbf{e}_{t-1}^{(k)}\|_2}$  and  $s = \frac{\|\Delta \mathbf{z}_t^{(k)}\|_2}{\|\mathbf{e}_{t-1}^{(k)}\|_2}$ . Note that when  $\|\mathbf{e}_{t-1}^{(k)}\|_2 = 0$  (i.e.,  $s \rightarrow \infty$ ), the algorithm already converges to the optimum  $\mathbf{w}_*$ . When  $\|\Delta \mathbf{z}_t^{(k)}\|_2 = 0$  (i.e.,  $s = 0$ ), it is clear that  $\Omega = 0$ . So, we only need to consider the case of  $0 < s < \infty$ .

First, we recall that for  $\mathbf{x}$  and  $\mathbf{y}$  satisfying  $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 = 1$ , it follows  $2\mathbf{x}^T \mathbf{A} \mathbf{y} = \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{y}^T \mathbf{A} \mathbf{y} - (\mathbf{x} - \mathbf{y})^T \mathbf{A} (\mathbf{x} - \mathbf{y}) \geq 2\lambda_{\min} - \lambda_{\max}$  where  $\lambda_{\max}$  and  $\lambda_{\min}$  are the largest and the smallest eigenvalues of  $\mathbf{A}$ . Since  $\mathbf{A} = \mathbf{I} - \alpha \nabla^2 F_k(\boldsymbol{\xi}_{t-1}^{(k)})$ , we have  $\lambda_{\max} \leq 1 - \alpha\mu$  and  $\lambda_{\min} \geq 1 - \alpha\lambda$ . Therefore, we have

$$\rho \geq \frac{1}{2}(2\lambda_{\min}(\mathbf{A}) - \lambda_{\max}(\mathbf{A})) \geq \frac{1}{2}(1 - \alpha(2\lambda - \mu)). \tag{39}$$

Form (38) and (39), we have

$$\begin{aligned}
\Omega &\geq \|\mathbf{e}_{t-1}^{(k)}\|_2^2 ((1-p(1-\epsilon))(1-\alpha(2\lambda-\mu))s \\
&\quad - (1-p(1-\epsilon^2))s^2).
\end{aligned}$$

Next, we have

$$\begin{aligned}
&\frac{\Omega - \|\mathbf{e}_{t-1}^{(k)}\|_2^2 (1-p(1-\epsilon^2))(1-\alpha(2\lambda-\mu))^2}{\|\mathbf{e}_{t-1}^{(k)}\|_2^2 (s^2 + (1-\alpha(2\lambda-\mu))^2)} \\
&= (1-p(1-\epsilon)) \frac{(1-\alpha(2\lambda-\mu))s}{s^2 + (1-\alpha(2\lambda-\mu))^2} - (1-p(1-\epsilon^2)) \\
&\geq (1-p(1-\epsilon))c_1 - (1-p(1-\epsilon^2)) \\
&\geq g(\epsilon, p),
\end{aligned}$$

where  $c_1 = \min_s \frac{(1-\alpha(2\lambda-\mu))s}{s^2 + (1-\alpha(2\lambda-\mu))^2} (\leq \frac{1}{2})$  and  $g(\epsilon, p) = -p\epsilon^2 + pc_1\epsilon - (1-c_1)(1-p)$ . Noting that  $\frac{1}{2}\lambda\|\mathbf{e}_{t-1}^{(k)}\|_2^2 \geq f(\mathbf{w}_{t-1}^{(k)}) - f(\mathbf{w}_*) \geq \frac{1}{2}\mu\|\mathbf{e}_{t-1}^{(k)}\|_2^2$  and  $f(\mathbf{w}_*) = 0$ , we have

$$\begin{aligned}
c_1 &= \frac{(1-\alpha(2\lambda-\mu))s}{s^2 + (1-\alpha(2\lambda-\mu))^2} \\
&\geq \min \left( \frac{c_2\sqrt{\mu}(1-\alpha(2\lambda-\mu))}{c_2^2(1-\alpha(2\lambda-\mu))^2 + \mu}, \right. \\
&\quad \left. \frac{c_2\sqrt{\lambda}(1-\alpha(2\lambda-\mu))}{c_2^2(1-\alpha(2\lambda-\mu))^2 + \lambda} \right), \tag{40}
\end{aligned}$$

where  $c_2 = \frac{1}{\|\Delta \mathbf{z}_t^{(k)}\|_2} \sqrt{2f(\mathbf{w}_{t-1}^{(k)})}$ .

Now, what remains is to show  $g(\epsilon, p) \geq 0$ . In fact, we have

$$g(\epsilon, p) = \frac{c_1^2}{4} - \frac{(1-c_1)(1-p)}{p} - \left(\epsilon - \frac{c_1}{2}\right)^2. \quad (41)$$

It is not difficult to check that  $g(\epsilon, p) \geq 0$  if

$$\epsilon = \frac{2(1-c_1)(1-p)}{c_1 p}, \quad (42)$$

$$p \geq \frac{4(1-c_1)}{(2-c_1)^2}. \quad (43)$$

Since  $g(\epsilon, p) \geq 0$ , we have  $\Omega - \|\mathbf{e}_{t-1}^{(k)}\|_2^2(1-p(1-\epsilon^2))(1-\alpha(2\lambda-\mu))^2 \geq 0$  which is the desired result. □

## APPENDIX E

### PROOF OF COROLLARY III.3

*Proof.* Using Jensen's inequality, we have

$$\begin{aligned} E[\|\widehat{\mathbf{w}}_t - \mathbf{w}_*\|_2^2] &= E\left[\left\|\sum_k \eta_k \mathbf{w}_t^{(k)} - \mathbf{w}_*\right\|_2^2\right] \\ &\leq \sum_k \eta_k E[\|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2] \\ &\stackrel{(a)}{\leq} \max_k E[\|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2], \end{aligned}$$

where (a) is because  $\sum_k \eta_k = 1$ . What remains is to show that if  $\alpha_t = \frac{\alpha_0}{t+1}$ , then

$$E[\|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2] \leq \frac{c}{t+1}. \quad (44)$$

We will prove (44) using the mathematical induction on  $t$ . First, since  $c \geq \max_k E[\|\mathbf{w}_0^{(k)} - \mathbf{w}_*\|_2^2]$ , it is clear that (44) holds true for  $t = 0$ .

Now we assume the induction hypothesis that (44) holds true for  $t-1$  and check if it also holds true for the case of  $t$ . Letting  $\mathbf{e}_t^{(k)} = \mathbf{w}_t^{(k)} - \mathbf{w}_*$ , and substituting  $\alpha_t$  instead of  $\alpha$  in (32),

we have

$$\begin{aligned}
E[\|\mathbf{e}_t^{(k)}\|_2^2] &\leq (1 - \alpha_t \mu)^2 E[\|\mathbf{e}_{t-1}^{(k)}\|_2^2] + \alpha_t^2 \sigma_k^2 \\
&\leq \left(1 - \frac{\alpha_0 \mu}{t+1}\right)^2 \frac{c}{t} + \frac{\alpha_0^2 \sigma_k^2}{(t+1)^2} \\
&= \frac{c}{t+1} \left( \left(1 - \frac{\alpha_0 \mu}{t+1}\right)^2 \frac{t+1}{t} + \frac{\alpha_0^2 \sigma_k^2}{c(t+1)} \right) \\
&\stackrel{(a)}{\leq} \frac{c}{t+1} \left( \left(1 - \frac{\alpha_0 \mu}{t+1}\right)^2 \frac{t+1}{t} + \frac{2 - (2 - \alpha_0 \mu)^2}{2(t+1)} \right) \\
&= \frac{c}{t+1} \frac{2t^2 + t(2 - \alpha_0^2 \mu^2) + 2(1 - \alpha_0 \mu)^2}{2t(t+1)} \\
&\stackrel{(b)}{<} \frac{c}{t+1} \frac{2t^2 + t(2 - \alpha_0^2 \mu^2) + \alpha_0^2 \mu^2}{2t(t+1)} \\
&\stackrel{(c)}{\leq} \frac{c}{t+1},
\end{aligned}$$

where (a) is because  $c \geq \frac{2\alpha_0^2 \sigma_k^2}{2 - (2 - \alpha_0 \mu)^2}$ , (b) is because  $\frac{2 - \sqrt{2}}{\mu} < \alpha_0 < \frac{2 + \sqrt{2}}{\mu}$ , and (c) is because  $t \geq 1$ , which is the desired result.  $\square$

## APPENDIX F

### PROOF OF THEOREM III.6

*Proof.* We first find a recursive expression of  $E[\sum_k \eta_k \|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2]$  and then prove by induction that

$$E\left[\sum_k \eta_k \|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2\right] \leq \frac{c}{t+1}. \quad (45)$$

From (7) and (10), we have

$$\begin{aligned}
E_{\mathcal{U}}\left[\sum_k \eta_k \|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2\right] &= E_{\mathcal{U}}\left[\sum_k \eta_k \|\delta_t \mathbf{u}_t^{(k)} \odot \bar{\mathbf{z}}_t + (\mathbf{1} - \delta_t \mathbf{u}_t^{(k)}) \odot \mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2\right] \\
&= E_{\mathcal{U}}\left[\sum_k \eta_k \|\mathbf{z}_t^{(k)} - \mathbf{w}_* - \delta_t \mathbf{u}_t^{(k)} \odot (\mathbf{z}_t^{(k)} - \bar{\mathbf{z}}_t)\|_2^2\right] \\
&= \sum_k \eta_k \|\mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2 + \delta_t E_{\mathcal{U}}\left[\sum_k \eta_k \|\mathbf{u}_t^{(k)} \odot (\mathbf{z}_t^{(k)} - \bar{\mathbf{z}}_t)\|_2^2\right] \\
&\quad - \delta_t E_{\mathcal{U}}\left[\sum_k \eta_k \langle 2(\mathbf{z}_t^{(k)} - \mathbf{w}_*), \mathbf{u}_t^{(k)} \odot (\mathbf{z}_t^{(k)} - \bar{\mathbf{z}}_t) \rangle\right] \\
&= \sum_k \eta_k \|\mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2
\end{aligned}$$

$$\begin{aligned}
& + \delta_t E_{\mathcal{U}} \left[ \sum_k \eta_k \langle \mathbf{u}_t^{(k)} \odot (\mathbf{z}_t^{(k)} - \bar{\mathbf{z}}_t) \right. \\
& \quad \left. - 2(\mathbf{z}_t^{(k)} - \mathbf{w}_*), \mathbf{u}_t^{(k)} \odot (\mathbf{z}_t^{(k)} - \bar{\mathbf{z}}_t) \rangle \right] \\
& = \sum_k \eta_k \|\mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2 \\
& \quad + \delta_t E_{\mathcal{U}} \left[ \sum_k \eta_k \langle \mathbf{u}_t^{(k)} \odot \mathbf{u}_t^{(k)} \odot (\mathbf{z}_t^{(k)} - \bar{\mathbf{z}}_t) \right. \\
& \quad \left. - 2\mathbf{u}_t^{(k)} \odot (\mathbf{z}_t^{(k)} - \mathbf{w}_*), \mathbf{z}_t^{(k)} - \bar{\mathbf{z}}_t \rangle \right] \\
& = \sum_k \eta_k \|\mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2 \\
& \quad + \delta_t \sum_k \eta_k \langle (p\epsilon^2 + 1 - p)(\mathbf{z}_t^{(k)} - \bar{\mathbf{z}}_t) \\
& \quad \quad - 2(p\epsilon + 1 - p)(\mathbf{z}_t^{(k)} - \mathbf{w}_*), \mathbf{z}_t^{(k)} - \bar{\mathbf{z}}_t \rangle \\
& \stackrel{(a)}{\leq} \sum_k \eta_k \|\mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2 \\
& \quad + \delta_t (p\epsilon + 1 - p) \sum_k \eta_k (\|\mathbf{z}_t^{(k)} - \bar{\mathbf{z}}_t - (\mathbf{z}_t^{(k)} - \mathbf{w}_*)\|_2^2 \\
& \quad \quad - \|\mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2) \\
& = \sum_k \eta_k \|\mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2 \\
& \quad + \delta_t (p\epsilon + 1 - p) (\|\bar{\mathbf{z}}_t - \mathbf{w}_*\|_2^2 - \sum_k \eta_k \|\mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2) \\
& \stackrel{(b)}{\leq} \sum_k \eta_k \|\mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2,
\end{aligned}$$

where  $E_{\mathcal{U}}[\mathbf{x}]$  is the expected value of  $\mathbf{x}$  with respect to  $\mathcal{U} = \{\mathbf{u}_t^{(k)}\}_k$ , (a) is because  $\epsilon \leq 1$  and  $\langle \mathbf{a} - 2\mathbf{b}, \mathbf{a} \rangle = \|\mathbf{a} - \mathbf{b}\|_2^2 - \|\mathbf{b}\|_2^2$ , and (b) is because  $\|\bar{\mathbf{z}}_t - \mathbf{w}_*\|_2^2 = \|\sum_k \eta_k \mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2 \leq \sum_k \eta_k \|\mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2$ .

Taking expectation again, we have

$$\begin{aligned}
E \left[ \sum_k \eta_k \|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2 \right] & \leq E \left[ \sum_k \eta_k \|\mathbf{z}_t^{(k)} - \mathbf{w}_*\|_2^2 \right] \\
& = \sum_k \eta_k E \left[ \|\mathbf{w}_{t-1}^{(k)} - \mathbf{w}_* - \alpha_t \nabla F_k(\mathbf{w}_{t-1}^{(k)}, A_t^{(i_k)})\|_2^2 \right]
\end{aligned}$$

$$\begin{aligned}
&= \sum_k \eta_k (E[\|\mathbf{w}_{t-1}^{(k)} - \mathbf{w}_*\|_2^2 \\
&\quad - 2\alpha_t \langle \nabla F_k(\mathbf{w}_{t-1}^{(k)}, A_t^{(i_k)}), \mathbf{w}_{t-1}^{(k)} - \mathbf{w}_* \rangle] \\
&\quad + \alpha_t^2 E[\|\nabla F_k(\mathbf{w}_{t-1}^{(k)}, A_t^{(i_k)})\|_2^2]) \\
&= \sum_k \eta_k (E[\|\mathbf{w}_{t-1}^{(k)} - \mathbf{w}_*\|_2^2 \\
&\quad - 2\alpha_t \langle E[\nabla F_k(\mathbf{w}_{t-1}^{(k)}, A_t^{(i_k)}) | \mathbf{w}_{t-1}^{(k)}], \mathbf{w}_{t-1}^{(k)} - \mathbf{w}_* \rangle] \\
&\quad + \alpha_t^2 E[\|\nabla F_k(\mathbf{w}_{t-1}^{(k)}, A_t^{(i_k)})\|_2^2]) \\
&\stackrel{(a)}{=} \sum_k \eta_k (E[\|\mathbf{w}_{t-1}^{(k)} - \mathbf{w}_*\|_2^2 \\
&\quad - 2\alpha_t \langle \nabla F_k(\mathbf{w}_{t-1}^{(k)}), \mathbf{w}_{t-1}^{(k)} - \mathbf{w}_* \rangle] \\
&\quad + \alpha_t^2 E[\|\nabla F_k(\mathbf{w}_{t-1}^{(k)}, A_t^{(i_k)})\|_2^2]) \\
&\stackrel{(b)}{\leq} \sum_k \eta_k ((1 - \alpha_t \mu) E[\|\mathbf{w}_{t-1}^{(k)} - \mathbf{w}_*\|_2^2] \\
&\quad + \alpha_t^2 E[\|\nabla F_k(\mathbf{w}_{t-1}^{(k)}, A_t^{(i_k)})\|_2^2]) \\
&\stackrel{(c)}{\leq} (1 - \alpha_t \mu) E[\sum_k \eta_k \|\mathbf{w}_{t-1}^{(k)} - \mathbf{w}_*\|_2^2] + \alpha_t^2 \sum_k \eta_k \sigma_k^2, \quad (46)
\end{aligned}$$

where (a) is because  $E[\nabla F_k(\mathbf{w}_{t-1}^{(k)}, A_t^{(i_k)}) | \mathbf{w}_{t-1}^{(k)}] = \nabla F_k(\mathbf{w}_{t-1}^{(k)})$ , (b) is due to **A4**, and (c) is due to **A5**.

Letting  $\mathbf{e}_t = E[\sum_k \eta_k \|\mathbf{w}_t^{(k)} - \mathbf{w}_*\|_2^2]$ , we will prove (45) using induction on  $t$ . First, since  $c \geq E[\sum_k \eta_k \|\mathbf{w}_0^{(k)} - \mathbf{w}_*\|_2^2]$ , it is clear that (45) holds true for  $t = 0$ .

Now we assume the induction hypothesis that (45) holds true for  $t - 1$  and check if it also holds true for the case of  $t$ . We have

$$\begin{aligned}
\mathbf{e}_t &\leq (1 - \alpha_t \mu) \mathbf{e}_{t-1} + \alpha_t^2 \sigma^2 \\
&\leq \left(1 - \frac{\alpha_0 \mu}{t+1}\right) \frac{c}{t} + \frac{\alpha_0^2 \sigma^2}{(t+1)^2} \\
&= \frac{c}{t+1} \left( \left(1 - \frac{\alpha_0 \mu}{t+1}\right) \frac{t+1}{t} + \frac{\alpha_0^2 \sigma^2}{c(t+1)} \right) \\
&\stackrel{(a)}{\leq} \frac{c}{t+1} \left( \left(1 - \frac{\alpha_0 \mu}{t+1}\right) \frac{t+1}{t} + \frac{\alpha_0 \mu - 1}{t+1} \right)
\end{aligned}$$

$$\begin{aligned}
&= \frac{c}{t+1} \frac{t(t+1) - (\alpha_0\mu - 1)}{t(t+1)} \\
&\stackrel{(b)}{<} \frac{c}{t+1},
\end{aligned}$$

where  $\sigma^2 = \sum_k \eta_k \sigma_k^2$ , (a) is because  $c \geq \frac{\alpha_0^2 \sigma^2}{\alpha_0 \mu - 1}$  and (b) is because  $\alpha_0 > \frac{1}{\mu}$ , which is the desired result. □

## REFERENCES

- [1] L. T. Nguyen and B. Shim, “Gradual Federated Learning Using Simulated Annealing,” in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, 2021.
- [2] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, “Communication-Efficient Learning of Deep Networks from Decentralized Data,” in *Proc. Int. Conf. Artificial Intell. Stat. (AISTATS)*, vol. 54, 2017, pp. 1273–82.
- [3] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, “Federated learning: Challenges, methods, and future directions,” *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 50–60, May 2020.
- [4] M. Chen, O. Semiari, W. Saad, X. Liu, and C. Yin, “Federated echo state learning for minimizing breaks in presence in wireless virtual reality networks,” *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 177–191, Jan. 2019.
- [5] L. T. Nguyen, J. Kim, and B. Shim, “Low-Rank Matrix Completion: A Contemporary Survey,” *IEEE Access*, vol. 7, no. 1, pp. 94215–94237, July 2019.
- [6] S. Samarakoon, M. Bennis, W. Saad, and M. Debbah, “Distributed federated learning for ultra-reliable low-latency vehicular communications,” *IEEE Trans. Commun.*, vol. 68, no. 2, pp. 1146–1159, Feb. 2019.
- [7] Q. Yang, Y. Liu, T. Chen, and Y. Tong, “Federated machine learning: Concept and applications,” *ACM Trans. Intell. Syst. Technol.*, vol. 10, no. 2, pp. 1–19, Jan. 2019.
- [8] S. Wang, T. Tuor, T. Salonidis, K. K. Leung, C. Makaya, T. He, and K. Chan, “Adaptive federated learning in resource constrained edge computing systems,” *IEEE J. Sel. Areas Commun.*, vol. 37, no. 6, pp. 1205–1221, Jun. 2019.
- [9] F. Sattler, S. Wiedemann, K.-R. Müller, and W. Samek, “Robust and communication-efficient federated learning from non-i.i.d. data,” *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 31, no. 9, pp. 3400–3413, Sept. 2020.
- [10] X. Wang, Y. Han, C. Wang, Q. Zhao, X. Chen, and M. Chen, “In-edge AI: Intelligentizing mobile edge computing, caching and communication by federated learning,” *IEEE Netw. Mag.*, vol. 33, no. 5, pp. 156–165, Sept. 2019.
- [11] Y. Zhang, M. J. Wainwright, and J. C. Duchi, “Communication-efficient algorithms for statistical optimization,” in *Proc. Int. Conf. Neural Inform. Process. Syst.*, Lake Tahoe, NV, USA, Dec. 2012, pp. 1502–1510.
- [12] Y. Arjevani and O. Shamir, “Communication complexity of distributed convex learning and optimization,” in *Proc. NIPS*, pp. 1756–1764, Dec. 2015.
- [13] M. Zinkevich, M. Weimer, L. Li, and A. J. Smola, “Parallelized stochastic gradient descent,” in *Proc. NIPS*, pp. 2595–2603, 2010.
- [14] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi, “Optimization by simulated annealing,” *Science*, vol. 220, no. 4598, pp. 671–680, 1983.
- [15] R. W. Eglese, “Simulated annealing: A tool for operation research,” *Eur. J. Oper. Res.*, vol. 46, no. 3, pp. 271–281, 1990.

- [16] M. Locatelli, “Simulated annealing algorithms for continuous global optimization: Convergence conditions,” *J. Optim. Theory Appl.*, vol. 104, no. 1, pp. 121–133, 2000.
- [17] J. Chen, W. Zhu, and M. M. Ali, “A hybrid simulated annealing algorithm for nonslicing VLSI floorplanning,” *IEEE Trans. Syst., Man, Cybern. C, Appl. Rev.*, vol. 41, no. 4, pp. 544–553, Jul. 2011.
- [18] X. Han, Y. Dong, L. Yue, and Q. Xu, “State transition simulated annealing algorithm for discrete-continuous optimization problems,” *IEEE Access*, vol. 7, pp. 44391–44403, 2019.
- [19] L. Bottou, F. E. Curtis, and J. Nocedal, “Optimization methods for large-scale machine learning,” *Siam Review*, vol. 60, pp. 223–311, 2018.
- [20] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, “Federated optimization in heterogeneous networks,” in *Proc. Conf. Mach. Learn. Syst.*, 2020.
- [21] P. Jiang and G. Agrawal, “A linear speedup analysis of distributed deep learning with sparse and quantized communication,” In *Proc. 32nd Int. Conf. Neural Inform. Process. Syst.*, 2018, pp. 2530-2541.
- [22] H. Yu, R. Jin, and S. Yang, “On the linear speedup analysis of communication efficient momentum SGD for distributed non-convex optimization,” In *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 7184-7193.
- [23] R. M. Gower, N. Loizou, X. Qian, A. Sailanbayev, E. Shulgin, and P. Richtarik, “SGD: General analysis and improved rates,” in *Proc. 36th Int. Conf. Mach. Learn.*, Long Beach, California, USA: PMLR, 09–15 Jun 2019, vol. 97, pp. 5200–5209.
- [24] I. Necoara, Y. Nesterov, and F. Glineur, “Linear convergence of first order methods for nonstrongly convex optimization,” *Mathematical Programming*, 175(1-2):69–107, 2019.
- [25] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, “How to backdoor federated learning,” arXiv preprint arXiv:1807.00459, 2018.
- [26] Y. LeCun. (1998). The MNIST Database of Handwritten Digits. [Online]. Available: <http://yann.lecun.com/exdb/mnist/>
- [27] H. Xiao, K. Rasul, and R. Vollgraf, “Fashion-MNIST: A novel image dataset for benchmarking machine learning algorithms,” arXiv preprint arXiv:1708.07747.
- [28] A. Krizhevsky, V. Nair, and G. Hinton. (2014). The CIFAR-10 dataset. [Online]. Available: <https://www.cs.toronto.edu/~kriz/cifar.html>
- [29] P. Warden, “Launching the speech commands dataset,” Google Research Blog, 2017.
- [30] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov. 1998.
- [31] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *Proc. Int. Conf. Learn. Representat.*, San Diego, CA, USA, pp. 1–14, 2015.
- [32] V. Smith, C. K. Chiang, M. Sanjabi, and A. S. Talwalkar, “Federated multi-task learning,” in *Proc. NIPS*, pp. 4424-4434, 2017.
- [33] Y. Chen, X. Sun, and Y. Jin, “Communication-Efficient Federated Deep Learning with Layerwise Asynchronous Model Update and Temporally Weighted Aggregation,” *IEEE Trans. Neural Net. Learn. Sys.*, vol. 31, no. 10, pp. 4229–4238, Oct. 2020.
- [34] Y. Yeganeh, A. Farshad, N. Navab, S. Albarqouni, “Inverse distance aggregation for federated learning with non-iid data,” in *Proc. DCL Workshop at MICCAI*, 2020, pp. 150–159.
- [35] J. Chen, R. Monga, S. Bengio, and R. Jozefowicz, “Revisiting distributed synchronous SGD,” in *Proc. ICLR Workshop Track*, 2016. [Online]. Available: <https://arxiv.org/abs/1604.00981>.

- [36] H. Wang, M. Yurochkin, Y. Sun, D. Papailiopoulos, and Y. Khazaeni, “Federated learning with matched averaging,” in *Proc. Int. Conf. Learn. Represent.*, 2020. [Online]. Available: <http://arxiv.org/abs/2002.06440>.
- [37] S. Ji, T. Saravirta, S. Pan, G. Long, and A. Walid, “Emerging trends infederated learning: From model fusion to federated X learning,” *arXivpreprint arXiv:2102.12920*, 2021. [Online]. Available: <https://arxiv.org/abs/2102.12920>.
- [38] S. Ji, S. Pan, G. Long, X. Li, J. Jiang, and Z. Huang, “Learning private neural language modeling with attentive aggregation,” in *Proc. IEEE Int. Joint Conf. Neural Netw. (IJCNN)*, Budapest, Hungary, Jul. 2019, pp. 1-8.
- [39] J. Jiang, S. Ji, G. Long, “Decentralized knowledge acquisition for mobile internet applications,” *World Wide Web* (2020).
- [40] X. Wu, Z. Liang, and J. Wang, “FedMed: A federated learning framework for language modeling,” *Sensors*, vol. 20, no. 14, p. 4048, Jul. 2020.
- [41] Y. Huang, L. Chu, Z. Zhou, L. Wang, J. Liu, J. Pei, Y. Zhang, “Personalized cross-silo federated learning on non-IID data,” in *Proc. Assoc. Adv. Artif. Intell.*, 2021, pp. 7865–7873.