

SAFARI: Sparsity enabled Federated Learning with Limited and Unreliable Communications

Yuzhu Mao*, Zihao Zhao*, Meilin Yang, Le Liang, Yang Liu, Wenbo Ding, Tian Lan, *Senior Member, IEEE*, Xiao-Ping Zhang, *Fellow, IEEE*

Abstract

Federated learning (FL) enables edge devices to collaboratively learn a model in a distributed fashion. Many existing researches have focused on improving communication efficiency of high-dimensional models and addressing bias caused by local updates. However, most of FL algorithms are either based on reliable communications or assume fixed and known unreliability characteristics. In practice, networks could suffer from dynamic channel conditions and non-deterministic disruptions, with time-varying and unknown characteristics. To this end, in this paper we propose a sparsity enabled FL framework with both communication efficiency and bias reduction, termed as SAFARI. It makes novel use of a similarity among client models to rectify and compensate for bias that is resulted from unreliable communications. More precisely, sparse learning is implemented on local clients to mitigate communication overhead, while to cope with unreliable communications, a similarity-based compensation method is proposed to provide surrogates for missing model updates. We analyze SAFARI under bounded dissimilarity and with respect to sparse models. It is demonstrated that SAFARI under unreliable communications is guaranteed to converge at the same rate as the standard FedAvg with perfect communications. Implementations and

* These authors contribute equally.

Y. Mao, Z. Zhao, M. Yang, W. Ding, and X.-P. Zhang are with Tsinghua-Berkeley Shenzhen Institute, Tsinghua Shenzhen International Graduate School, Tsinghua University, China. W. Ding is the corresponding author. E-mail: ({myz20, zhaozh21, yml21} @mails.tsinghua.edu.cn, ding.wenbo@sz.tsinghua.edu.cn)

W. Ding is also with RISC-V International Open Source Laboratory, Shenzhen, China, 518055.

L. Liang is with the National Mobile Communications Research Laboratory, Southeast University, Nanjing 210096, and also with the Purple Mountain Laboratories, Nanjing 211111, China. E-mail: (lliang@seu.edu.cn)

Y. Liu is with the Institute for AI Industry Research (AIR), Tsinghua University, China. E-mail: (liuy03@air.tsinghua.edu.cn)

T. Lan is with the Department of Electrical and Computer Engineering, George Washington University, DC, USA. Email: (tlan@gwu.edu)

X.-P. Zhang is also with the Department of Electrical, Computer and Biomedical Engineering, Ryerson University, Toronto, ON M5B 2K3, Canada. E-mail: (xzhang@ee.ryerson.ca)

evaluations on CIFAR-10 dataset validate the effectiveness of SAFARI by showing that it can achieve the same convergence speed and accuracy as FedAvg with perfect communications, with up to 80% of the model weights being pruned and a high percentage of client updates missing in each round.

I. INTRODUCTION

With rapid deployment of mobile sensing and computing devices, there are growing interests in fully exploiting distributed computing resources, as well as huge volumes of data generated at network edge, for efficient learning [1]. To this end, federated learning (FL) [2], [3] enables distributed edge devices to collaboratively learn a model while maintaining data privacy [4]–[6], by allowing a central server and distributed clients to exchange updated model parameters and performing global aggregations. As wireless communications in practice often have limited network capacity [1], [2], a number of proposals have been made on the communication-efficient FL. Examples include model pruning and sparsity to exploit the structural redundancy of dense models [7] and leveraging multiple local training epochs before periodical global aggregation in order to mitigate communication overhead [8]–[10].

Nevertheless, most of existing FL algorithms either are based on reliable communications [9], [10] or assume fixed and known unreliability characteristics [11], [12]. These assumptions may not hold in real-world FL applications. Protocols for data-intensive communications like the lightweight User Datagram Protocol (UDP) tend to focus on best effort delivery without mechanisms for detecting failures and re-transmission. Reliable transmission of local updates cannot be guaranteed [11]. Further, an underlying wireless network could suffer from dynamic channel conditions and non-deterministic disruptions, whose characteristics are often unknown and time-varying. This raises serious challenges in FL – unpredictable absences of local updates with time-varying characteristics would lead to non-homogeneous bias under non-IID data distribution, potentially introducing an unknown drift and causing slow and unstable convergence.

In this paper, we propose a Sparsity enAbleD Federated leArning framework under limited and unReliable commUnications, termed as SAFARI. When unreliability characteristics are unknown and potentially time-varying, we show that it is possible to rectify the resulting bias in global model aggregation by leveraging similarity among different client models. More precisely, once distributed clients locally train their models with sparse algorithms, the central server (i) updates a similarity matrix tracking the similarity among different clients based on received sparse models, and (ii) for any absent update in the current round, substitutes it with an available update

received from the most similar client. Intuitively, these similarity-based surrogates provide an optimal way of compensating for any missing local updates on the fly. This compensation works even if sparse algorithms are employed, as we show that similarity properties are preserved under sparsity. We formally analyze the impact of such compensations in FL and prove that under bounded dissimilarity (i.e., the difference among sparse models produced by different clients are bounded) and a sufficiently small learning rate, the proposed SAFARI algorithm is guaranteed to converge. Extensive evaluations over several popular sparse algorithms (including MAG and Synflow [13]) are conducted. The experiment results validate our theoretical analysis showing that the proposed SAFARI algorithm under unreliable communications achieves the same asymptotic convergence rate as standard FedAvg with reliable communications, even if 80% of the model weights are pruned and a large percentage (up to 70%) of client updates are lost in each round. SAFARI consistently achieves faster convergence than that without compensation under unreliable communications.

The contributions of this paper are summarized as follows.

- A sparsity enabled robust FL framework, SAFARI, is proposed to simultaneously reduce communication overhead and cope with unreliable communications in FL.
- SAFARI leverages a novel similarity-based compensation scheme that actively tracks client similarity and substitutes any missing update (due to unreliable communications) on the fly with an available model update received from the most similar client.
- We theoretically analyze the impact of such compensation with respect to sparse algorithms and prove that similarity properties are preserved under the use of sparse models.
- We establish global convergence analysis for SAFARI and demonstrate that even with limited and unreliable communications, SAFARI can achieve the same convergence rate of vanilla FedAvg with perfectly reliable communications.
- Experiments on CIFAR-10 dataset validate our theoretical analysis. SAFARI demonstrates fast and stable convergence under unreliable communications and outperforms baselines without compensation.

The rest of this paper is organized as follows. Section II introduces the background and related work as well as the motivation. In Section III, the proposed method is described in details. The theoretical analysis and the experimental results are provided in Section IV and V, respectively. Finally, conclusion remarks are summarized in Section VI.

II. BACKGROUND AND RELATED WORK

A. Federated Learning

Assume an FL system with one central server and m distributed clients. Each client i in the client set $\mathbb{M} = \{1, \dots, m\}$ has a local dataset D_i of n_i data samples. The goal of federated training is to solve the following optimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}_i(\mathbf{x}), \quad (1)$$

where $\mathcal{L}_i(\mathbf{x}) = \sum_{i=1}^m \frac{1}{n_i} \sum_{z \in D_i} \ell_i(\mathbf{x}, z)$ is the local objective function at the i -th client. Specifically, z represents a data sample from D_i and $\ell_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is the local loss function based on the learning model \mathbf{x} and client i 's own data.

In the t -th communication round, the server first broadcasts the global model \mathbf{x}^t to clients. Then each client independently runs τ local iterations by optimization solver such as the stochastic gradient descent (SGD) from the current global model \mathbf{x}^t to optimize its own local objective function $\mathcal{L}_i(\mathbf{x})$. Take the SGD for example and the local iterations are as follows,

$$\begin{cases} \mathbf{x}_{i,0}^t = \mathbf{x}^t, \\ \mathbf{x}_{i,k}^t = \mathbf{x}_{i,k-1}^t - \frac{\eta}{\tau} g_i(\mathbf{x}_{i,k-1}^t | \xi_{i,k}), k \in \mathbb{K}, \end{cases} \quad (2)$$

where η is the learning rate, $g_i(\mathbf{x}_{i,k-1}^t | \xi_{i,k})$ is the stochastic gradient computed with the data batch $\xi_{i,k} \sim D_i$, $\mathbf{x}_{i,k}^t$ is the local model after k local iterations and $\mathbb{K} = \{1, \dots, \tau\}$.

After completing τ iterations of local training, each client i will send the new model $\mathbf{x}_{i,\tau}^t$ back to the central server, and the server will aggregate the received client models to update the global model by:

$$\mathbf{x}^{t+1} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_{i,\tau}^t. \quad (3)$$

B. Practical Issues and Related Work

In the design and application of the FL system, there are some practical issues needed to be considered. According to a recent survey [1], the major issues in FL are summarized as SGD, robust aggregation, upload frequency, privacy leakage and wireless communications. Among these five issues, robust aggregation, upload frequency, and wireless communications are all related to the capability and reliability of communications, and the bias caused by consecutive local SGD steps cannot be neglected in all scenarios. Since privacy leakage is a separate line of

research subjects that can always be combined with other works, in this paper we mainly focus on previous works that aim to address limited communication resources, unreliable communication links and local bias.

Limited Communication Resources. Edge devices in wireless networks usually have limited resources, especially for frequent communications. To reduce the transmission burden at each communication round, gradient or model weight compression is a mainstream technique, including quantization and sparsification. Gradient quantization maps each real-valued gradient/model element to a constant number of bits with lower-precision [14]–[16]. As another line of work, sparsification prunes the dense gradient/model with a large amount of non-zero elements to a sparser one. In practice, these two compression techniques can be jointly used, and sparsification is usually the first step to reduce the number of weights for further quantization and transmission. The simplest way to sparsify a model is to keep only the coordinates with large magnitudes exceeding a selected threshold [17]. More sophisticated methods like unbiased sparsification and variance-reduced sparsified SGD have also been developed for training in a distributed fashion [18]–[20]. One remaining question is that such sparsification operates after the local training completes, which provides no reduction on the computation and memory cost during training.

As the training model becomes larger along with the growth of training data in recent years, sparse learning that pre-conducts sparsification and maintains sparse structure throughout training has been intensively investigated. In [21], fully-connected layers were replaced by sparse ones achieved from an initial sparse topology with evolutionary algorithm before training. The connection sensitive had been investigated in [22] for Single-Shot Network Pruning (SNIP). In [23], the exponentially smoothed gradients was utilized to identify model layers and weights which reduced the error efficiently. You *et al.* proposed to use the change of mask distances between epochs to identify a small sub-network at the early training stage, which could restore the comparable test accuracy to the dense network when being trained independently [24]. Moreover, the sparse topology’s updates based on parameter magnitudes and infrequent gradient calculations in [25] loosened the limitation on the size relationship between sparse model and the corresponding dense model, which further reduced the computation cost for sparse learning. However, despite the success empirical performance of the above sparse learning methods, theoretical analysis of the sparse model’s property is still limited.

Unreliable Communications and Local Bias. Due to the limited capability of distributed

clients and communication channels, the communication reliability cannot be guaranteed in the FL system, especially with wireless networks [1]. Previous work has proposed to address unreliable issues by optimizing the aggregation weights according to the link reliability matrix of communication links [11]. Thus, it requires the knowledge of reliability matrix in advance, which is sometimes infeasible in real-world systems.

To tackle the bias caused by local training steps, methods like drift-reduced SCAFFOLD [26] and Inexact DANE [27] with local approximate sub-solver have proved to be effective when the heterogeneity of local objectives is small enough. Recently, the Bias-Variance Reduced Local SGD algorithm surpasses non-local methods under a more relaxed second-order heterogeneity assumption [28]. But the existing bias-reduction techniques still rely on reliable communications and do not consider the client bias caused by the random update loss.

C. Motivating Applications

In this section, we provide several examples to explain some useful properties in practical FL systems that can be utilized to address the aforementioned issues.

Local Computing Resources. In FL collaborative systems, local clients are always equipped with a certain degree of computing power, no matter small edge devices as smartphones, wearables and sensors, or distributed medical/financial institutes. It makes local sparse learning feasible, and reveals great potentials to achieve highly efficient local training with limited distributed resources.

Clusterable Clients. Although the non-IID data distribution and unstable clients of large amount remain challenging in FL systems, it is useful to notice that the clients in quite a few real-world systems tend to be clusterable in terms of data distribution. For example, in an Internet of vehicles system, vehicles within a certain area tend to record similar transportation information. Besides, the devices within the same smart home system usually collect the features of the same person. In these examples, the dissimilarity between client data in a certain group may be negligible, or even follow IID data distribution for the same learning task.

Note that although the pervasiveness of clusterable clients is demonstrated, the following analysis of our method is built upon the standard assumption on data dissimilarity as previous works [29].

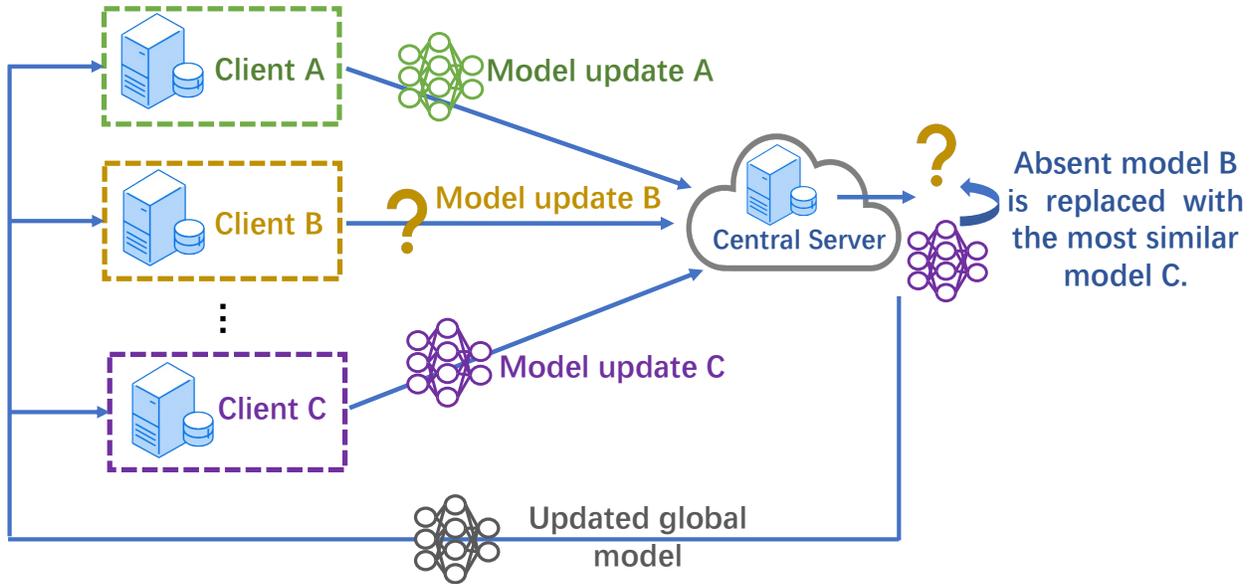


Fig. 1: The schematic illustration of SAFARI.

III. METHODS

In this section, we introduce our approach and elaborate on details of the proposed SAFARI algorithm as illustrated in Fig. 1.

A. Core Concepts and Approach

Here we first describe the two building blocks of the proposed SAFARI algorithm, which are *the sparsity enabled communication efficiency* and *the similarity assisted bias reduction with unreliable wireless communications*. The target problem and the proposed solution are explained in details.

FL with Limited and Unreliable Communications: According to the previous work, the lightweight message based connectionless protocol UDP is commonly used in resource-limited wireless communications. Specifically, UDP reduces much overhead by omitting mechanisms like ACK message confirmation and lost package retransmission [11]. Therefore, despite the relatively low communication overhead, the transmission reliability can not be guaranteed in UDP transmissions. Assume a link reliability list $P = \{p_1^t, \dots, p_m^t\}$, where $0 \leq p_i^t \leq 1$ is the probability that the server successfully receives the local model $\mathbf{x}_{i,\tau}^t$ from client i at the communication round t . In real-world scenarios, each server-client link's reliability could depend

on several factors, i.e., the quality of the channel, the distance between the central server and the corresponding client, as well as the reliability of the client device.

Sparsity enabled Communication Efficiency and Similarity assisted Bias Reduction: To save computing resources and training/inference time, sparse learning on large neural networks has been widely deployed in the deep learning field [7], [21], [22], [25], [30]. When being introduced to FL scenarios, it can save the communication overload by reducing the amount of model weights to be sent. In this context, we propose to conduct the sparse learning at local clients, and utilize the similarity of sparse models to address the bias caused by unreliable communications. Concretely, the server keeps a record of the similarity across clients, which is measured by the sparse models they produce. The similarity record changes along with the training process according to the sparse models successfully received at each global round. With this record, for inactive client whose model has not been received by the server (client fails to participate in training or encounters network failure), the missing model is substituted by the model from the most similar active client.

We will show in the theoretical part that in such way, the bias caused by random loss of local updates can be entirely eliminated when the clients are clusterable, or at least be limited to the same order of the intrinsic data dissimilarity bound under more general scenarios. This enables us to keep the same asymptotic convergence rate as vanilla FedAvg with perfectly reliable communications.

B. The SAFARI Algorithm

The proposed SAFARI algorithm to address the limited and unreliable communication issue is summarized in Algorithm 1. As in vanilla FedAvg [2], the server first initializes an original global model \mathbf{x}^0 and broadcasts it through communication links. Due to the unreliability of communications, some clients may fail to receive the global model from the server. For each client i that successfully receives the global model, it first calculates a mask \mathcal{M}_i based on a specific sparse algorithm to sparsify the global model's structure, and then performs τ iterations' local SGD with the sparse structure. Once the local sparse training is completed, the client will send the sparse local model $\mathbf{x}_{i,\tau}^t$ back to the server, as illustrated in Algorithm 2.

Algorithm 1 SAFARI

Input: The number of communication rounds T , the learning rate η , the number of local steps τ .

Initialize: The initial dense global model \mathbf{x}^0 .

for $t = 0$ to $T - 1$ **do**

Server broadcasts \mathbf{x}^t to all clients.

for each client i receives the message **in parallel do**

Perform *Local Sparse Training*(\mathbf{x}^t, η, τ).

Send the updated sparse model $\mathbf{x}_{i,\tau}^t$ back to the server.

end for

Server performs *Bias Reduced Global Aggregation*.

end for

Finish the training with global model \mathbf{x}^T .

Algorithm 2 Local Sparse Training.

Input: The received global model \mathbf{x}^t , the learning rate η , the number of local steps τ .

Calculate mask \mathcal{M}_i based on a specific sparse algorithm.

Prune the model for a sparser structure: $\mathbf{x}_{i,0}^t = \mathbf{x}^t \odot \mathcal{M}_i$.

for $k = 1$ to τ **do**

Sample a mini-batch $\xi_{i,k}$ from local dataset D_i .

Compute the local gradient $g_i(\mathbf{x}_{i,0}^t | \xi_{i,k})$.

Local SGD step: $\mathbf{x}_{i,k}^t = \mathbf{x}_{i,k-1}^t - \frac{\eta}{\tau} g_i(\mathbf{x}_{i,k-1}^t | \xi_{i,k})$.

end for

return $\mathbf{x}_{i,\tau}^t$.

Again, since the communications are unreliable, not all of the updated local models can be received by the server. To address the potential bias caused by such random loss of client updates, the server will determine the active client group \mathbb{M}_a based on the received client models. Before the aggregation, the server will update the similarity matrix among active clients, and then replace the model from each missing client j with the received model from the most similar active client

j' , as shown in Algorithm 3. After the total T global rounds, the FL training is completed with a trained global model \mathbf{x}^T .

Algorithm 3 Global Aggregation with Similarity-based Compensation.

Input: The received client models, the whole client set \mathbb{M} , the active client group $\mathbb{M}_a = \emptyset$, and s similarity function.

for each client i whose model has been received **do**

$$\mathbb{M}_a = \mathbb{M}_a \cup \{i\}.$$

end for

Server updates the similarity matrix $\rho \in \mathbb{R}^{m \times m}$ with $\rho_{u,v} = s(\mathbf{x}_{u,\tau}^t, \mathbf{x}_{v,\tau}^t), \forall u, v \in \mathbb{M}_a, u \neq v$.

for each client $j \in \mathbb{M} \setminus \mathbb{M}_a$ **do**

$$j' \leftarrow i \in \mathbb{M}_a \text{ that maximizes } \rho_{i,j}.$$

end for

Server performs global aggregation:

$$\mathbf{x}^{t+1} = \frac{1}{m} (\sum_{i \in \mathbb{M}_a} \mathbf{x}_{i,\tau}^t + \sum_{j \in \mathbb{M} \setminus \mathbb{M}_a} \mathbf{x}_{j',\tau}^t).$$

return \mathbf{x}^{t+1} .

IV. THEORETICAL ANALYSIS

In this section, we analyze the convergence of our method and theoretically prove that it can achieve the same convergence rate as the vanilla FedAvg with reliable communications [29].

Notation. In the following part, we use $\|\mathbf{x}\|$, $\|\mathbf{x}\|_1$ and $[\mathbf{x}]_n$ to denote the l_2 , l_1 norms and the n -th element of a vector \mathbf{x} , respectively.

A. Assumptions

1) **Functions:** We first adopt the following three standard assumptions on functions, which are widely used in the non-convex federated optimization field:

- **Smoothness.** The local objective functions are L -smooth, i.e., $\forall i \in \mathbb{M}$:

$$\|\nabla \mathcal{L}_i(\mathbf{x}) - \nabla \mathcal{L}_i(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d. \quad (4)$$

- **Unbiased Gradient and Bounded Variance.** $\forall i \in \mathbb{M}$, the stochastic gradient $g_i(\mathbf{x}|\xi)$ calculated with local data batch ξ is an unbiased estimator of the local gradient: $\mathbb{E}_{\xi \sim D_i} [g_i(\mathbf{x}|\xi)] =$

$\nabla \mathcal{L}_i(\mathbf{x})$, and the variance is bounded by: $\mathbb{E}_{\xi \sim D_i} \|g_i(\mathbf{x}|\xi) - \nabla \mathcal{L}_i(\mathbf{x})\|^2 \leq \sigma^2$, $\forall \mathbf{x} \in \mathbb{R}^d$, $\sigma^2 \geq 0$.

- **Bounded Dissimilarity.** *There exist constants $\beta^2 \geq 1$ and $\zeta^2 \geq 0$ such that:*

$$\frac{1}{m} \sum_{i=1}^m \|\nabla \mathcal{L}_i(\mathbf{x})\|^2 \leq \beta^2 \left\| \frac{1}{m} \sum_{i=1}^m \nabla \mathcal{L}_i(\mathbf{x}) \right\|^2 + \zeta^2. \quad (5)$$

Particularly, $\beta^2 = 1$ and $\zeta^2 = 0$ indicate the IID situation where all the local functions are identical.

2) **Sparse Models:** To analyze the property of local training with sparse models, a common assumption on the mask-induced error is also adopted from sparsification-related literature [7].

- **Mask-induced Error.** *It is assumed that $\forall \mathbf{x} \in \mathbb{R}^d$, the corresponding binary mask $\mathcal{M} \in \{0, 1\}^d$ satisfy*

$$\|\mathbf{x} \odot \mathcal{M} - \mathbf{x}\|^2 \leq \delta^2 \|\mathbf{x}\|^2, 0 < \delta < 1, \quad (6)$$

where \odot denotes the Hadamard product.

Note that the above assumption is a quite relaxed one, which is not limited to any specific sparse algorithms. Furthermore, to analyze the impact of sparse learning in distributed fashion, we make an assumption on the similarity between local training with sparse structures.

- **Similarity Preservation.** *Under the bounded dissimilarity assumption, $\forall \mathbf{x} \in \mathbb{R}^d$, $\forall i, j \in \mathbb{M}$ and local model mask $\{\mathcal{M}_i\}_{i=1}^m$:*

$$\|\nabla \mathcal{L}_i(\mathbf{x} \odot \mathcal{M}_i)\|^2 \leq \beta^2 \|\nabla \mathcal{L}_j(\mathbf{x} \odot \mathcal{M}_j)\|^2 + \zeta^2. \quad (7)$$

The above assumption indicates the rationality behind the compensation based on the similarity among sparse models produced by different clients. The theoretical analysis that demonstrates this assumption will hold for most existing sparse algorithms is provided in the appendix.

3) **Communication Networks:** Similar to previous work, we also make an additional assumption on the unreliable communication network [11]. But compared with the *independent and stable links* assumption made by Ye *et al.*, we extend the condition to cover *independent and unstable links*. In other words, the algorithm proposed in this paper does not require the link reliability to be known in advance or keep stable during training.

- **Independent and Unstable Links.** *The transmissions on different client links are independent and each link's reliability may change during training process.*

B. Descent Lemma with Sparsification

Lemma 1. (Descent Lemma with Sparsification) *With the above assumption on function smoothness, unbiased gradient and bounded variance, as well as sparsification, if $\eta \leq \tau/(6L)$, it holds $\forall i \in \mathbb{M}$, $t \in \mathbb{T} = \{0, \dots, T-1\}$, $k \in \mathbb{K}$ that,*

$$\begin{aligned} \mathbb{E} [\mathcal{L}_i(\mathbf{x}_{i,k}^t)] &\leq \mathbb{E} [\mathcal{L}_i(\mathbf{x}_{i,k-1}^t)] - \frac{\eta}{3\tau} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t)\|^2 \\ &\quad + \frac{\eta^2 L \sigma^2}{2\tau^2} + \frac{2\eta L^2 \delta^2}{3\tau} \mathbb{E} \|\mathbf{x}_{i,k-1}^t\|^2. \end{aligned} \quad (8)$$

From Lemma 1, with the appropriate learning rate, the local objective value will decrease by $\frac{\eta}{3\tau} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t)\|^2$ after every local step. The lemma also meets the expectation that the training will suffer from stochastic gradient variance σ and weight pruning error δ . To the best of the authors' knowledge, rigorous analysis to quantify such error is still uncertain in related research fields, and is also beyond the scope of this work. However, the empirical success of the popular sparse algorithms implies that the tolerance to such error can be quite large in practice [7], [21], [22], [25], [30], which enables us to implement sparse training in FL for communication efficiency, and meanwhile utilizes the properties of sparse models to address the unreliable communications.

C. Global Convergence

To keep consistent and fair comparison with existing FL researches, we build our analysis within the general analysis framework for heterogeneous federated optimization algorithm proposed by [29]. The lemma below points out the influence of unreliable communication links and the following remark explains how our proposed method perfectly addresses such influence with similarity-maintaining sparse models in clusterable scenarios.

Lemma 2. *Under the above assumptions, if $\eta \leq \frac{1}{2\tau L}$, then the optimization error will be bounded as follows:*

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 &\leq \frac{4[\mathcal{L}(\mathbf{x}^0) - \mathcal{L}_{\text{inf}}]}{3\eta\tau T} \\ &\quad + \frac{16\tau\eta L\sigma^2}{3m} + 2\eta^2\sigma^2 L^2(\tau-1) + 4\eta^2 L^2\tau(\tau-1)\zeta^2 \\ &\quad + (2\tau\eta L - 2/3)\frac{1}{m^2}\varphi, \end{aligned} \quad (9)$$

where $\varphi = \sum_{i=1}^m (1 - p_i^t)^2 \mathbb{E} \|\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}\|^2$. Specifically, $\mathbf{h}_i^{(t)} = \frac{1}{\tau} \sum_{k=1}^{\tau} \nabla \mathcal{L}_i(\mathbf{x}_{i,k}^t)$, and i' is the index of the most similar client used for replacing client i in case it is lost. See the appendix for the definition and proof details.

Remark 1. In Lemma 2, φ captures the influence of the random communication network. If we keep the constraints on the learning rate unchanged as in the classic analysis [29], then $(2\tau\eta L - 2/3) \geq 0$ is possible and the global convergence will suffer from φ with any unreliable communication link, i.e., $p_i \leq 1$. However, with the proposed method, if the clients are clusterable, i.e., $\mathbb{E} \|\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}\|^2 = 0$, then φ will be zero and thus the impact of unreliable communication links is entirely eliminated regardless of the unreliability. In this case, the convergence property of vanilla FedAvg with reliable communications will be perfectly preserved. It justifies our strategy to determine the alternative $\mathbf{h}_{i'}^{(t)}$ based on sparse model similarity: the property that the drift between sparse models will be bounded by ζ enables us to reduce the variance caused by unreliable communication into that caused by non-IID data distribution, which is intrinsic in FL and can be addressed by a separate line of research works. See the appendix for proof details.

Furthermore, the following theorem indicates the proposed method converge at the asymptotic rate as vanilla FedAvg with reliable communications even with unclusterable clients.

Theorem 1 *Under the above assumptions, if $\eta = \sqrt{\frac{m}{\tau T}}$, the optimization error after total T iterations is bounded as follows:*

$$\begin{aligned} \min_{t \in \mathbb{T}} \mathbb{E} \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 &\leq \mathcal{O} \left(\frac{1}{\sqrt{m\tau T}} \right) + \mathcal{O} \left(\frac{A\sigma^2}{\sqrt{m\tau T}} \right) \\ &+ \mathcal{O} \left(\frac{mB\sigma^2}{\tau T} \right) + \mathcal{O} \left(\frac{mC\zeta^2}{\tau T} \right), \end{aligned} \quad (10)$$

where $A = \tau, B = \tau - 1, C = \tau(\tau - 1)$, and all other constants are subsumed in \mathcal{O} . Please refer to appendix for details.

Comparison with vanilla FedAvg. Compared with the convergence analysis of FedAvg in [29], the above theorem theoretically indicates that SAFARI with unreliable communications can achieve the same asymptotic convergence rate as FedAvg with reliable communication network under the same parameter setting. Hence, the negative influence of communication unreliability is effectively controlled. In the next section, the experiment results that confirm our theoretical analysis are provided.

V. EXPERIMENTS

We evaluate the proposed framework with different sparse algorithms with 10 clients. We train the ResNet-20 model [31] on the CIFAR-10 dataset, which contains 50,000 images for training and 10,000 images for testing. Specifically, the models are trained using Adam [32] optimizer with learning rate of 0.001, batch size of 64 and tested using batch size of 256. All of our experimental results are trained and evaluated using two NVIDIA-3090 GPUs with 24GB GPU RAM.

A. Performance of SAFARI on Non-IID Data Distribution

To evaluate the generalization of our framework, we have compared the performance of five representative neural network pruning algorithms with SAFARI. The sparsity level α is set to 80%, where 80% of model parameters will be pruned to 0. The selected pruning algorithms include: (1) Rand [13]: randomly prunes 80% parameters; (2) MAG [13]: prunes the 80% smallest absolute values of the model parameters; (3) GraSP [33]: preserves 80% of gradient flow through the network; (4) Synflow [13]: uses the synaptic saliency score to determine the importance of parameters in the network; (5) SNIP [22]: refers to the discussion in Section II.

Following the balanced non-IID data partition setting [34] in FL, 10 total clients are divided into 2 groups equally, and each client contains 5 labels in CIFAR-10. Besides, local steps $\tau = 5$ and local learning rate $\eta = 0.001$ are set to perform the local sparse training in Algorithm 2.

In addition, as addressed in Section III, the successful transmission probability for each link P are chosen as $\{1, 0.3, 0.3, 0.3, 0.3, 1, 0.3, 0.3, 0.3, 0.3\}$, where at least one client in each group will participate in every training process, while the other links has the failure probability of 0.3.

The results of MAG are shown in Fig. 2. Without the proposed similarity-based compensation scheme for bias reduction, the unreliable communication channel will cause huge concussion during the global model training procedure. However, by introducing the compensation based on the similarity between sparse models, the model performance becomes more stable and accurate, which could reach 98% (top-5) accuracy in CIFAR-10.

We also investigate the performance of SAFARI with Synflow. Fig. 3 compares the convergence performance with respect to the number of iterations of Synflow training with and without compensation, as well as the original experiments with no dropouts. It can be seen that the training of Synflow with compensation has achieved nearly identical rate of convergence and convergence stability of Synflow with no dropouts, which is far superior to Synflow without

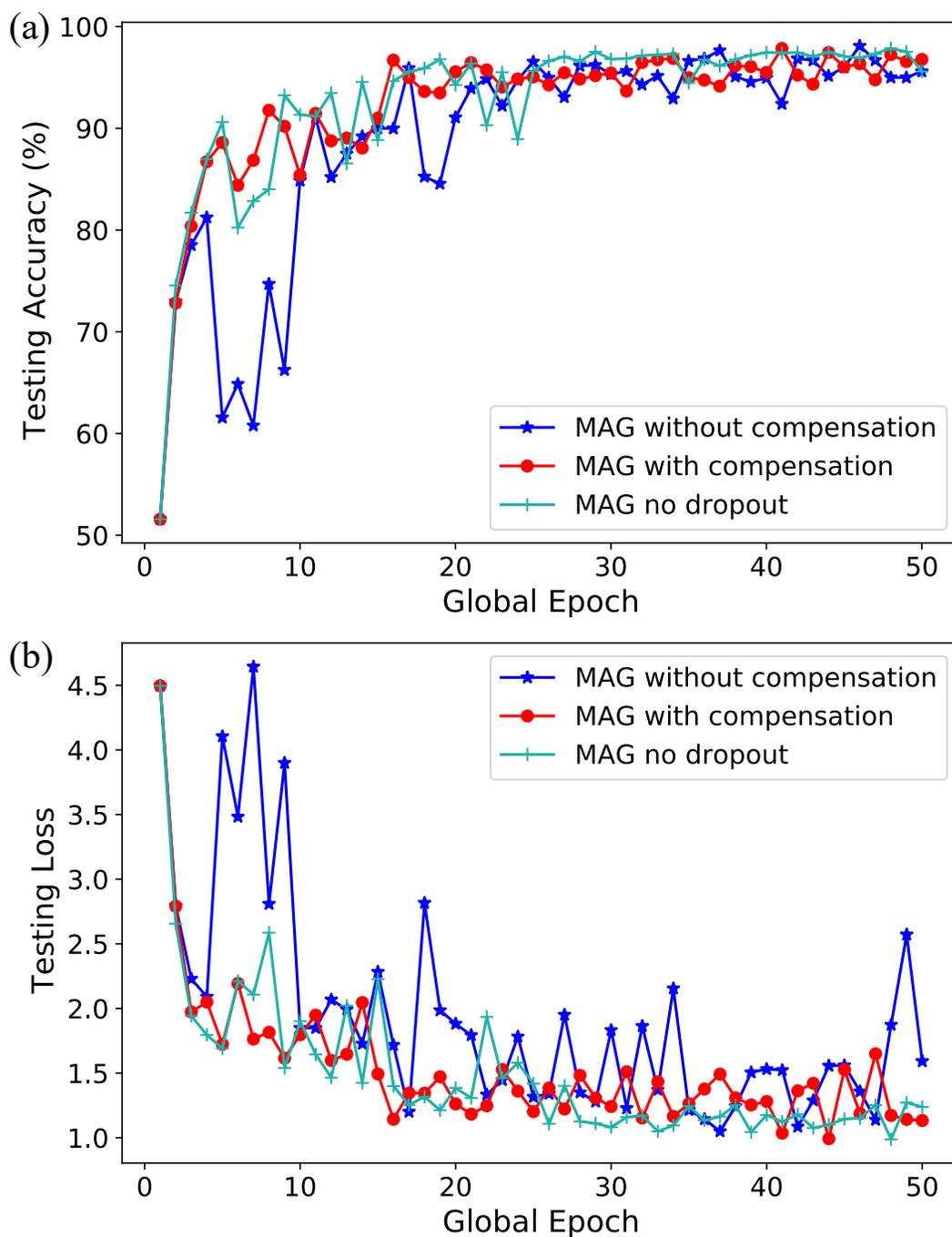


Fig. 2: Performance of SAFARI with MAG: (a) Testing accuracy; (b) Testing loss.

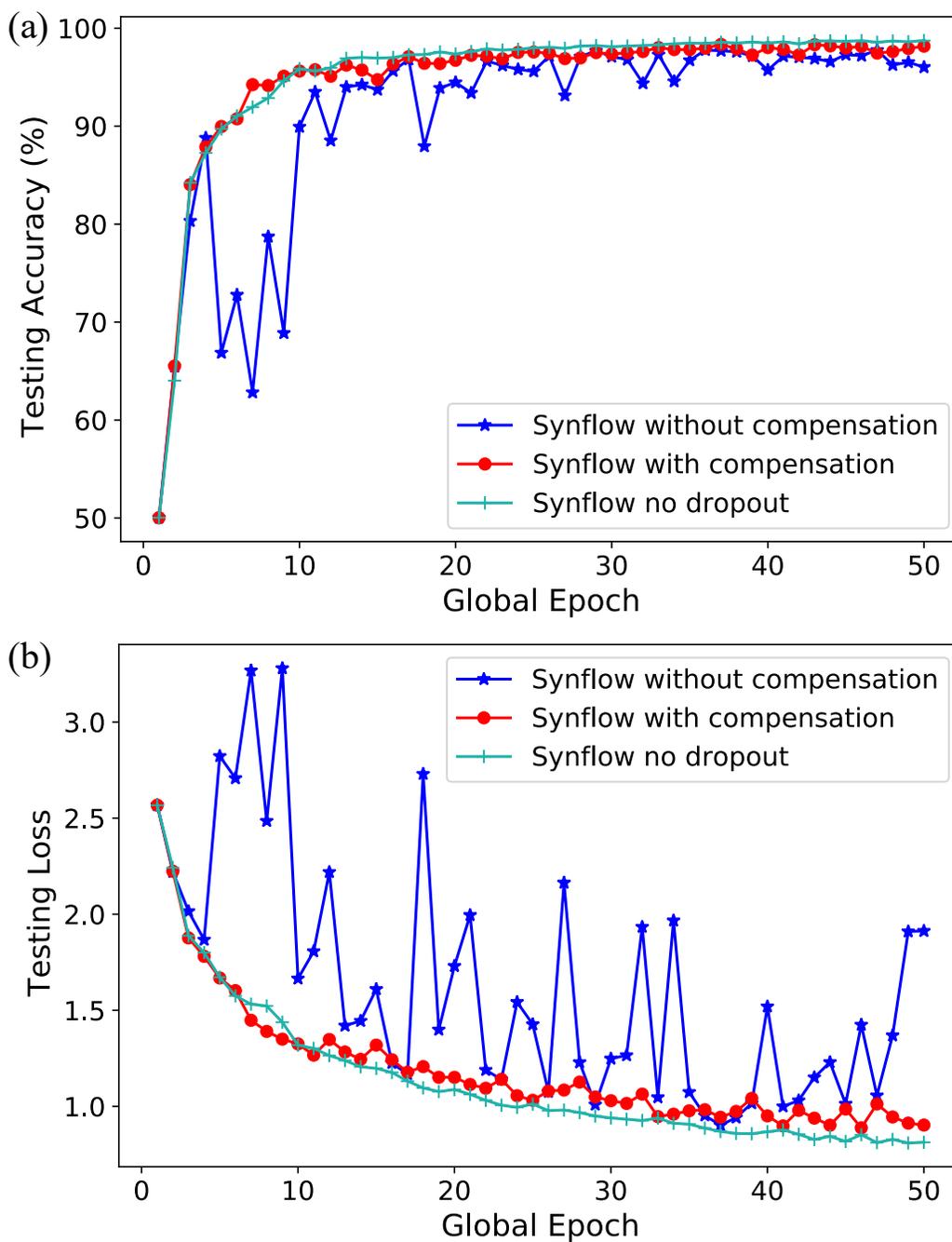


Fig. 3: Performance of SAFARI with Synflow: (a) Testing accuracy; (b) Testing loss.

compensation. The results of Fig. 3 exhibit consistent performance with Fig. 2, which indicates the similarity-based compensation scheme will significantly improve the stability and speed of convergence.

B. Validity of Similarity-based Compensation Scheme

In this section, the experiments are conducted to verify the validity of the proposed similarity-based compensation scheme. Following the lemmas in Section IV, the l_2 -norm based distance of model parameters of two clients u, v is adopted in our experiment as the similarity function $s(\mathbf{x}_u, \mathbf{x}_v)$ in Algorithm 3:

$$s(\mathbf{x}_u, \mathbf{x}_v) := \|\mathbf{x}_u - \mathbf{x}_v\|. \quad (11)$$

Particularly, we display the final similarity matrix ρ after the whole training is completed, as plotted in Fig. 4. For this experiment, following the basic setting, clients 0 to 4 are in Group 1 and have the same label split, while clients 5 to 9 are in Group 2. The lighter colored areas in the upper left and lower right corners indicate that the similarity between members of each group is relatively low. However, the areas in the lower left and upper right corners of this figure represent the similarity of clients between the two groups, and the dark colors between them indicate the high similarity. In the proposed scheme, since the client with the lightest colors in each line will be chosen to compensate, this figure proves that most of the time one client tends to select the model parameters in exactly the same group.

C. Evaluation of Stability

In this section, we evaluate the validity and stability of SAFARI by changing two decisive hyperparameters: the successful transmission probability P and the sparsity level α . Specifically, only these two hyperparameters are changed and the other experiment settings are following the Synflow test.

1) *Successful Transmission Probability*: In order to explore the performance of SAFARI under different P , five different P are selected and the result is shown in Fig. 5. We can observe that as the successful transmission probability P changes, the testing accuracy will tend to be consistent, but the deduction in P can significantly increase the volatility of testing loss.

2) *Sparsity Level*: Fig. 6 shows the the training accuracy and training loss of SAFARI with different sparsity levels from 0.2 to 0.5. When the sparsity level α varies, We can find that the change in both the testing accuracy and testing loss is relatively small, and will eventually converge to a consistent interval.

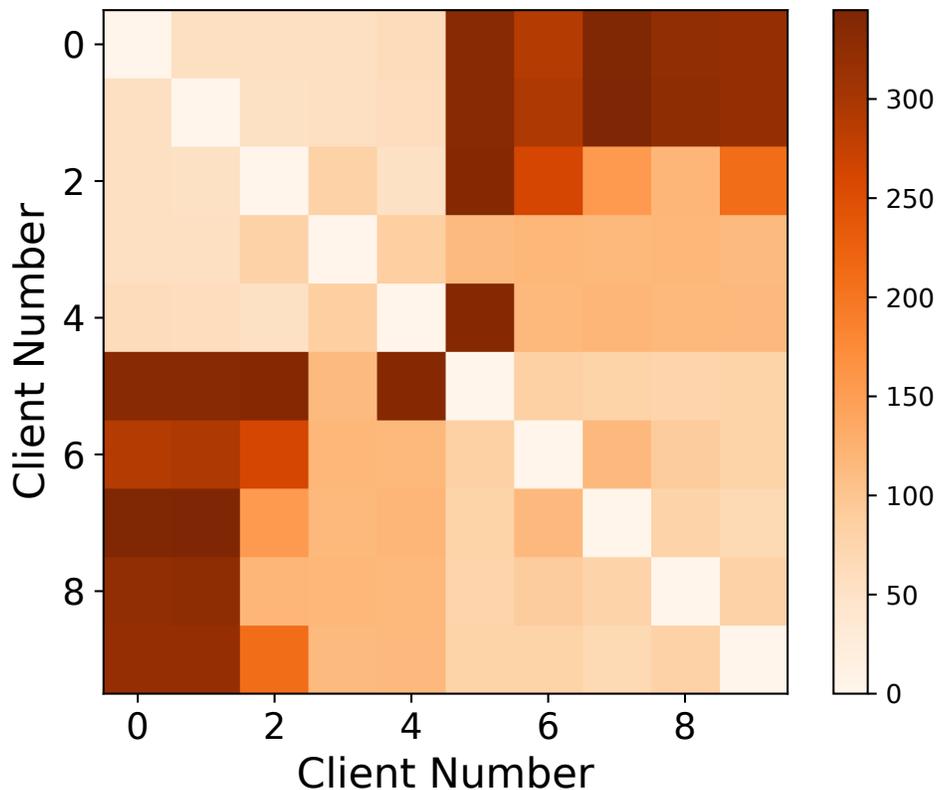


Fig. 4: The final similarity matrix.

VI. CONCLUSION

In this paper, we propose a sparsity enabled robust FL framework, named as SAFARI, which can reduce communication overhead by local sparse learning, and meanwhile rectify the aggregation bias resulted from unreliable communications with unknown and potentially time-varying unreliability characteristics. Our theoretical analysis with respect to sparse model demonstrates that the similarity properties of client models are preserved under sparsity, and thus the proposed SAFARI algorithm with the similarity-based compensation can achieve the same asymptotic convergence rate as FedAvg with reliable communications. The experiments with CIFAR10 dataset and several representative sparse algorithms show that SAFARI can not only save up to 80% communication overhead but also consistently outperforms baselines by achieving fast and stable convergence under unreliable communications. Future work includes more sophisticated algorithm designs for more complex FL scenarios.

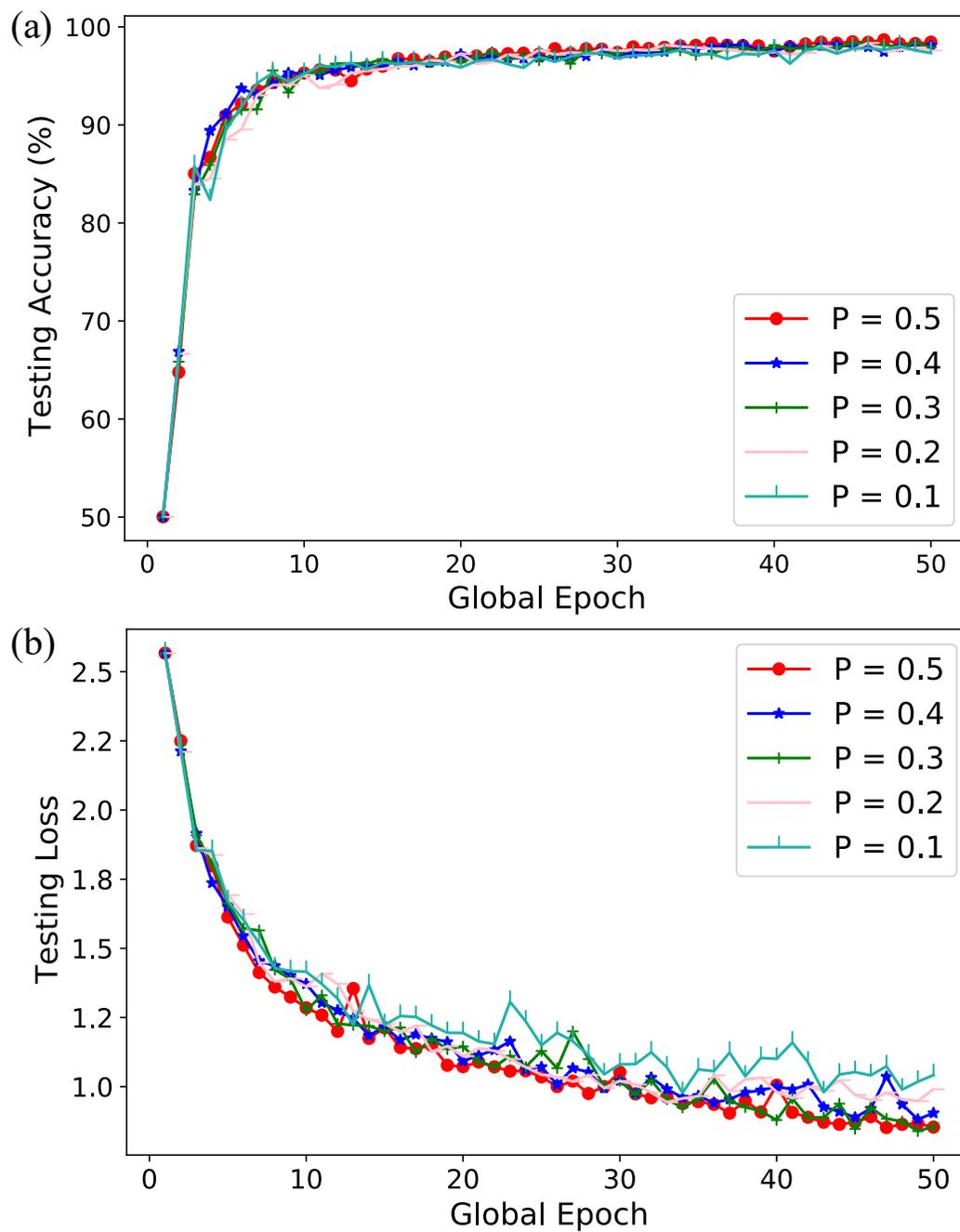


Fig. 5: Performance under different transmission probability settings: (a) Testing accuracy; (b) Testing loss.

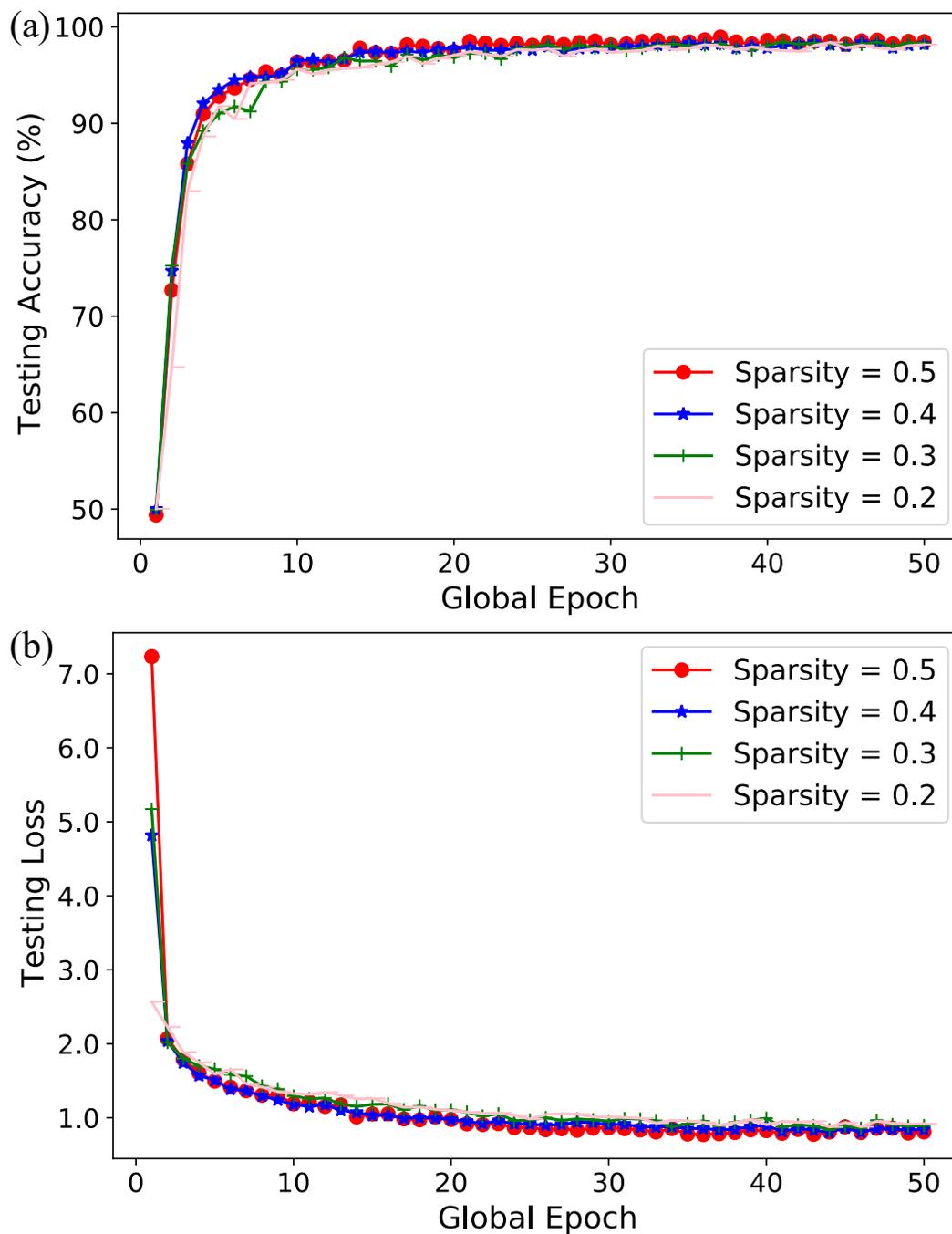


Fig. 6: Performance under different sparsity settings: (a) Testing accuracy; (b) Testing loss.

APPENDIX

In this section, we prove that under the bounded dissimilarity assumption, the local sparse models calculated by sparse learning methods maintain the relationship between local data distribution which is reflected by ζ .

Here we first analyze the masks calculated by sparse learning algorithms. Consider a client i with data input z_i sampled from local dataset D_i and a dense neural network $\mathbf{x} \in \mathbb{R}^d$. Take the mask calculated with SNIP for example, the mask $\mathcal{M}_i \in \mathbb{R}^d$ for the given model \mathbf{x} comes from the connection sensitivity of each model weight $[\mathbf{x}]_n, n = 1, \dots, d$, which is defined as the effect of removing the connection,

$$[\Delta\mathcal{L}_i(\mathbf{x}; D_i)]_n = \mathcal{L}_i(\mathbf{1} \odot \mathbf{x}; D_i) - \mathcal{L}_i((\mathbf{1} - \mathbf{e}_n) \odot \mathbf{x}; D_i),$$

where \mathbf{e}_n is the one-hot indicating vector of the n -th element (i.e. zeros everywhere except at the index n) and $\mathbf{1}$ is a all-one vector with the same length of \mathbf{x} .

To avoid the expensive $d+1$ forward passes over the dataset to calculate the precise $\Delta\mathcal{L}(\mathbf{x}; D_i) \in \mathbb{R}^d$, the connection sensitivity is estimated by its infinitesimal version,

$$[\Delta\mathcal{L}_i(\mathbf{x}; D_i)]_n \approx \lim_{\delta \rightarrow 0} \frac{\mathcal{L}_i(\mathbf{1} \odot \mathbf{x}; D_i) - \mathcal{L}_i((\mathbf{1} - \delta \mathbf{e}_n) \odot \mathbf{x}; D_i)}{\delta}.$$

Thus, the connection sensitivity is actually measured by the change in loss with respect to an infinitesimal multiplicative perturbation δ in weight $[\mathbf{x}]_n$. According to the SNIP algorithm design, this is computed by automatic differentiation in one forward-backward pass [22]. Note that the gradient with respect to the weight $[\mathbf{x}]_n$ is defined as numerical differentiation measured with respect to an additive change as follows,

$$[\nabla\mathcal{L}_i(\mathbf{x}; D_i)]_n \approx \lim_{\delta \rightarrow 0} \frac{\mathcal{L}_i(\mathbf{x}; D_i) - \mathcal{L}_i(\mathbf{x} + \delta \mathbf{e}_n; D_i)}{\delta}. \quad (12)$$

Considering that automatic differentiation is usually used to avoid the error of numerical differentiation during back propagation, the connection sensitivity of the model weights can be regarded as being equivalent to the gradients calculated with the same dataset,

$$\Delta\mathcal{L}_i(\mathbf{x}; D_i) \approx \nabla\mathcal{L}_i(\mathbf{x}; D_i). \quad (13)$$

Intuitively speaking, the weight with a higher connection sensitivity value has a considerable effect on the loss, and thus should be preserved. Therefore, the probability that the weight $[\mathbf{x}]_n$ is preserved after pruning ($[\mathcal{M}_i]_n = 1$) is measured by the corresponding sensitivity magnitude,

$$s_n = \frac{|[\Delta\mathcal{L}_i(\mathbf{x}; D_i)]_n|}{\sum_{k=1}^d |[\Delta\mathcal{L}_i(\mathbf{x}; D_i)]_k|}, \quad (14)$$

and there is,

$$\mathbb{E} \|g_i(\mathbf{x} \odot \mathcal{M}_i | z_i)\|_1 = \frac{\mathbb{E} \|g_i(\mathbf{x} | z_i) \odot \Delta \mathcal{L}_i(\mathbf{x}; D_i)\|_1}{\mathbb{E} \|\Delta \mathcal{L}_i(\mathbf{x}; D_i)\|_1}. \quad (15)$$

The (13) and (15) indicate that the effect of mask depends only on the dense model weights and local data. It should be noted that in addition to SNIP, the masks from other sparse algorithms are also based only on data and dense models, except with much more direct algorithm designs. Next, assume the initialization of each model weight is independent, and we have

$$\|\nabla \mathcal{L}_i(\mathbf{x} \odot \mathcal{M}_i)\|_1 = \frac{\|\Delta \mathcal{L}_i(\mathbf{x})\|_1 \|\nabla \mathcal{L}_i(\mathbf{x})\|_1}{\|\Delta \mathcal{L}_i(\mathbf{x})\|_1} = \|\nabla \mathcal{L}_i(\mathbf{x})\|_1,$$

where $\nabla \mathcal{L}_i(\mathbf{x}; D_i)$ and $\Delta \mathcal{L}_i(\mathbf{x}; D_i)$ are abbreviated as $\nabla \mathcal{L}_i(\mathbf{x})$ and $\Delta \mathcal{L}_i(\mathbf{x})$, respectively.

With the bounded dissimilarity assumption on the dense gradients, it holds that

$$\frac{1}{m} \sum_{i=1}^m \|\nabla \mathcal{L}_i(\mathbf{x} \odot \mathcal{M}_i)\|^2 \leq \beta^2 \left\| \frac{1}{m} \sum_{i=1}^m \nabla \mathcal{L}_i(\mathbf{x} \odot \mathcal{M}_i) \right\|^2 + \zeta^2.$$

The above result implies that the gradients calculated with well-designed sparse structure maintains the relationship between the dense gradients. Specifically, within a client cluster with similar data distributions and the same global model \mathbf{x} , for any two clients, i.e. i and j , the following conditions are satisfied,

$$\begin{aligned} \mathbb{E}_{z_i \sim D_i} \|g_i(\mathbf{x} \odot \mathcal{M}_i | z_i)\| &= \mathbb{E}_{z_j \sim D_j} \|g_j(\mathbf{x} \odot \mathcal{M}_j | z_j)\|, \\ \|\nabla \mathcal{L}_i(\mathbf{x} \odot \mathcal{M}_i)\| &= \|\nabla \mathcal{L}_j(\mathbf{x} \odot \mathcal{M}_j)\|, \end{aligned}$$

where \mathcal{M}_i and \mathcal{M}_j are the masks computed for clients i and j . Obviously, if all the local functions are identical, the sparse gradients make the bounded dissimilarity assumption hold perfectly with $\beta = 1$ and $\zeta = 0$,

$$\frac{1}{m} \sum_{i=1}^m \|\nabla \mathcal{L}_i(\mathbf{x} \odot \mathcal{M}_i)\|^2 = \left\| \frac{1}{m} \sum_{i=1}^m \nabla \mathcal{L}_i(\mathbf{x} \odot \mathcal{M}_i) \right\|^2. \quad (16)$$

A. Proof of Lemma 1

Here we simplify the local gradient $g_i(\mathbf{x}_{i,k-1}^t | \xi_{i,k})$ calculated with batch $\xi_{i,k}$ from client i as $g_i(\mathbf{x}_{i,k-1}^t)$. Since \mathcal{L}_i is L -smooth (smoothness assumption), at the global iteration t and local iteration k , it holds that for each client i ,

$$\begin{aligned} \mathcal{L}_i(\mathbf{x}_{i,k}^t) &\leq \mathcal{L}_i(\mathbf{x}_{i,k-1}^t) + \frac{L}{2} \|\mathbf{x}_{i,k}^t - \mathbf{x}_{i,k-1}^t\|^2 \\ &\quad + \langle \nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t), \mathbf{x}_{i,k}^t - \mathbf{x}_{i,k-1}^t \rangle \end{aligned} \quad (17)$$

$$\begin{aligned}
&= \mathcal{L}_i(\mathbf{x}_{i,k-1}^t) + \frac{\eta^2 L}{2\tau^2} \|g_i(\mathbf{x}_{i,k-1}^t \odot \mathcal{M}_i)\|^2 \\
&\quad - \frac{\eta}{\tau} \langle \nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t), g_i(\mathbf{x}_{i,k-1}^t \odot \mathcal{M}_i) \rangle.
\end{aligned} \tag{18}$$

Besides, with the unbiased gradient and bounded variance assumption, it is obvious that

$$\begin{aligned}
&\mathbb{E} \langle \nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t), g_i(\mathbf{x}_{i,k-1}^t \odot \mathcal{M}_i) \rangle \\
&= \mathbb{E} \langle \nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t), \nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t \odot \mathcal{M}_i) \rangle,
\end{aligned} \tag{19}$$

and

$$\mathbb{E} \|g_i(\mathbf{x}_{i,k-1}^t \odot \mathcal{M}_i)\|^2 \leq \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t \odot \mathcal{M}_i)\|^2 + \sigma^2. \tag{20}$$

By taking the expectation over (18) and substituting (19) and (20) into (18), we achieve

$$\begin{aligned}
\mathbb{E} [\mathcal{L}_i(\mathbf{x}_{i,k}^t)] &\leq \mathbb{E} [\mathcal{L}_i(\mathbf{x}_{i,k-1}^t)] + \frac{\eta^2 L \sigma^2}{2\tau^2} \\
&\quad + \frac{\eta^2 L}{2\tau^2} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t \odot \mathcal{M}_i)\|^2 \\
&\quad - \frac{\eta}{\tau} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t)\|^2 - \frac{\eta}{\tau} \mathbb{E} \langle \nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t), \mathcal{T}_0 \rangle,
\end{aligned} \tag{21}$$

where $\mathcal{T}_0 = \nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t \odot \mathcal{M}_i) - \nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t)$.

Note that,

$$\begin{aligned}
&\mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t \odot \mathcal{M}_i)\|^2 \\
&\leq 2\mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t)\|^2 \\
&\quad + 2\mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t) - \nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t \odot \mathcal{M}_i)\|^2
\end{aligned} \tag{22}$$

$$\begin{aligned}
&\leq 2\mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t)\|^2 \\
&\quad + 2L^2 \mathbb{E} \|\mathbf{x}_{i,k-1}^t \odot \mathcal{M}_i - \mathbf{x}_{i,k-1}^t\|^2
\end{aligned} \tag{23}$$

$$\leq 2\mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t)\|^2 + 2L^2 \delta^2 \mathbb{E} \|\mathbf{x}_{i,k-1}^t\|^2, \tag{24}$$

where the last two inequalities come from the smoothness assumption and the relaxed assumption on mask-induced error. Next, with $2\langle \mathbf{a}, \mathbf{b} \rangle \leq \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2$, we have,

$$\begin{aligned}
-\mathbb{E} \langle \nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t), \mathcal{T}_0 \rangle &\leq \frac{1}{2} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t)\|^2 \\
&\quad + \frac{1}{2} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t \odot \mathcal{M}_i) - \nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t)\|^2
\end{aligned} \tag{25}$$

$$\leq \frac{1}{2} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t)\|^2 + \frac{L^2 \delta^2}{2} \mathbb{E} \|\mathbf{x}_{i,k-1}^t\|^2. \tag{26}$$

Substitute (24) and (26) into(21), and by setting $\eta \leq \tau/(6L)$, we have

$$\begin{aligned} \mathbb{E} [\mathcal{L}_i(\mathbf{x}_{i,k}^t)] &\leq \mathbb{E} [\mathcal{L}_i(\mathbf{x}_{i,k-1}^t)] + \frac{\eta^2 L \sigma^2}{2\tau^2} \\ &\quad - \frac{\eta(\tau - 2\eta L)}{2\tau^2} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t)\|^2 \\ &\quad + \frac{\eta L^2 \delta^2 (\tau + 2\eta L)}{2\tau^2} \mathbb{E} \|\mathbf{x}_{i,k-1}^t\|^2 \end{aligned} \quad (27)$$

$$\begin{aligned} &\leq \mathbb{E} [\mathcal{L}_i(\mathbf{x}_{i,k-1}^t)] - \frac{\eta}{3\tau} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}_{i,k-1}^t)\|^2 \\ &\quad + \frac{\eta^2 L \sigma^2}{2\tau^2} + \frac{2\eta L^2 \delta^2}{3\tau} \mathbb{E} \|\mathbf{x}_{i,k-1}^t\|^2. \end{aligned} \quad (28)$$

The proof of Lemma 1 is completed.

B. Proof of Lemma 2

From the global point of view, $\mathbf{x}_{i,\tau}^t$ represents the local model sent to server after client i 's local iterations, which is supposed to be a sparse one. Recall that the global model is updated by the following rule under reliable communications:

$$\mathbf{x}^{t+1} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_{i,\tau}^t = \mathbf{x}^t - \tau\eta \sum_{i=1}^m \frac{1}{m} \mathbf{d}_i^{(t)}, \quad (29)$$

where $\mathbf{d}_i^{(t)} = \frac{1}{\tau} \sum_{k=1}^{\tau} g_i(\mathbf{x}_{i,k}^t)$ is the normalized stochastic gradient at client i . Correspondingly, the normalized gradient at each client is defined as

$$\mathbf{h}_i^{(t)} = \frac{1}{\tau} \sum_{k=1}^{\tau} \nabla \mathcal{L}_i(\mathbf{x}_{i,k}^t), i \in \mathbb{M}. \quad (30)$$

To solve the problem caused by unreliable communications, the global model is updated with the proposed compensation based on sparse model similarity. Therefore, the expectation of global model update can be written as

$$\mathbb{E} [\mathbf{x}^{t+1} - \mathbf{x}^t] = -\tau\eta \sum_{i=1}^m \frac{1}{m} \left[p_i^t \mathbf{d}_i^{(t)} + (1-p_i^t) \mathbf{d}_{i'}^{(t)} \right], \quad (31)$$

where i' is the index of the most similar client used for replacing client i in case it is lost, and p_i^t is the reliability of the channel between client i and the server at round t .

According to the smoothness assumption, there is,

$$\begin{aligned} &\mathbb{E} [\mathcal{L}(\mathbf{x}^{t+1})] - \mathcal{L}(\mathbf{x}^t) \\ &\leq \mathbb{E} \langle \nabla \mathcal{L}(\mathbf{x}^t), \mathbf{x}^{t+1} - \mathbf{x}^t \rangle + \frac{L}{2} \mathbb{E} \|\mathbf{x}^{t+1} - \mathbf{x}^t\|^2 \end{aligned} \quad (32)$$

$$\begin{aligned}
&\leq -\tau\eta \underbrace{\mathbb{E} \left\langle \nabla \mathcal{L}(\mathbf{x}^t), \sum_{i=1}^m \frac{1}{m} [p_i^t \mathbf{d}_i^{(t)} + (1-p_i^t) \mathbf{d}_{i'}^{(t)}] \right\rangle}_{T_1} \\
&\quad + \frac{\tau^2 \eta^2 L}{2} \underbrace{\mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} [p_i^t \mathbf{d}_i^{(t)} + (1-p_i^t) \mathbf{d}_{i'}^{(t)}] \right\|^2}_{T_2}. \tag{33}
\end{aligned}$$

Bounding the first term T_1 . For the first term on the right hand side of (33), there is

$$\begin{aligned}
T_1 &= \mathbb{E} \left\langle \nabla \mathcal{L}(\mathbf{x}^t), \sum_{i=1}^m \frac{1}{m} p_i^t (\mathbf{d}_i^{(t)} - \mathbf{h}_i^{(t)} + \mathbf{h}_i^{(t)}) \right\rangle \\
&\quad + \mathbb{E} \left\langle \nabla \mathcal{L}(\mathbf{x}^t), \sum_{i=1}^m \frac{1}{m} (1-p_i^t) (\mathbf{d}_{i'}^{(t)} - \mathbf{h}_{i'}^{(t)} + \mathbf{h}_{i'}^{(t)}) \right\rangle \\
&= \mathbb{E} \left\langle \nabla \mathcal{L}(\mathbf{x}^t), \sum_{i=1}^m \frac{1}{m} p_i^t \mathbf{h}_i^{(t)} \right\rangle \\
&\quad + \mathbb{E} \left\langle \nabla \mathcal{L}(\mathbf{x}^t), \sum_{i=1}^m \frac{1}{m} (1-p_i^t) \mathbf{h}_{i'}^{(t)} \right\rangle, \tag{34}
\end{aligned}$$

where the second equality comes from the unbiased gradient assumption which implies $\mathbb{E}(\mathbf{d}_i^{(t)} - \mathbf{h}_i^{(t)}) = 0$.

$$\begin{aligned}
T_1 &= \mathbb{E} \left\langle \nabla \mathcal{L}(\mathbf{x}^t), \sum_{i=1}^m \frac{1}{m} p_i^t \mathbf{h}_i^{(t)} + \sum_{i=1}^m \frac{1}{m} (1-p_i^t) \mathbf{h}_i^{(t)} \right\rangle \\
&\quad + \mathbb{E} \left\langle \nabla \mathcal{L}(\mathbf{x}^t), \sum_{i=1}^m \frac{1}{m} (1-p_i^t) (\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}) \right\rangle \\
&= \mathbb{E} \left\langle \nabla \mathcal{L}(\mathbf{x}^t), \sum_{i=1}^m \frac{1}{m} \mathbf{h}_i^{(t)} \right\rangle \\
&\quad + \mathbb{E} \left\langle \nabla \mathcal{L}(\mathbf{x}^t), \sum_{i=1}^m \frac{1}{m} (1-p_i^t) (\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}) \right\rangle \\
&= \frac{1}{2} \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + \frac{1}{2} \mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} \mathbf{h}_i^{(t)} \right\|^2 \\
&\quad + \mathbb{E} \left\langle \nabla \mathcal{L}(\mathbf{x}^t), \sum_{i=1}^m \frac{1}{m} (1-p_i^t) (\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}) \right\rangle \\
&\quad - \frac{1}{2} \mathbb{E} \left\| \nabla \mathcal{L}(\mathbf{x}^t) - \sum_{i=1}^m \frac{1}{m} \mathbf{h}_i^{(t)} \right\|^2 \tag{35}
\end{aligned}$$

$$\begin{aligned}
&\leq \frac{1}{2} \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + \frac{1}{2} \mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} \mathbf{h}_i^{(t)} \right\|^2 \\
&\quad + \frac{1}{2} \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + \frac{1}{2} \mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} (1-p_i^t) (\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}) \right\|^2 \\
&\quad - \frac{1}{2} \mathbb{E} \left\| \nabla \mathcal{L}(\mathbf{x}^t) - \sum_{i=1}^m \frac{1}{m} \mathbf{h}_i^{(t)} \right\|^2, \tag{36}
\end{aligned}$$

where the last inequality follows from $2\langle \mathbf{a}, \mathbf{b} \rangle = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 - \|\mathbf{a} - \mathbf{b}\|^2$.

Bounding the second term T_2 . For the second term on the right hand side of(33), there is,

$$\begin{aligned}
T_2 &= \mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} \mathbf{d}_i^{(t)} \right\|^2 + \mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} [(1-p_i^t) (\mathbf{d}_{i'}^{(t)} - \mathbf{d}_i^{(t)})] \right\|^2 \\
&\leq 2\mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} \mathbf{h}_i^{(t)} \right\|^2 + 2\mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} (\mathbf{d}_i^{(t)} - \mathbf{h}_i^{(t)}) \right\|^2 \\
&\quad + 3\mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} (1-p_i^t) (\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}) \right\|^2 \\
&\quad + 3\mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} (1-p_i^t) (\mathbf{h}_i^{(t)} - \mathbf{d}_i^{(t)}) \right\|^2 \\
&\quad + 3\mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} (1-p_i^t) (\mathbf{d}_{i'}^{(t)} - \mathbf{h}_{i'}^{(t)}) \right\|^2 \tag{37}
\end{aligned}$$

$$\begin{aligned}
&= 2\mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} \mathbf{h}_i^{(t)} \right\|^2 + 2 \sum_{i=1}^m \frac{1}{m^2} \mathbb{E} \|\mathbf{d}_i^{(t)} - \mathbf{h}_i^{(t)}\|^2 \\
&\quad + 3 \sum_{i=1}^m \frac{1}{m^2} (1-p_i^t)^2 \mathbb{E} \|\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}\|^2 \\
&\quad + 3 \sum_{i=1}^m \frac{1}{m^2} (1-p_i^t)^2 \mathbb{E} \|\mathbf{h}_i^{(t)} - \mathbf{d}_i^{(t)}\|^2 \\
&\quad + 3 \sum_{i=1}^m \frac{1}{m^2} (1-p_i^t)^2 \mathbb{E} \|\mathbf{d}_{i'}^{(t)} - \mathbf{h}_{i'}^{(t)}\|^2, \tag{38}
\end{aligned}$$

where the second equality is based on the independence of the surrogate selection, and the following inequalities follow from: $\|\sum_{i=1}^n \mathbf{w}_i\|_2^2 \leq n \sum_{i=1}^n \|\mathbf{w}_i\|_2^2$.

With assumption on gradient variance, the second term can be bounded as,

$$T_2 \leq \sum_{i=1}^m \frac{1}{m^2} \left[2 + 6(1-p_i^t)^2 \right] \sigma^2 + 2\mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} \mathbf{h}_i^{(t)} \right\|^2$$

$$+ 3 \sum_{i=1}^m \frac{1}{m^2} (1 - p_i^t)^2 \mathbb{E} \|\mathbf{h}_i^{(t)} - \mathbf{h}_{i'}^{(t)}\|^2. \quad (39)$$

Bounding the objective reduction (33). Plugging (36) and (39) back into (33), there is,

$$\begin{aligned} & \mathbb{E} [\mathcal{L}(\mathbf{x}^{t+1})] - \mathcal{L}(\mathbf{x}^t) \\ & \leq -\tau\eta \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + \left(\tau^2\eta^2 L - \frac{\tau\eta}{2}\right) \mathbb{E} \left\| \sum_{i=1}^m \frac{1}{m} \mathbf{h}_i^{(t)} \right\|^2 \\ & \quad + \frac{\tau\eta}{2} \mathbb{E} \left\| \nabla \mathcal{L}(\mathbf{x}^t) - \sum_{i=1}^m \frac{1}{m} \mathbf{h}_i^{(t)} \right\|^2 \\ & \quad + \left(\frac{3}{2}\tau^2\eta^2 L - \frac{\tau\eta}{2}\right) \sum_{i=1}^m \frac{1}{m^2} (1 - p_i^t)^2 \mathbb{E} \|\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}\|^2 \\ & \quad + \tau^2\eta^2 L \sum_{i=1}^m \frac{1}{m^2} [1 + 3(1 - p_i^t)^2] \sigma^2. \end{aligned} \quad (40)$$

If $\tau\eta L \leq \frac{1}{2}$, we have,

$$\begin{aligned} & \frac{1}{\tau\eta} (\mathbb{E} [\mathcal{L}(\mathbf{x}^{t+1})] - \mathcal{L}(\mathbf{x}^t)) \\ & \leq -\|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + \tau\eta L \sum_{i=1}^m \frac{1}{m^2} [1 + 3(1 - p_i^t)^2] \sigma^2 \\ & \quad + \left(\frac{3}{2}\tau\eta L - \frac{1}{2}\right) \sum_{i=1}^m \frac{1}{m^2} (1 - p_i^t)^2 \mathbb{E} \|\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}\|^2 \\ & \quad + \frac{1}{2} \mathbb{E} \left\| \nabla \mathcal{L}(\mathbf{x}^t) - \sum_{i=1}^m \frac{1}{m} \mathbf{h}_i^{(t)} \right\|^2 \end{aligned} \quad (41)$$

$$\begin{aligned} & \leq -\|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + \tau\eta L \sum_{i=1}^m \frac{1}{m^2} [1 + 3(1 - p_i^t)^2] \sigma^2 \\ & \quad + \left(\frac{3}{2}\tau\eta L - \frac{1}{2}\right) \sum_{i=1}^m \frac{1}{m^2} (1 - p_i^t)^2 \mathbb{E} \|\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}\|^2 \\ & \quad + \frac{1}{2} \sum_{i=1}^m \frac{1}{m} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}^t) - \mathbf{h}_i^{(t)}\|^2, \end{aligned} \quad (42)$$

where the last inequality comes from Jensen's Inequality $\left\| \sum_{i=1}^m \frac{1}{m} \mathbf{w}_i \right\|^2 \leq \sum_{i=1}^m \frac{1}{m} \|\mathbf{w}_i\|^2$.

Bounding the difference between global gradient and normalized client gradient. According to the definition of $\mathbf{h}_i^{(t)}$ in (30) and the smoothness assumption, there is

$$\mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}^t) - \mathbf{h}_i^{(t)}\|^2$$

$$= \mathbb{E} \left\| \frac{1}{\tau} \sum_{k=1}^{\tau} [\nabla \mathcal{L}_i(\mathbf{x}^t) - \nabla \mathcal{L}_i(\mathbf{x}_{i,k}^t)] \right\|^2 \quad (43)$$

$$\leq \frac{1}{\tau} \sum_{k=1}^{\tau} \mathbb{E} \left\| \nabla \mathcal{L}_i(\mathbf{x}^t) - \nabla \mathcal{L}_i(\mathbf{x}_{i,k}^t) \right\|^2 \quad (44)$$

$$\leq \frac{L^2}{\tau} \sum_{k=1}^{\tau} \mathbb{E} \left\| \mathbf{x}^t - \mathbf{x}_{i,k}^t \right\|^2. \quad (45)$$

In terms of bounding the difference between the global model \mathbf{x}^t and the local model $\mathbf{x}_{i,k}^t$ after the k -th local iteration, we have the following results based on the model update rule as well as the assumption on the mask-induced error,

$$\begin{aligned} \mathbb{E} \left\| \mathbf{x}^t - \mathbf{x}_{i,k}^t \right\|^2 &= \eta^2 \mathbb{E} \left\| \sum_{s=0}^{k-1} g_i(\mathbf{x}_{i,s}^t \odot \mathcal{M}_i) \right\|^2 \\ &\leq 2\eta^2 \mathbb{E} \left\| \sum_{s=0}^{k-1} [g_i(\mathbf{x}_{i,s}^t \odot \mathcal{M}_i) - \nabla \mathcal{L}_i(\mathbf{x}_{i,s}^t \odot \mathcal{M}_i)] \right\|^2 \\ &\quad + 2\eta^2 \mathbb{E} \left\| \sum_{s=0}^{k-1} \nabla \mathcal{L}_i(\mathbf{x}_{i,s}^t \odot \mathcal{M}_i) \right\|^2 \end{aligned} \quad (46)$$

$$\leq 2\eta^2 \delta^2 \mathbb{E} \left\| \sum_{s=0}^{k-1} \nabla \mathcal{L}_i(\mathbf{x}, \mathbf{s}_i^t) \right\|^2 + 2\eta^2 \sigma^2 k \quad (47)$$

$$\leq 2\eta^2 k \sum_{s=0}^{\tau-1} \mathbb{E} \left\| \nabla \mathcal{L}_i(\mathbf{x}_{i,s}^t) \right\|^2 + 2\eta^2 \sigma^2 k. \quad (48)$$

Taking the average over τ local iterations, we get,

$$\frac{1}{\tau} \sum_{k=1}^{\tau} \mathbb{E} \left\| \mathbf{x}^t - \mathbf{x}_{i,k}^t \right\|^2 \quad (49)$$

$$\leq 2\eta^2 \sigma^2 (\tau - 1) + 2\eta^2 (\tau - 1) \sum_{k=0}^{\tau-1} \mathbb{E} \left\| \nabla \mathcal{L}_i(\mathbf{x}_{i,k}^t) \right\|^2. \quad (50)$$

Moreover, the local normalized gradient can be bounded by,

$$\begin{aligned} &\mathbb{E} \left\| \nabla \mathcal{L}_i(\mathbf{x}_{i,k}^t) \right\|^2 \\ &\leq 2\mathbb{E} \left\| \nabla \mathcal{L}_i(\mathbf{x}_{i,k}^t) - \nabla \mathcal{L}_i(\mathbf{x}^t) \right\|^2 + 2\mathbb{E} \left\| \nabla \mathcal{L}_i(\mathbf{x}^t) \right\|^2 \\ &\leq 2L^2 \mathbb{E} \left\| \mathbf{x}^t - \mathbf{x}_{i,k}^t \right\|^2 + 2\mathbb{E} \left\| \nabla \mathcal{L}_i(\mathbf{x}^t) \right\|^2. \end{aligned} \quad (51)$$

Plug (51) back into (50), we achieve,

$$\frac{1}{\tau} \sum_{k=1}^{\tau} \mathbb{E} \left\| \mathbf{x}^t - \mathbf{x}_{i,k}^t \right\|^2 \leq 2\eta^2 \sigma^2 (\tau - 1)$$

$$\begin{aligned}
& + 4\eta^2 L^2 (\tau - 1) \sum_{k=0}^{\tau-1} \mathbb{E} \|\mathbf{x}^t - \mathbf{x}_{i,k}^t\|^2 \\
& + 4\eta^2 (\tau - 1) \sum_{k=0}^{\tau-1} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}^t)\|^2.
\end{aligned} \tag{52}$$

After rearranging, the difference between the global gradient and client gradient can be bounded by

$$\begin{aligned}
\frac{1}{\tau} \sum_{k=1}^{\tau} \mathbb{E} \|\mathbf{x}^t - \mathbf{x}_{i,k}^t\|^2 & \leq \frac{2\eta^2 \sigma^2 (\tau - 1)}{1 - 4\eta^2 L^2 \tau (\tau - 1)} \\
& + \frac{4\eta^2 \tau (\tau - 1)}{1 - 4\eta^2 L^2 \tau (\tau - 1)} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}^t)\|^2.
\end{aligned} \tag{53}$$

Similarly as in [29] where this analysis framework is proposed, we define $\gamma = 4\eta^2 L^2 \tau (\tau - 1) \leq 1$, and then (53) can be simplified as,

$$\frac{L^2}{\tau} \sum_{k=1}^{\tau} \mathbb{E} \|\mathbf{x}^t - \mathbf{x}_{i,k}^t\|^2 \leq \frac{2\eta^2 \sigma^2 L^2 (\tau - 1)}{1 - \gamma} + \frac{\gamma}{1 - \gamma} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}^t)\|^2. \tag{54}$$

Taking the average across all clients, there is

$$\begin{aligned}
& \sum_{i=1}^m \frac{1}{m} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}^t) - \mathbf{h}_i^{(t)}\|^2 \\
& \leq \frac{2\eta^2 \sigma^2 L^2 (\tau - 1)}{1 - \gamma} + \frac{\gamma}{1 - \gamma} \sum_{i=1}^m \frac{1}{m} \mathbb{E} \|\nabla \mathcal{L}_i(\mathbf{x}^t)\|^2
\end{aligned} \tag{55}$$

$$\leq \frac{2\eta^2 \sigma^2 L^2 (\tau - 1)}{1 - \gamma} + \frac{\gamma \beta^2}{1 - \gamma} \mathbb{E} \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + \frac{\gamma \zeta^2}{1 - \gamma}. \tag{56}$$

Global Convergence Property. Based on the above analysis, we can bound the learning progress with (56) and (42),

$$\begin{aligned}
& \frac{1}{\tau \eta} (\mathbb{E} [\mathcal{L}(\mathbf{x}^{t+1})] - \mathcal{L}(\mathbf{x}^t)) \\
& \leq -\|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + \tau \eta L \sum_{i=1}^m \frac{1}{m^2} \left[1 + 3(1 - p_i^t)^2 \right] \sigma^2 \\
& + \frac{\gamma \beta^2}{2(1 - \gamma)} \mathbb{E} \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + \frac{\gamma \zeta^2}{2(1 - \gamma)} \\
& + \left(\frac{3}{2} \tau \eta L - \frac{1}{2} \right) \sum_{i=1}^m \frac{1}{m^2} (1 - p_i^t)^2 \mathbb{E} \|\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}\|^2 \\
& + \frac{\eta^2 \sigma^2 L^2 (\tau - 1)}{1 - \gamma}.
\end{aligned} \tag{57}$$

If $\gamma \leq \frac{1}{2\beta^2+1}$, then we have $\frac{1}{1-\gamma} \leq 1 + \frac{1}{2\beta^2}$ and $\frac{\gamma\beta^2}{1-\gamma} \leq \frac{1}{2}$. Therefore the above result can be simplified as,

$$\begin{aligned}
& \frac{1}{\tau\eta} (\mathbb{E} [\mathcal{L}(\mathbf{x}^{t+1})] - \mathcal{L}(\mathbf{x}^t)) \\
& \leq -\frac{3}{4} \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + 4\tau\eta L \sum_{i=1}^m \frac{1}{m^2} \sigma^2 \\
& \quad + \eta^2 \sigma^2 L^2 (\tau - 1) \left(1 + \frac{1}{2\beta^2}\right) \\
& \quad + \left(\frac{3}{2}\tau\eta L - \frac{1}{2}\right) \sum_{i=1}^m \frac{1}{m^2} (1 - p_i^t)^2 \mathbb{E} \|\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}\|^2 \\
& \quad + \left[2\eta^2 L^2 \tau (\tau - 1) \left(1 + \frac{1}{2\beta^2}\right)\right] \zeta^2. \tag{58}
\end{aligned}$$

Taking the average across all T communication rounds,

$$\begin{aligned}
& \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 \leq \frac{4[\mathcal{L}(\mathbf{x}^0) - \mathcal{L}_{\text{inf}}]}{3\eta\tau T} \\
& \quad + \frac{16\tau\eta L \sigma^2}{3m} + 2\eta^2 \sigma^2 L^2 (\tau - 1) + 4\eta^2 L^2 \tau (\tau - 1) \zeta^2 \\
& \quad + \left(2\tau\eta L - \frac{2}{3}\right) \sum_{i=1}^m \frac{1}{m^2} (1 - p_i^t)^2 \mathbb{E} \|\mathbf{h}_{i'}^{(t)} - \mathbf{h}_i^{(t)}\|^2.
\end{aligned}$$

The proof of Lemma 2 is completed.

C. Proof of Theorem 1

From (57), due to the bounded similarity assumption on sparse models, we have

$$\begin{aligned}
& \frac{1}{\tau\eta} (\mathbb{E} [\mathcal{L}(\mathbf{x}^{t+1})] - \mathcal{L}(\mathbf{x}^t)) \\
& \leq -\|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + \tau\eta L \sum_{i=1}^m \frac{1}{m^2} \left[1 + 3(1 - p_i^t)^2\right] \sigma^2 \\
& \quad + \left(\frac{3}{2}\tau\eta L - \frac{1}{2}\right) \beta^2 \mathbb{E} \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + \left(\frac{3}{2}\tau\eta L - \frac{1}{2}\right) \zeta^2 \\
& \quad + \frac{\eta^2 \sigma^2 L^2 (\tau - 1)}{1 - \gamma} \\
& \quad + \frac{\gamma\beta^2}{2(1 - \gamma)} \mathbb{E} \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 + \frac{\gamma\zeta^2}{2(1 - \gamma)}. \tag{59}
\end{aligned}$$

With the above constraint on the learning rate, we have,

$$\frac{1}{\tau\eta} (\mathbb{E} [\mathcal{L}(\mathbf{x}^{t+1})] - \mathcal{L}(\mathbf{x}^t))$$

$$\begin{aligned}
&\leq -\|\nabla\mathcal{L}(\mathbf{x}^t)\|^2 + 4\tau\eta L \sum_{i=1}^m \frac{1}{m^2}\sigma^2 \\
&+ \left[\frac{\gamma\beta^2}{2(1-\gamma)} + \frac{1}{4} \right] \mathbb{E}\|\nabla\mathcal{L}(\mathbf{x}^t)\|^2 + \left[\frac{\gamma}{2(1-\gamma)} + \frac{1}{4} \right] \zeta^2 \\
&+ \frac{\eta^2\sigma^2 L^2(\tau-1)}{1-\gamma}. \tag{60}
\end{aligned}$$

Similarly, with the same constraint on γ and β in (58), the above result can be simplified as,

$$\begin{aligned}
&\frac{1}{\tau\eta} (\mathbb{E}[\mathcal{L}(\mathbf{x}^{t+1})] - \mathcal{L}(\mathbf{x}^t)) \\
&\leq -\frac{1}{2}\|\nabla\mathcal{L}(\mathbf{x}^t)\|^2 + 4\tau\eta L \sum_{i=1}^m \frac{1}{m^2}\sigma^2 \\
&+ \eta^2\sigma^2 L^2(\tau-1) \left(1 + \frac{1}{2\beta^2}\right) \\
&+ \left[2\eta^2 L^2 \tau(\tau-1) \left(1 + \frac{1}{2\beta^2}\right) + \frac{1}{4} \right] \zeta^2 \tag{61}
\end{aligned}$$

$$\begin{aligned}
&\leq -\frac{1}{2}\|\nabla\mathcal{L}(\mathbf{x}^t)\|^2 + 4\tau\eta L \sum_{i=1}^m \frac{1}{m^2}\sigma^2 \\
&+ \frac{3}{2}\eta^2\sigma^2 L^2(\tau-1) + \left[3\eta^2 L^2 \tau(\tau-1) + \frac{1}{4} \right] \zeta^2. \tag{62}
\end{aligned}$$

Taking the average across all rounds,

$$\begin{aligned}
\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\|\nabla\mathcal{L}(\mathbf{x}^t)\|^2 &\leq \frac{2[\mathcal{L}(\mathbf{x}^0) - \mathcal{L}_{\text{inf}}]}{\eta\tau T} \\
&+ \frac{8\tau\eta L\sigma^2}{m} + 3\eta^2\sigma^2 L^2(\tau-1) + \left[6\eta^2 L^2 \tau(\tau-1) + \frac{1}{2} \right] \zeta^2. \tag{63}
\end{aligned}$$

For the ease of writing, we define $A = \tau$, $B = \tau - 1$ and $C = \tau(\tau - 1)$, and then we derive

$$\begin{aligned}
\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\|\nabla\mathcal{L}(\mathbf{x}^t)\|^2 &\leq \frac{2[\mathcal{L}(\mathbf{x}^0) - \mathcal{L}_{\text{inf}}]}{\eta\tau T} \\
&+ \frac{8\eta L\sigma^2 A}{m} + 3\eta^2\sigma^2 L^2 B + \left(6\eta^2 L^2 C + \frac{1}{2} \right) \zeta^2. \tag{64}
\end{aligned}$$

Since there is

$$\min_{t \in \mathbb{T}} \mathbb{E}\|\nabla\mathcal{L}(\mathbf{x}^t)\|^2 \leq \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\|\nabla\mathcal{L}(\mathbf{x}^t)\|^2, \tag{65}$$

it holds that,

$$\min_{t \in \mathbb{T}} \mathbb{E}\|\nabla\mathcal{L}(\mathbf{x}^t)\|^2 \leq \frac{2[\mathcal{L}(\mathbf{x}^0) - \mathcal{L}_{\text{inf}}]}{\eta\tau T}$$

$$+ \frac{8\eta L\sigma^2 A}{m} + 3\eta^2\sigma^2 L^2 B + \left(6\eta^2 L^2 C + \frac{1}{2}\right) \zeta^2. \quad (66)$$

By setting $\eta = \sqrt{\frac{m}{\tau T}}$, we have

$$\begin{aligned} \min_{t \in \mathbb{T}} \mathbb{E} \|\nabla \mathcal{L}(\mathbf{x}^t)\|^2 &\leq \mathcal{O}\left(\frac{1}{\sqrt{m\tau T}}\right) + \mathcal{O}\left(\frac{A\sigma^2}{\sqrt{m\tau T}}\right) \\ &+ \mathcal{O}\left(\frac{mB\sigma^2}{\tau T}\right) + \mathcal{O}\left(\frac{mC\zeta^2}{\tau T}\right). \end{aligned}$$

The proof of Theorem 1 is completed.

REFERENCES

- [1] Z. Qin, G. Y. Li, and H. Ye, “Federated learning and wireless communications,” *IEEE Wireless Communications*, vol. 28, no. 5, pp. 134–140, 2021.
- [2] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *International Conference on Artificial Intelligence and Statistics (AISTATS)*, Fort Lauderdale, USA, April 2017, pp. 1273–1282.
- [3] Q. Yang, Y. Liu, T. Chen, and Y. Tong, “Federated machine learning: Concept and applications,” *ACM Transactions on Intelligent Systems and Technology*, vol. 10, no. 2, pp. 1–19, 2019.
- [4] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, “Deep learning with differential privacy,” in *ACM SIGSAC Conference on Computer and Communications Security*, Vienna, Austria, May 2016, p. 308–318.
- [5] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, “Practical secure aggregation for privacy-preserving machine learning,” in *ACM SIGSAC Conference on Computer and Communications Security*, Dallas, Texas, USA, May 2017, p. 1175–1191.
- [6] Y. Liu, Y. Kang, C. Xing, T. Chen, and Q. Yang, “A secure federated transfer learning framework,” *IEEE Intelligent Systems*, vol. 35, no. 4, pp. 70–82, 2020.
- [7] X. Ma, M. Qin, F. Sun, Z. Hou, K. Yuan, Y. Xu, Y. Wang, Y.-K. Chen, R. Jin, and Y. Xie, “Effective model sparsification by scheduled grow-and-prune methods,” in *International Conference on Learning Representations*, April 2022.
- [8] B. Woodworth, K. K. Patel, S. Stich, Z. Dai, B. Bullins, B. McMahan, O. Shamir, and N. Srebro, “Is local SGD better than minibatch SGD?” in *International Conference on Machine Learning*, July 2020, pp. 10 334–10 343.
- [9] A. Koloskova, N. Loizou, S. Boreiri, M. Jaggi, and S. Stich, “A unified theory of decentralized SGD with changing topology and local updates,” in *International Conference on Machine Learning*, July 2020, pp. 5381–5393.
- [10] A. Khaled, K. Mishchenko, and P. Richtárik, “Tighter theory for local SGD on identical and heterogeneous data,” in *International Conference on Artificial Intelligence and Statistics*, 2020, pp. 4519–4529.
- [11] H. Ye, L. Liang, and G. Y. Li, “Decentralized federated learning with unreliable communications,” *IEEE Journal of Selected Topics in Signal Processing*, pp. 1–13, 2022.
- [12] Y. Fraboni, R. Vidal, L. Kamení, and M. Lorenzi, “Clustered sampling: Low-variance and improved representativity for clients selection in federated learning,” in *International Conference on Machine Learning*, July 2021, pp. 3407–3416.
- [13] H. Tanaka, D. Kunin, D. L. Yamins, and S. Ganguli, “Pruning neural networks without any data by iteratively conserving synaptic flow,” in *Advances in Neural Information Processing Systems*, vol. 33, December 2020, pp. 6377–6389.
- [14] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic, “QSGD: Communication-efficient SGD via gradient quantization and encoding,” in *Advances in Neural Information Processing Systems*, vol. 30, CA, USA, December 2017, pp. 1709–1720.

- [15] D. Jhunjhunwala, A. Gadhikar, G. Joshi, and Y. C. Eldar, "Adaptive quantization of model updates for communication-efficient federated learning," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, June 2021, pp. 3110–3114.
- [16] J. Bernstein, Y.-X. Wang, K. Azizzadenesheli, and A. Anandkumar, "signSGD: Compressed optimisation for non-convex problems," in *International Conference on Machine Learning*, Stockholm, Sweden, July 2018, pp. 560–569.
- [17] N. Strom, "Scalable distributed dnn training using commodity gpu cloud computing," in *Annual Conference of the International Speech Communication Association*, Dresden, Germany, September 2015.
- [18] J. Wangni, J. Wang, J. Liu, and T. Zhang, "Gradient sparsification for communication-efficient distributed optimization," in *Advances in Neural Information Processing Systems*, vol. 31, Montreal, CANADA, December 2018, pp. 1299–1309.
- [19] S. U. Stich, J.-B. Cordonnier, and M. Jaggi, "Sparsified SGD with memory," in *Advances in Neural Information Processing Systems*, vol. 31, Montreal, CANADA, December 2018, pp. 1–12.
- [20] S. P. Karimireddy, Q. Rebjock, S. Stich, and M. Jaggi, "Error feedback fixes signsgd and other gradient compression schemes," in *International Conference on Machine Learning*, California, USA, June 2019, pp. 3252–3261.
- [21] D. C. Mocanu, E. Mocanu, P. Stone, P. H. Nguyen, M. Gibescu, and A. Liotta, "Scalable training of artificial neural networks with adaptive sparse connectivity inspired by network science," *Nature Communications*, vol. 9, no. 1, pp. 1–12, 2018.
- [22] N. Lee, T. Ajanthan, and P. Torr, "Snip: Single-shot network pruning based on connection sensitivity," in *International Conference on Learning Representations*, New Orleans, USA, May 2019.
- [23] T. Dettmers and L. Zettlemoyer, "Sparse networks from scratch: Faster training without losing performance," *arXiv preprint arXiv:1907.04840*, 2019.
- [24] H. You, C. Li, P. Xu, Y. Fu, Y. Wang, X. Chen, R. G. Baraniuk, Z. Wang, and Y. Lin, "Drawing early-bird tickets: Towards more efficient training of deep networks," *arXiv preprint arXiv:1909.11957*, 2019.
- [25] U. Evci, T. Gale, J. Menick, P. S. Castro, and E. Elsen, "Rigging the lottery: Making all tickets winners," in *International Conference on Machine Learning*, July 2020, pp. 2943–2952.
- [26] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "Scaffold: Stochastic controlled averaging for federated learning," in *International Conference on Machine Learning*, July 2020, pp. 5132–5143.
- [27] S. J. Reddi, J. Konečný, P. Richtárik, B. Póczós, and A. Smola, "Aide: Fast and communication efficient distributed optimization," *arXiv preprint arXiv:1608.06879*, 2016.
- [28] T. Murata and T. Suzuki, "Bias-variance reduced local SGD for less heterogeneous federated learning," *arXiv preprint arXiv:2102.03198*, 2021.
- [29] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, "Tackling the objective inconsistency problem in heterogeneous federated optimization," in *Advances in Neural Information Processing Systems*, vol. 33, Vancouver, Canada, Dec. 2020, pp. 7611–7623.
- [30] H. You, C. Li, P. Xu, Y. Fu, Y. Wang, X. Chen, Z. Wang, R. G. Baraniuk, and Y. Lin, "Drawing early-bird tickets: Towards more efficient training of deep networks," in *International Conference on Learning Representations*, Addis Ababa, Ethiopia, April 2020.
- [31] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, Las Vegas, NV, USA, June 2016, pp. 770–778.
- [32] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations*, San Diego, CA, USA, May 2015.
- [33] C. Wang, G. Zhang, and R. Grosse, "Picking winning tickets before training by preserving gradient flow," in *International Conference on Learning Representations*, New Orleans, USA, May 2019.

- [34] E. Diao, J. Ding, and V. Tarokh, "Heterofl: Computation and communication efficient federated learning for heterogeneous clients," in *International Conference on Learning Representations*, 2020.