

# Design and Analysis of Uplink and Downlink Communications for Federated Learning

Sihui Zheng    Cong Shen    Xiang Chen

## Abstract

Communication has been known to be one of the primary bottlenecks of federated learning (FL), and yet existing studies have not addressed the efficient communication design, particularly in wireless FL where both uplink and downlink communications have to be considered. In this paper, we focus on the design and analysis of physical layer quantization and transmission methods for wireless FL. We answer the question of *what* and *how* to communicate between clients and the parameter server and evaluate the impact of the various quantization and transmission options of the updated model on the learning performance. We provide new convergence analysis of the well-known FEDAVG under non-i.i.d. dataset distributions, partial clients participation, and finite-precision quantization in uplink and downlink communications. These analyses reveal that, in order to achieve an  $\mathcal{O}(1/T)$  convergence rate with quantization, transmitting the weight requires increasing the quantization level at a *logarithmic* rate, while transmitting the weight differential can keep a constant quantization level. Comprehensive numerical evaluation on various real-world datasets reveals that the benefit of a FL-tailored uplink and downlink communication design is enormous – a carefully designed quantization and transmission achieves more than 98% of the floating-point baseline accuracy with fewer than 10% of the baseline bandwidth, for majority of the experiments on both i.i.d. and non-i.i.d. datasets. In particular, 1-bit quantization (3.1% of the floating-point baseline bandwidth) achieves 99.8% of the floating-point baseline accuracy at almost the same convergence rate on MNIST, representing the best known bandwidth-accuracy tradeoff to the best of the authors’ knowledge.

## Index Terms

Wireless federated learning; Convergence analysis; Communication design.

S. Zheng and X. Chen are with School of Electronics and Information Technology, Sun Yat-sen University, China. (e-mail: zhengsh28@mail2.sysu.edu.cn; chenxiang@mail.sysu.edu.cn).

C. Shen is with the Charles L. Brown Department of Electrical and Computer Engineering, University of Virginia, USA. (e-mail: cong@virginia.edu).

## I. INTRODUCTION

Wireless federated learning (FL) [1], [2] is an emerging edge artificial intelligence framework [3]. FL has many attractive properties that cater to the growing trend of how data is generated and how machine learning (ML) model is trained. Empowered by the growing storage and computational capabilities of mobile devices and motivated by the increasing concern over transmitting private information to a central server, FL has become an attractive ML paradigm that trains ML models locally on each device where data never leaves the device [4], [5].

While FL offers many important benefits, it also faces several critical challenges including but not limited to significant communication cost, handling client heterogeneity (both dataset and the computation and communication capabilities) and the straggler problem, preserving the privacy of user data, improving robustness to adversarial attacks and failures, and ensuring fairness. A comprehensive review of these challenges can be found in [6]. In particular, despite being recognized as one of the primary bottlenecks of FL [2], [6], [7], research on the *communication* aspect in the FL pipeline has not been on par with the *learning* component, particularly in a wireless environment. Early research on communication-efficient FL largely focuses on reducing the number of communication rounds and the amount of information for communication, while assuming that the underlying communication “tunnel” has been established by existing wireless protocols. More recent research starts to fill this void from a communication and signal processing point of view. In general, the principle is to balance learning performance and communication efficiency via, e.g., device selection, bandwidth allocation, and power control; see Section II for an overview. There are also recent studies that focus on the communication system design [8]–[10], but they are either system-specific (e.g., cellular networks) or with high complexity beyond the current implementation capability (e.g., very high dimensional vector quantization).

While the early studies provide a glimpse of the potential of optimizing communication for learning, the actual implementation of the communication algorithms has not been tailored to the unique characteristics of FL. In particular, it is often taken for granted that standard signal processing and communication techniques can be directly applied to FL. We show in this paper that this can be highly suboptimal because they are mostly designed for independent and identically distributed (i.i.d.) sources over time, while the communicated model update in FL represents a long-term process consisting of many progressive learning rounds that collectively determine the final learning outcome. This phenomenon is known in the machine learning literature and has been leveraged to optimize the learning hyperparameters, e.g., decaying the learning rate over time [11], but has not been considered for the communication algorithms. To further complicate the matter, the overall FL performance is determined by both local model weight

(i.e., parameters of the ML model) upload and global model weight download over multiple learning rounds, suggesting that both uplink and downlink communications have to be considered.

In this paper, we study FL-tailored communication designs for training ML models locally at mobile devices and aggregation at the base station, where the information for communication (both uplink and downlink) is the model weight (or its update) after each learning round. The design goal is to maximize the learning accuracy and convergence rate, two prime objectives in FL. We answer the questions of *what* and *how* to transmit the updated model in each round between clients (mobile devices) and parameter server (base station), and study practical quantization and transmission methods that leverage the inherent structure of the machine learning model. Our main contributions are as follows.

- 1) We study practical quantization schemes for FL and show that the dynamic range of the weight needs to be taken into account, and the choice of rounding has a profound impact on the performance. For uplink, we demonstrate that transmitting only the weight differential is beneficial if the practical constraint allows, while pointing out that this differential transmission cannot be utilized for downlink when only *partial* clients participate each round. We also propose an enhancement called *layered quantization* for downlink, in which the quantization gain is adjusted to match the dynamic range of the weights in different network layers.
- 2) We rigorously prove convergence rate upper bounds of the well-known FEDAVG [4] under finite-precision quantization in both uplink and downlink communications. The theoretical analysis reveals a novel conclusion: in order to maintain the  $\mathcal{O}(1/T)$  convergence rate of the floating-point FEDAVG, the uplink or downlink quantization for direct weight transmission should increase the quantization precision at a *logarithmic* rate  $\mathcal{O}(\log(t))$ , while transmitting the weight differential can maintain a *constant* (i.e.,  $\mathcal{O}(1)$ ) quantization precision throughout the FL process. This result holds for non-i.i.d. dataset and partial (randomly selected) clients participation in each learning round, which is more general and matches the unique characteristics of FL [4].
- 3) Comprehensive numerical evaluation on four widely adopted datasets with increasing learning difficulties, *MNIST*, *F-EMNIST*, *CIFAR-10* and *Shakespeare* are done. We design a series of experiments to show the impact of each step in the quantization including quantization gain, rounding method and the relationship between the quantization and hyperparameters like batch size, local epoch, etc. In particular, we corroborate the theoretical conclusion that the quantization precision needs to increase at a logarithmic rate for direct weight transmission via numerical experiments. The results also reveal that the benefit of a FL-tailored uplink and downlink communication design is significant. In majority of the experiments, we see that a carefully designed quantization and transmission achieves more than

98% of the floating-point baseline accuracy with fewer than 10% of the baseline bandwidth, for both i.i.d. and non-i.i.d. datasets. As a final exclamation point, a 1-bit quantization (3.1% of the floating-point baseline bandwidth) achieves 99.8% of the floating-point baseline accuracy at almost the same convergence rate in the MNIST experiment, representing the best known bandwidth-accuracy tradeoff to the best of the authors' knowledge.

The rest of this paper is organized as follows. A brief overview of the related literature is provided in Section II. Section III describes the wireless FL system model. Uplink and downlink communication designs, together with the theoretical convergence analyses, are presented in Section IV and V, respectively. Experimental results are reported in Section VI. Section VII concludes the paper, and the technical proofs of the main theorems are provided in the Appendices.

## II. RELATED WORKS

Federated learning [4] is an emerging distributed machine learning [12] paradigm that addresses several new features created by modern ML applications. It has been extensively studied in recent years in the machine learning community, which aims to address various questions around improving machine learning efficiency and effectiveness [11], [13]–[15], preserving the privacy of user data [16]–[18], robustness to attacks and failures [19], [20], and ensuring fairness and addressing sources of bias [21], [22]. However, these works mostly focus on the machine learning aspect of FL and largely consider over-simplified communication models.

Recently, researchers have started looking into the communication design of FL, particularly the communication algorithms, protocols, and systems. Joint radio and computation resource management is another active research topic. Existing works [23] study the inherent trade-off between local model update and global model aggregation, to optimize over transmission power/rate and training time. To enable FL at scale and address the straggler problem, client selection is essential. In this regard, various joint radio resource allocation and client selection policies [24]–[28] have been proposed to minimize the learning loss or the training time.

Communication-efficient design has been another active research topic in FL [6], where the attempts have largely focused on either reducing the total number of communication rounds, or reducing the size of the exchanged messages in each round. One of the representative approaches for reducing the communication rounds is FEDAVG [4], which allows periodic model aggregation and local model updates and thus enables flexible communication-computation tradeoff [29]. Theoretical understanding of this tradeoff has been an active research area and, depending on the underlying assumptions (e.g., i.i.d. or non-i.i.d. local datasets, convex or non-convex loss functions, gradient descent or stochastic gradient

descent), rigorous analysis of the convergence behavior has been carried out [11], [13], [14], [30]. For model compression, general discussions on sparsification, subsampling, and quantization are given in [5]. Particularly, sparsification methods reduce the number of non-zero entries in the stochastic gradient [31]. Structured and sketched updates are proposed in [32] to reduce the size of model updates, which are further extended by lossy compression and federated dropout [33]. There have been recent efforts in developing quantization and source coding to reduce the communication cost [10], [31], [34]–[39]. However, most of the quantizers studied in these papers do not consider practical constraints and are not widely used in practice. Reference [37] only considers i.i.d. datasets and uplink quantization, and [39] focuses on downlink quantization of the model differential and does not apply to partial clients participation, which is an important feature of FL.

### III. SYSTEM MODEL

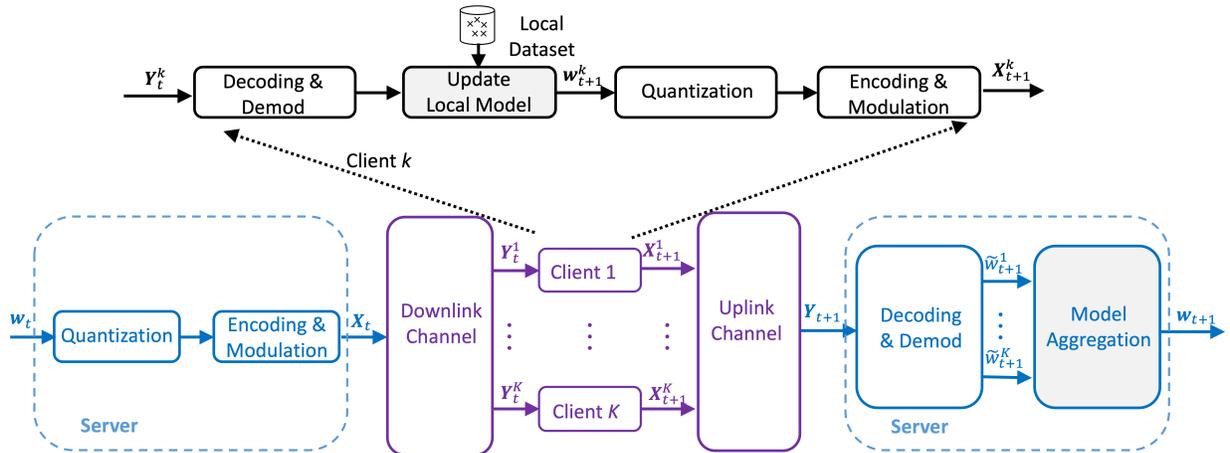


Fig. 1. Wireless FL system model. The  $t$ th to the  $(t + 1)$ th round of operations at both clients (mobile devices, shown in black) and server (BS, shown in blue) are illustrated. The shaded boxes correspond to the learning operations and others refer to communications.

The wireless federated learning system is illustrated in Fig. 1. We assume a federated learning task of collaboratively training a ML model (e.g., logistic regression or deep neural network (DNN)) as in [4]. In particular, there is a central parameter server (e.g., base station) and a set of at most  $N$  clients (e.g., mobile devices). Client  $k$  stores a (disjoint) local dataset  $\mathcal{D}_k = \{\mathbf{z}_i\}_{i=1}^{D_k}$ , with its size denoted by  $D_k$ , that never leaves the device. Datasets across devices are assumed to be non-i.i.d., which is an important feature of FL [4], [7]. The maximum data size when all devices participate in FL is  $D = \sum_{k=1}^N D_k$ . Each data sample  $\mathbf{z}$  is given as an input-output pair  $\{\mathbf{x}, y\}$ . The loss function  $f(\mathbf{w}, \mathbf{z})$  measures how well

a ML model with parameter  $\mathbf{w} \in \mathbb{R}^d$  fits one particular data sample  $\mathbf{z}$ . For the  $k$ th device, its local loss function  $F_k(\cdot)$  is defined by

$$F_k(\mathbf{w}) \triangleq \frac{1}{D_k} \sum_{\mathbf{z} \in \mathcal{D}_k} f(\mathbf{w}, \mathbf{z}).$$

Then, the global optimization objective over all  $N$  clients is given by

$$F(\mathbf{w}) \triangleq \sum_{k=1}^N \frac{D_k}{D} F_k(\mathbf{w}) = \frac{1}{D} \sum_{k=1}^N \sum_{\mathbf{z} \in \mathcal{D}_k} f(\mathbf{w}, \mathbf{z}). \quad (1)$$

The global loss function measures how well the model fits the entire corpus of data on average. As a result, the objective is to find the best model parameter  $\mathbf{w}^*$  that minimizes the global loss function:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\arg \min} F(\mathbf{w}).$$

Let  $F^*$  and  $F_k^*$  be the minimum value of  $F$  and  $F_k$ , respectively. Then,  $\Gamma = F^* - \frac{1}{N} \sum_{k=1}^N F_k^*$  quantifies the degree of non-i.i.d. [11]. We note that using  $\Gamma$  to measure the degree of non-i.i.d. is more meaningful when the dataset size is large, in which the minimum loss function values of  $F_k$  approach the true expected minimum loss function values with respect to the individual dataset distributions.

This work considers a generic FL framework where partial client participation and non-i.i.d. local datasets, two critical features that separate FL from distributed SGD, are explicitly captured. More specifically, the FL pipeline works by iteratively executing the following steps at the  $t$ th learning round.

- 1) **Downlink communication for model download.** The centralized server broadcasts the current global ML model, which is described by the latest weight vector  $\mathbf{w}_t$ , to a set of randomly selected clients denoted as  $\mathcal{S}_t$  with  $|\mathcal{S}_t| = K$ . The detailed communication mechanism for this phase will be described in Section V.
- 2) **Local computation.** Each client uses its local data to train a local ML model improved upon the received global ML model. In this work, we assume that mini-batch stochastic gradient descent (SGD) is used in training, where the weight  $\mathbf{w}_t^k$  is updated iteratively (for  $E$  steps in the current learning round) at device  $k$  as:

$$\begin{aligned} \mathbf{w}_{t,0}^k &= \mathbf{w}_t^k, \\ \mathbf{w}_{t,\tau}^k &= \mathbf{w}_{t,\tau-1}^k - \eta_t \nabla F(\mathbf{w}_{t,\tau-1}^k, \xi_\tau^k), \quad \forall \tau = 1, \dots, E, \\ \mathbf{w}_{t+1}^k &= \mathbf{w}_{t,E}^k, \end{aligned}$$

where  $\xi_\tau^k$  is a mini batch of data points that are independently sampled uniformly at random from the local dataset of client  $k$ .

- 3) **Uplink communication for model upload.** The selected  $K$  devices upload their latest local models to the server synchronously. The communication mechanism for this phase will be described in Section IV.
- 4) **Global aggregation.** The server aggregates the received local models to generate a new global ML model:

$$\mathbf{w}_{t+1} = \sum_{k \in \mathcal{S}_t} \frac{D_k}{\sum_{i \in \mathcal{S}_t} D_i} \mathbf{w}_{t+1}^k. \quad (2)$$

The server then moves on to the  $(t + 1)$ th round. The process completes after  $T$  rounds.

By and large, the above process is followed by majority of the existing FL formulations. There are some variants, such as adapting the client selection [40], allowing for varying number of local updates [6], or improving the model learning by distributed primal-dual methods [41]. Our work nevertheless focuses on the communication aspect (both uplink and downlink) of FL and can incorporate these enhancements.

#### IV. UPLINK COMMUNICATION DESIGN

The task of a particular uplink communication round (e.g.  $t$ th round) is to deliver  $\mathbf{w}_t^k$  for client  $k$  to the BS as accurately and efficiently as possible. However, since the FL process involves  $T$  rounds of model update/communication operations, which are inherently correlated over time, there exist opportunities to improve the communication design. In particular, uplink design involves answering two questions from the communication perspective: *what* to transmit, and *how* to transmit.

##### A. What to Transmit: Weight versus Weight Differential

If we treat the design of the  $t$ th uplink communication round at client  $k$  as an isolated task, i.e., we ignore the operations in the past both at client  $k$  and at the server, we can directly transmit the latest local weight vector  $\mathbf{w}_t^k$ . A different choice, which leverages the past information, is to transmit the *weight update* (also called *weight differential* in this paper)  $\mathbf{d}_t^k = \mathbf{w}_t^k - \mathbf{w}_{t-1}^k$  as opposed to the weight itself.

From a pure learning perspective, there is no difference whether the updated model itself ( $\mathbf{w}_t^k$ ) or its differential ( $\mathbf{d}_t^k$ ) is communicated from clients to the server. As long as the server can reconstruct  $\mathbf{w}_t^k$ , this aspect does not impact the learning performance [4]. Thus, it seems that the choice is insignificant and boils down to other practical considerations. For example, transmitting weight differential relies on the server keeping the previous global model  $\mathbf{w}_{t-1}$ , so that the new local models can be reconstructed from the differential. This however may not always be true if the server deletes intermediate model aggregation for privacy preservation [21], which makes reconstruction from the model differential infeasible. As another example, transmitting weight differential implicitly assumes  $\mathbf{w}_{t-1}^k = \mathbf{w}_{t-1}$ , i.e., the previously received

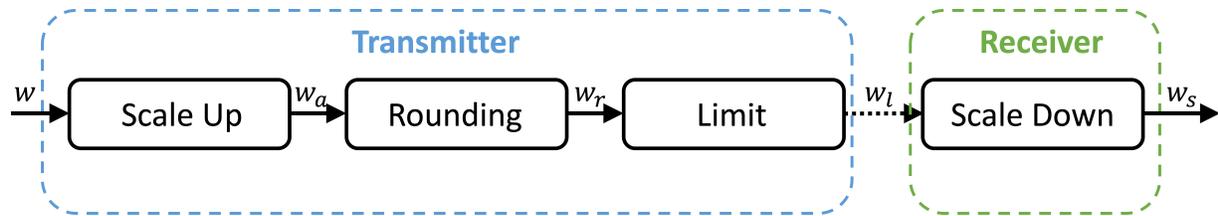


Fig. 2. Illustration of the adopted quantization structure.

global model is accurate. This may not be true in a practical communication system where decoding error is inevitable. In both examples, transmitting the weight vector itself is more preferable.

However, in addition to these considerations, we show in Section IV-C that the choice between weight or weight differential in the uplink communication phase has a more profound impact to the learning performance (in particular the convergence), when imperfect reconstruction due to quantization is captured.

### B. How to Transmit: Quantization Designs

In order to transmit the  $d$ -dimensional source message (either weight or weight differential) to the server, we first quantize the source vector into discrete values, and then apply coding (both source and channel) and modulation in the baseband, as shown in Fig. 1. The coding and modulation operation can leverage existing wireless system designs [42] and is not the focus of this paper. The quantization method, however, bears some consideration as discussed below.

In ML models such as deep neural networks (DNN), weights are usually represented in floating point format<sup>1</sup>. A *quantizer* is designed to reduce the necessary bit-width for each weight and hence decrease the message size for communications. It is worth noting that the quantization design in FL is very different to DNN model compression [43]–[45], which focuses on reducing model storage and simplifying inference computation. Also for DNN model compression, the impact of quantization is reflected in the final model after the training is complete. Quantization design for FL, on the other hand, aims at reducing the communication bandwidth, and has to be carried out in every round such that all the quantizations collectively affect the performance of FL.

We focus on a practical quantizer design that is suitable for communication system implementations. For this reason, we do not consider vector quantization [10] which is highly complex and not well used in practice despite its theoretical advantages. The adopted quantizer design is illustrated in Fig. 2. We note

<sup>1</sup>A 32-bit representation is a common choice in practice.

that this diagram is different from the existing literature [9], [31], [35], [37], which only has the *rounding* and *limit* operations. Inspired by the classical quantization methods in communication system [46], we add the *scaling up* and *scaling down* steps. This is because the dynamic range of the weights may not match a pre-determined rounding strategy, and a proper quantization gain control<sup>2</sup> is often applied to handle this issue. Specifically, for a full-precision weight  $w$ , the quantizer output  $Q(w)$  can be obtained via the following steps:

- 1) **Scale Up.**  $w$  is first amplified with a scaling factor called *quantization gain*. Denoting the quantization gain as  $G$ , the amplified value is  $w_a = wG$ .
- 2) **Rounding.**  $w_a$  is truncated to only retain its integer part:  $w_r = R(w_a)$  where  $R(\cdot)$  denotes the rounding function.
- 3) **Limit.** The range of integer  $w_r$  is further limited to  $B$  bits:

$$w_l = \begin{cases} 2^{B-1} - 1 & \text{if } w_r > 2^{B-1} - 1, \\ w_r & \text{if } w_r \in [-2^{B-1}, 2^{B-1} - 1], \\ -2^{B-1} & \text{if } w_r < -2^{B-1}. \end{cases}$$

- 4) **Scale Down.** The receiver output  $w_s$  is obtained by scaling down  $w_l$ :  $w_s = w_l/G$ .

We now detail how different components are designed for a  $B$ -bit quantizer as in Fig. 2.

**Quantization Gain.** A large  $G$  preserves more decimal digits of  $w$  and hence improves the representation accuracy, but it also increases the percentage of overflow in the subsequent limit operation, which introduces quantization errors in a different way. It is worth mentioning that  $G$  is typically set as power of 2, which simplifies the implementation to bit shifting.

**Quantizer Structure.** For comparison, we consider two quantizer structures in this work. In *Native Quantization* (NQ), the scaling up is limited to  $G = 2^{B-1}$  (1 bit for sign and the rest for decimal), and thus the scaling down step can be done at the transmitter, which means the receiver does not need to know  $G$ . An alternative structure *Tuned Quantization* (TQ) allows for fine-tuning  $G$  to a more suitable value (usually greater than  $2^{B-1}$ ) but requires that the scaling down step be done at the receiver.

**Rounding method.** Two rounding functions are considered. The basic one is *nearest rounding* (NR):

$$R(x) = \begin{cases} \lfloor x \rfloor & \text{if } x - \lfloor x \rfloor < 0.5 \\ \lfloor x \rfloor + 1 & \text{otherwise} \end{cases}$$

<sup>2</sup>This can be implemented by the automatic gain control (AGC) module in the wireless transmitter, which is usually enforced before the analog-to-digital converter (ADC) so that the input signal can match the dynamic range of the ADC.

where  $\lfloor x \rfloor$  is the floor of  $x$ . The second method is *stochastic rounding* (SR) [43], which rounds  $x$  to  $\lfloor x \rfloor$  with a probability proportional to the proximity of  $x$  to  $\lfloor x \rfloor$  (w.p. is short for ‘with probability’):

$$R(x) = \begin{cases} \lfloor x \rfloor & \text{w.p. } 1 - (x - \lfloor x \rfloor) \\ \lfloor x \rfloor + 1 & \text{w.p. } x - \lfloor x \rfloor. \end{cases}$$

**Enhanced 1-bit quantizer.** For the special case of a 1-bit quantizer, the quantization operation can be simplified as following two steps, without following the scale up – rounding and limit– scale down operations. In particular, we first round  $w$  with either NR or SR as follows.

- Nearest Rounding:

$$R(w) = \begin{cases} +1, & \text{if } w \geq 0, \\ -1, & \text{if } w < 0. \end{cases}$$

- Stochastic Rounding:

$$R(w) = \begin{cases} +1, & \text{w.p. Pr,} \\ -1, & \text{w.p. (1 - Pr),} \end{cases}$$

where  $\text{Pr} = \min(1, \max(0, \frac{w+1/G}{2/G}))$ .

Then, the receiver performs scale down to get  $Q(w) = R(w)/G$ .

### C. Convergence Analysis for FEDAVG with Uplink Quantization

As stated in Section IV-A, both the weight itself  $\mathbf{w}_t^k$  and weight differential  $\mathbf{d}_t^k$  can be used for uplink communication. However, this section shows that the two options have very different convergence behaviors, which lead to different requirements on quantization.

1) *Analysis for weight transmission:* We first analyze directly transmitting weight  $\mathbf{w}_t^k$  in the uplink of FEDAVG with quantization. The main theoretical result is that this configuration converges to the global optimum at a rate of  $\mathcal{O}(\frac{1}{T})$ , which is the same scaling behavior of the vanilla FEDAVG, *if we gradually increase the quantization precision over  $t$ .*

To simplify the analysis, we assume in the remainder of the paper that the local dataset sizes at all devices are the same:  $D_i = D_j, \forall i, j \in [N]$ , and focus on the general case of randomly selected  $K$  out of  $N$  clients participating in the server aggregation with non-i.i.d. dataset<sup>3</sup>. Let the set  $\mathcal{S}_t \subset [N]$  denote

<sup>3</sup>As will become clear after Section V, the case of unbalanced datasets can be easily incorporated in the analysis of both uplink and downlink communications when there is *full* clients participation in FL. However, when combined with partial clients participation, the analysis of unbalanced dataset becomes nontrivial. In this case, the coefficients in Eqn. (2) vary in each round, which makes the random sampling of clients no longer an unbiased estimation of the full participation case. Also, the error of SGD and quantization becomes much more complex, since it not only depends on how many clients are selected but also on which clients are selected. We leave this case for future research.

the  $K$  randomly selected clients in the  $t$ th round. With quantization, these devices transmit  $\{Q(\mathbf{w}_t^k)\}_{k=1}^K$  in the uplink, and the server performs aggregation as

$$\mathbf{w}_t = \frac{1}{K} \sum_{k \in \mathcal{S}_t} Q(\mathbf{w}_t^k). \quad (3)$$

The following assumptions are made for the convergence analysis. Assumption 1 is fairly standard and has been widely used in the convergence analysis of FEDAVG; see [11], [13], [15], [37]. Assumption 2 simply upper bounds the largest value of the weight so that the error in quantization is bounded. In practice, this assumption almost always holds because of the limited bit-width of weights in storage and computation.

**Assumption 1.** 1) *L-smooth*:  $\forall \mathbf{v}$  and  $\mathbf{w}$ ,  $F_k(\mathbf{v}) \leq F_k(\mathbf{w}) + (\mathbf{v} - \mathbf{w})^T \nabla F_k(\mathbf{w}) + \frac{L}{2} \|\mathbf{v} - \mathbf{w}\|^2$ .

2)  *$\mu$ -strongly convex*:  $\forall \mathbf{v}$  and  $\mathbf{w}$ ,  $F_k(\mathbf{v}) \geq F_k(\mathbf{w}) + (\mathbf{v} - \mathbf{w})^T \nabla F_k(\mathbf{w}) + \frac{\mu}{2} \|\mathbf{v} - \mathbf{w}\|^2$ .

3) *Bounded variance for mini-batch SGD*: The variance of stochastic gradients satisfies:

$$\mathbb{E} \left\| \nabla F_k(\mathbf{w}_t^k, \xi_t^k) - \nabla F_k(\mathbf{w}_t^k) \right\|^2 \leq \sigma_k^2$$

for  $k = 1, \dots, N$ .

4) *Uniformly bounded gradient*:  $\mathbb{E} \left\| \nabla F_k(\mathbf{w}_t^k, \xi_t^k) \right\|^2 \leq H^2$  for all  $k = 1, \dots, N$ .

**Assumption 2.**  $\max_{k \in [N], t \in [T]} \|\mathbf{w}_t^k\|_\infty \leq M$ , for constant  $M \geq 0$ .

**Theorem 1.** Define  $\kappa = \frac{L}{\mu}$ ,  $\gamma = \max\{8\kappa, E\}$ . Choose learning rate  $\eta_t = \frac{2}{\mu(\gamma+t)}$  and quantization level  $B_t = \log_2 \left[ \frac{\mu(\gamma+t-1)}{2} + 1 \right]$ . Then, under Assumptions 1 and 2 and using stochastic rounding with quantization gain  $G = \frac{2^{B_t-1}}{M}$  on weight  $\mathbf{w}_t^k$ , the convergence of FEDAVG with non-i.i.d. local datasets and partial clients participation satisfies

$$\mathbb{E} [F(\mathbf{w}_T)] - F^* \leq \frac{2\kappa}{\gamma+T} \left[ \frac{D}{\mu} + \left( 2L + \frac{E\mu}{4} \right) \|\mathbf{w}_0 - \mathbf{w}^*\|^2 \right], \quad (4)$$

where the constant  $D$  is

$$D = \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 + \frac{N-K}{N-1} \frac{4}{K} E^2 H^2 + \frac{dM^2}{K}. \quad (5)$$

The complete proof of Theorem 1 can be found in Appendix A. We note that the expectation in Eqn. (4) is with respect to three random events: (a) stochastic gradient when updating the model; (b) stochastic rounding in quantization; and (c) random sampling when selecting clients in each round.

2) *Analysis for weight differential transmission*: We call the communication design for weight differential as *Differential Transmission* (DT). With quantized weight differentials, the  $K$  randomly selected

devices transmit  $\{Q(\mathbf{d}_t^k)\}_{k=1}^K$  in the uplink, and the server performs aggregation as

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \frac{1}{K} \sum_{k \in \mathcal{S}_t} Q(\mathbf{d}_t^k). \quad (6)$$

Intuitively, the global aggregation in Eqn. (6) may be better than Eqn. (3) for a given quantizer design as in Fig. 2. As stated previously, we have to strike a balance between representation range and accuracy when selecting the value of  $G$ . The range of the weight differential  $\mathbf{d}_t^k$  is typically smaller than the raw weight  $\mathbf{w}_t^k$ , in particular towards convergence. Hence, a larger  $G$  can be used for DT to achieve higher quantization precision while avoiding excessive overflow. This can also be interpreted as not wasting bits on the constant part of the weights, which improves communication efficiency.

In addition to the advantage in quantization precision, Theorem 2 shows that DT can converge to the global optimum at rate of  $\mathcal{O}(\frac{1}{T})$  without requiring an increasing quantization level, which is better than Theorem 1. The complete proof of Theorem 2 can be found in Appendix B.

**Theorem 2.** *Let Assumption 1 hold and  $\kappa, \gamma, \eta_t$  be defined in Theorem 1. Let  $B$  be a fixed quantization level. Using stochastic rounding on weight differential  $\mathbf{d}_t^k$ , the convergence bound in Eqn. (4) for FEDAVG with non-i.i.d. local datasets, partial clients participation, and uplink quantization still holds, with  $D$  being replaced by*

$$D = \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 + \frac{N-K}{N-1} \frac{4}{K} E^2 H^2 + \frac{4d}{K(2^B-1)^2} E^2 H^2.$$

Compared to the known convergence results without quantization [11], [13], [14], Theorems 1 and 2 state that the same convergence rate can be largely preserved if the quantization is carefully designed. Intuitively, errors introduced by uplink quantizations may be accumulated and reflected in the new global model, which is then used by clients for the next round of training. This could potentially lead to error propagation over rounds and affect the convergence of FL. The core idea behind these theorems, especially reflected in their proofs, is that the quantizer design should ensure the errors introduced by quantization are well controlled at a lower level comparing to the noise of SGD, such that the overall “noise level” is not increased and thus the convergence of SGD is not violated.

We can also see from Theorems 1 and 2 that the convergence bounds have certain monotonic relationships with several hyperparameters. The bounds increase with  $E$ , which is consistent with the result of [47]. Larger  $B$  and  $K$  reduce the bounds, which is intuitive. Finally, we note that the effect of non-i.i.d. datasets, which is captured by  $\Gamma$ , is reflected in both Theorem 1 and 2. Furthermore, the convergence upper bounds of these theorems are monotonically decreasing when  $\Gamma$  reduces. When  $\Gamma$  goes to zero as the dataset size increases asymptotically, we have the i.i.d. dataset and the convergence upper bounds

have the best results.

## V. DOWNLINK COMMUNICATION DESIGN

The task of downlink communication is to broadcast the latest global model  $\mathbf{w}_t$  to the selected clients at the beginning of each learning round. It is clear that the quantization design described in Section IV-B can still be adopted in downlink. However, the differential transmission scheme in uplink is no longer feasible in the downlink for *partial* clients participation, which is a key feature of FL particularly when massive amount of clients exist [4], [21]. This is because participating clients differ from round to round, and a newly participating client does not have the “base” model of the previous round to reconstruct the new global model based on weight differential. Thus, we only focus on transmitting the global model  $\mathbf{w}_t$  at round  $t$ , and develop an enhanced method call *Layered Quantization* (LQ) for downlink communication.

### A. Layered Quantization

Layered quantization is an enhancement that builds on the quantization design in Section IV-B. We have emphasized the importance of selecting an appropriate quantization gain  $G$  to match the dynamic range of the weight or the weight differential, depending on the specific transmission method. To empirically see this, Fig. 3 plots the statistics of different layers of a typical CNN model trained for CIFAR-10 dataset<sup>4</sup>. We can see that the dynamic ranges of different layers are indeed very different. This phenomenon of varying weight distributions across layers of the DNN model has been reported in the literature [45], and an automatic clip ranging tuning method has been proposed for the secure aggregation FL protocol [48]. Intuitively, if we apply different quantization gains to different layers, the overall performance can be improved over applying one global quantization gain.

To elaborate this approach, we first denote a quantization operation on weight  $w$  with gain  $G$  as  $Q(w; G)$ . Then, the quantization gain control on a particular layer can be written as  $G = G_b G_e$ , where  $G_b$  represents the *base quantization gain* that remains the same across different DNN layers, and  $G_e$  represents the *layer-specific quantization gain*. More specifically,  $G_b$  is determined by the overall available quantization bit-width  $B$ , and  $G_e$  is then applied to adjust the position of the remained digits for the specific layer. Then, for each training round, the server can implement LQ on the global model  $\mathbf{w}_t$  to be broadcasted according to the following steps:

- 1) **Determine the base gain.** Set  $G_b = 2^{B-1}$  for all layers.

<sup>4</sup>More details can be found in Section VI.

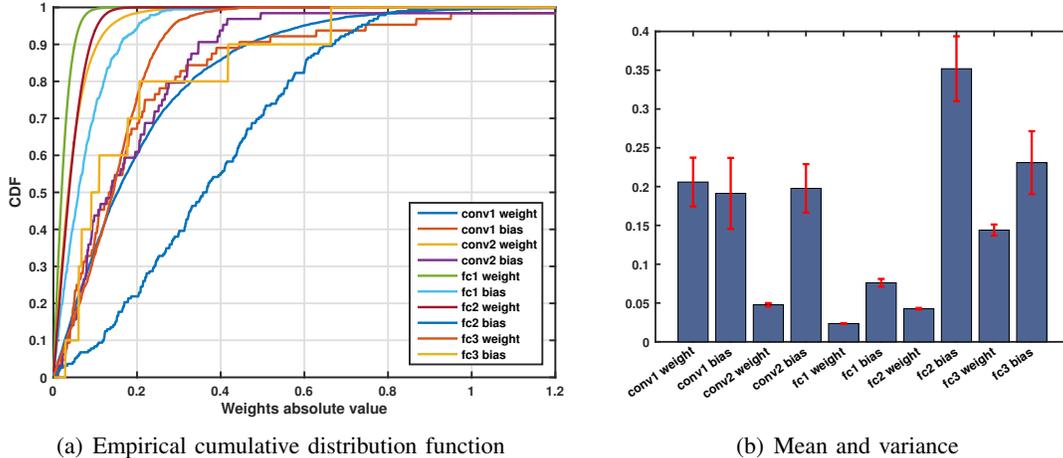


Fig. 3. Comparison of the dynamic range of weights in different layers of a typical CNN model for CIFAR-10 dataset.

2) **Determine the layer-specific gain.** For each layer of  $\mathbf{w}_t$ , calculate the empirical cumulative distribution function (CDF) of this layer, and then take the 90-percentile value  $\alpha$ . Set  $G_e = 2^\rho$ , where  $\rho = \lceil \log_2(1/\alpha) \rceil$ .

3) **Quantization.** Quantize the weights in this layer with  $Q(w; G)$  where  $G = G_b G_e$ .

The LQ design described above is “dynamic” in the sense that the layer-specific gain  $G_e$  is updated in every round of FL. As a result, the server needs to broadcast the current  $G_e$  for every layer to all participating clients, in conjunction with the latest (quantized) global model, so that the clients can properly scale down the receiver output. We note that this additional communication of broadcasting  $G_e$  for all layers is insignificant comparing to broadcasting the global model, and the overall communication overhead is not significantly increased. Furthermore, we can also adopt a “static” LQ design where the layer-specific gains  $G_e$  are determined in advance on a pre-trained model. Then,  $G_e$  can be fixed throughout the FL process (although still different across layers). This approach has the advantage of reduced computation (no need to compute the latest CDF and update  $G_e$  in each round) and reduced communication (no need to communicating the latest  $G_e$  in each round), at the expense of not tracking the dynamic range of weights in real time.

### B. Convergence Analysis for FEDAVG with Downlink Quantization

We now analyze the convergence behavior of FEDAVG with quantized downlink communication. In round  $t$ , the server first aggregates the uploaded weight update as  $\mathbf{w}_t = \frac{1}{K} \sum_{k \in \mathcal{S}_t} \mathbf{w}_t^k$  and then broadcasts a quantized version  $Q(\mathbf{w}_t)$  for round  $t + 1$ , as illustrated in Fig. 1. Suppose that we use the quantization scheme with tuned quantization and stochastic rounding as described in Section IV-B. The convergence

behavior for quantized downlink communication is characterized in Theorem 3. The proof, which is quite different from the uplink case, can be found in Appendix C.

**Theorem 3.** *Reuse the definitions of  $\kappa, \gamma, \eta_t$  in Theorem 1 and let*

$$B_t = \log_2 \left( 1 + \frac{\sqrt{1 - \eta_t \mu}}{\eta_t} \right) \quad (7)$$

*be the quantization level for the  $t$ -th iteration. With stochastic rounding on global weight  $\mathbf{w}_t$  and under Assumptions 1 and 2, the convergence bound in Eqn. (4) holds for non-i.i.d. local datasets, partial clients participation, and downlink quantization, with  $D$  being replaced by*

$$D = \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 + \frac{N-K}{N-1} \frac{4}{K} E^2 H^2 + dM^2.$$

Most of the dependencies on the hyperparameters still apply to the results in Theorem 3. As a final comment, we note that the quantization precision in Eqn. (7) suggests that  $B_t = \mathcal{O}(\log(t))$ , which matches the uplink analysis when the weight is directly transmitted. Since the downlink communication cannot adopt weight differential, it remains to be seen whether the  $\mathcal{O}(\log(t))$  requirement for quantization can be improved for the FL downlink.

## VI. EXPERIMENTS

We validate the uplink and downlink communication design and compare the performance against the floating-point baseline, which represents a natural performance upper bound. Following the setup in [4], [49], we have carried out FL experiments on four datasets: MNIST [50], CIFAR-10 [51], Shakespeare [52] and F-EMNIST [53]. Details of the setup are given in Section VI-A. Then, in Section VI-B, we focus on the performance of uplink communication and study the impact of parameters such as quantization gain and rounding. For downlink, we show in Section VI-C that a well-designed quantization scheme is critical to achieving good performance for downlink communication, and further demonstrate the performance improvement from layered quantization. Lastly, we combine both uplink and downlink designs and report the results in Section VI-D, which demonstrates that the proposed methods are capable of substantially improving the communication efficiency and, as a result, boosting the learning performance.

### A. Experiment Setup

1) *MNIST*: The training sets are evenly partitioned over  $N = 2000$  clients each containing 30 examples and we set  $K = 20$  per round (except in Fig.5, where we analysis the impact of  $K$ ) . For the **i.i.d.** case, the data is shuffled and randomly assigned to each client while for the **non-i.i.d.** case, the data is sorted

by labels, divided into 4000 shards, and each client is then assigned 2 shards randomly with 1 or 2 labels. The CNN model has two  $5 \times 5$  convolution layers, a fully connected layer with 512 units and ReLU activation, and a final output layer with softmax. The first convolution layer has 32 channels while the second one has 64 channels, and both are followed by  $2 \times 2$  max pooling. The following parameters are used for training: local batch size  $BS = 5$ , the number of local epochs  $E = 1$  for i.i.d. and  $E = 5$  for non-i.i.d., and learning rate  $\eta = 0.065$ .

2) *CIFAR-10*: The data partition is similar to the MNIST experiment for both i.i.d. and non-i.i.d. cases. We set  $N = 100$  and  $K = 10$  (except in Fig.5) for i.i.d while  $N = K = 10$  for non-i.i.d. We train a CNN model with two  $5 \times 5$  convolution layers (both with 64 channels), two fully connected layers (384 and 192 units respectively) with ReLU activation and a final output layer with softmax. The two convolution layers are both followed by  $2 \times 2$  max pooling and a local response norm layer. The training parameters are: (a) **i.i.d.**:  $BS = 50$ ,  $E = 5$ , learning rate initially sets to  $\eta = 0.15$  and decays every 10 rounds with rate 0.99; (b) **non-i.i.d.**:  $BS = 100$ ,  $E = 2$ ,  $\eta = 0.1$  and decay every round with rate 0.992.

3) *Shakespeare*: This dataset is built from *The Complete Works of William Shakespeare* and each speaking role is viewed as a device. Hence, the dataset is naturally unbalanced and **non-i.i.d.** since the number of lines and speaking habits of each role vary significantly. There are totally 1,129 roles in the dataset [53]. We randomly pick 300 of them and build a dataset with 794,659 training examples and 198,807 test examples. We also construct an i.i.d. dataset by shuffling the data and redistribute evenly to 300 roles and set  $K = 10$ . The task is the next-character prediction, and we use a classifier with an 8D embedding layer, two LSTM layers (each with 256 hidden units) and a softmax output layer with 86 nodes. The training parameters are:  $BS = 20$ ,  $E = 1$ , learning rate initially sets to  $\eta = 0.8$  and decays every 10 rounds with rate 0.99.

4) *F-EMNIST*: We use the federated version of the EMNIST dataset (F-EMNIST) [53] in this experiment. There are 3,400 clients with total 704,017 training examples and 79,952 test examples. It should be noted that F-EMNIST partitions the images of digits or English characters by their authors, thus the dataset is *naturally non-i.i.d.* since the writing style varies from person to person. We use the model recommended by [49], which is a CNN model with two convolutional layers, max pooling, and dropout, followed by a 128-unit linear layer. We set  $K = 10$ ,  $BS = 10$ ,  $E = 1$  and  $\eta = 0.03$  for training.

## B. Results for Uplink Communication

**Native quantization versus tuned quantization.** In Fig. 4, we compare the performance of Native Quantization (NQ) and Tuned Quantization (TQ), two different structures described in Section IV-B, on

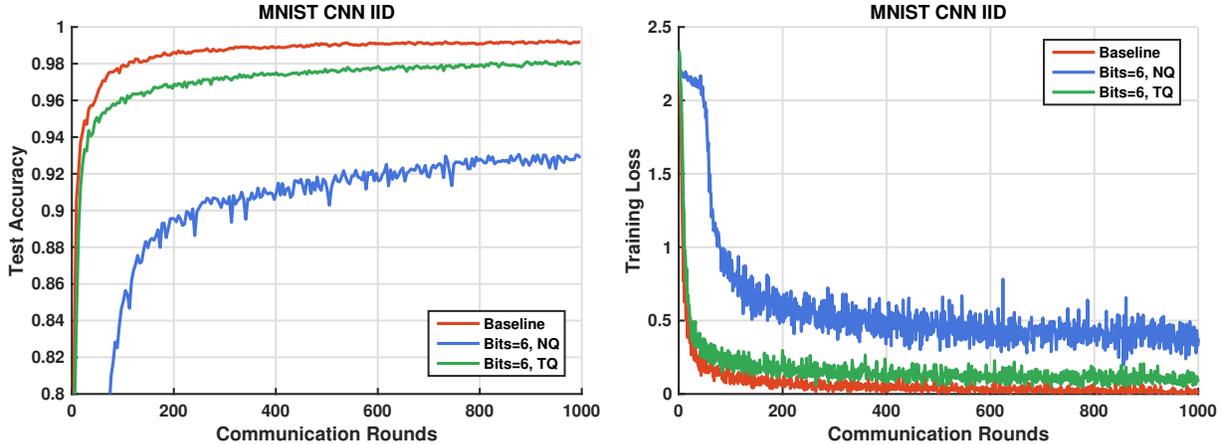


Fig. 4. Comparing the performance of Native Quantization (NQ) and Tuned Quantization (TQ).

the MNIST dataset<sup>5</sup>. The quantization gain for NQ is set to  $G = 2^{6-1} = 32$  (the maximum value for 6-bit), which has a 6.38% degradation in the test accuracy compared to the baseline. TQ on the other hand allows a larger and more suitable  $G$  (256 in this case) and achieves significantly better performance. This demonstrates the advantage of TQ.

**Nearest rounding versus stochastic rounding.** Our next experiment compares stochastic rounding (SR) and nearest rounding (NR). Although NR is widely used in communication systems, we see from Fig. 5 that SR is significantly better in both the final model accuracy and the convergence speed even when fewer bits are used. In addition, we observe an interesting phenomenon that the impact of  $K$  is different for NR and SR. For NR, having more clients participate in the model training may actually *degrade* the performance<sup>6</sup>, while this observation does not hold for SR, which is consistent with our theoretical results in Theorem 1, where a larger  $K$  reduces the value of  $D$  and leads to a reduced upper bound of convergence error. This observation is even more prominent in the results of CIFAR-10 in Fig. 5. Nevertheless, the message from the experiment is clear – one should adopt SR over NR when possible. We note that this is also consistent with the DNN compression literature [44], [45].

**Benefits of increasing quantization level.** The convergence analysis in Section IV-C indicates that to achieve an  $\mathcal{O}(\frac{1}{T})$  convergence rate with quantization, transmitting the weights without differential requires increasing the quantization level at a logarithmic rate. Our experimental results verify this conclusion. In Fig.6, the logarithmic approach increases the quantization bit-width according to

<sup>5</sup>Both model accuracy on the test set and the training loss are plotted for the remainder of this paper for all experiments.

<sup>6</sup>We hypothesize that this is because NR, which is not an unbiased quantizer, might lead to error accumulation with more clients participating in the aggregation, and this detrimental effect may outweigh the benefit of more clients. We plan to investigate this aspect in a future work.

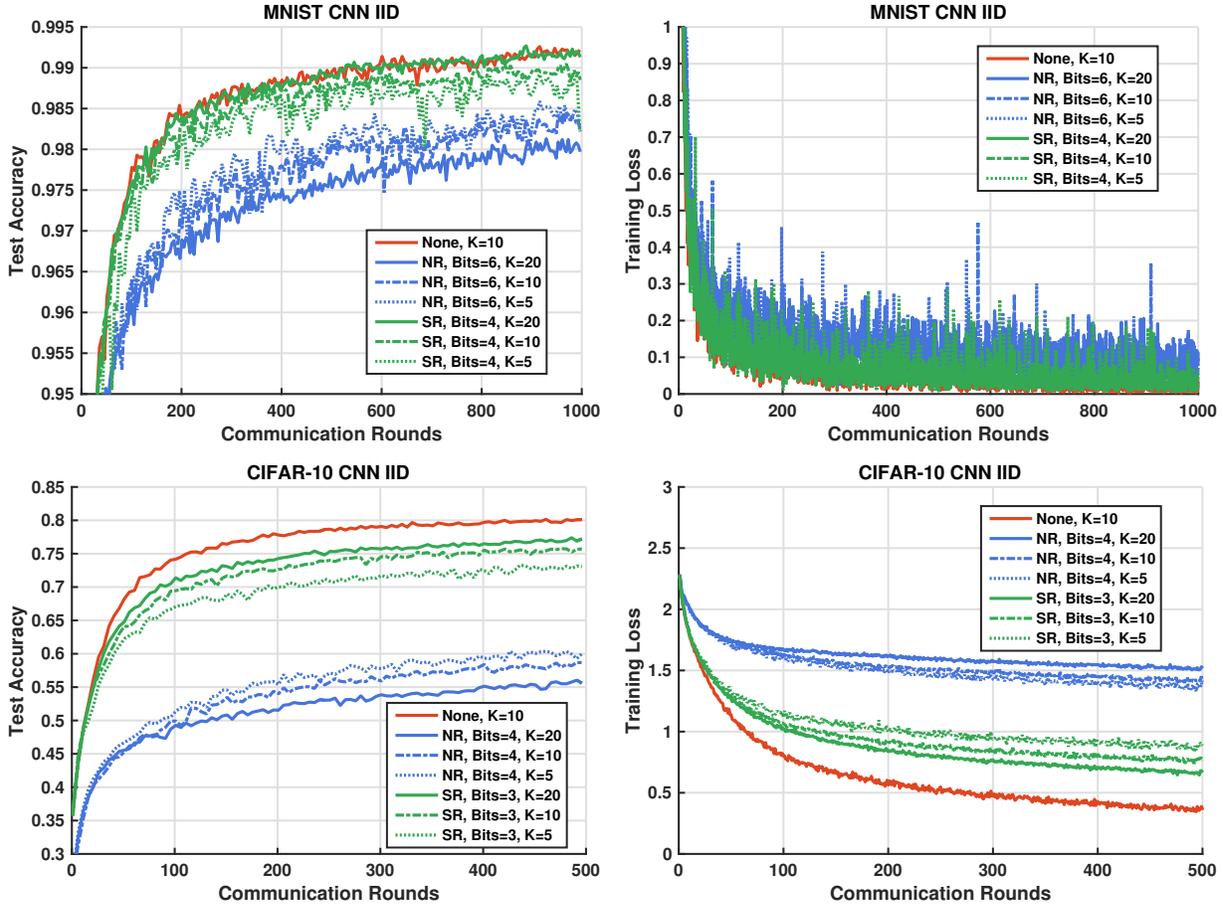


Fig. 5. Comparing the performance of Nearest Rounding (NR) and Stochastic Rounding (SR) on MNIST (top two subplots) and CIFAR-10 (bottom two subplots).

$B = \lfloor \log_2 [f + (r - 1)/p] \rfloor$ , where  $r = 1, 2, \dots$  is the index of training round. By contrast, the fixed approach keep a constant bit-width throughout. In Fig.6, the average bit-width for each round of the logarithmic approach on CIFAR-10 dataset is 2, but we can see that it outperforms the result with fixed 2-bit in the final convergence accuracy. The average bit-width of the logarithmic approach on Shakespeare dataset is 3, which also have better performance compared to the fix-3bit quantization.

**Advantages of differential transmission.** One of the key benefits in using DT is that the dynamic range of weight differential  $\mathbf{d}_{t+1}^k$  is much smaller than the weight  $\mathbf{w}_{t+1}^k$  itself, and thus quantization will be more precise with the same bit-width  $B$ . We now empirically validate this point by plotting the empirical cumulative distribution function (CDF) of both representations in Fig. 7. We can see that DT has a dynamic range that is an order-of-magnitude smaller than the weight itself, which suggests that the intuition is correct. We also see that distribution of the weight differential gradually concentrates and the support also decreases as the training progresses towards the end. At round 10, 35% of the weights is

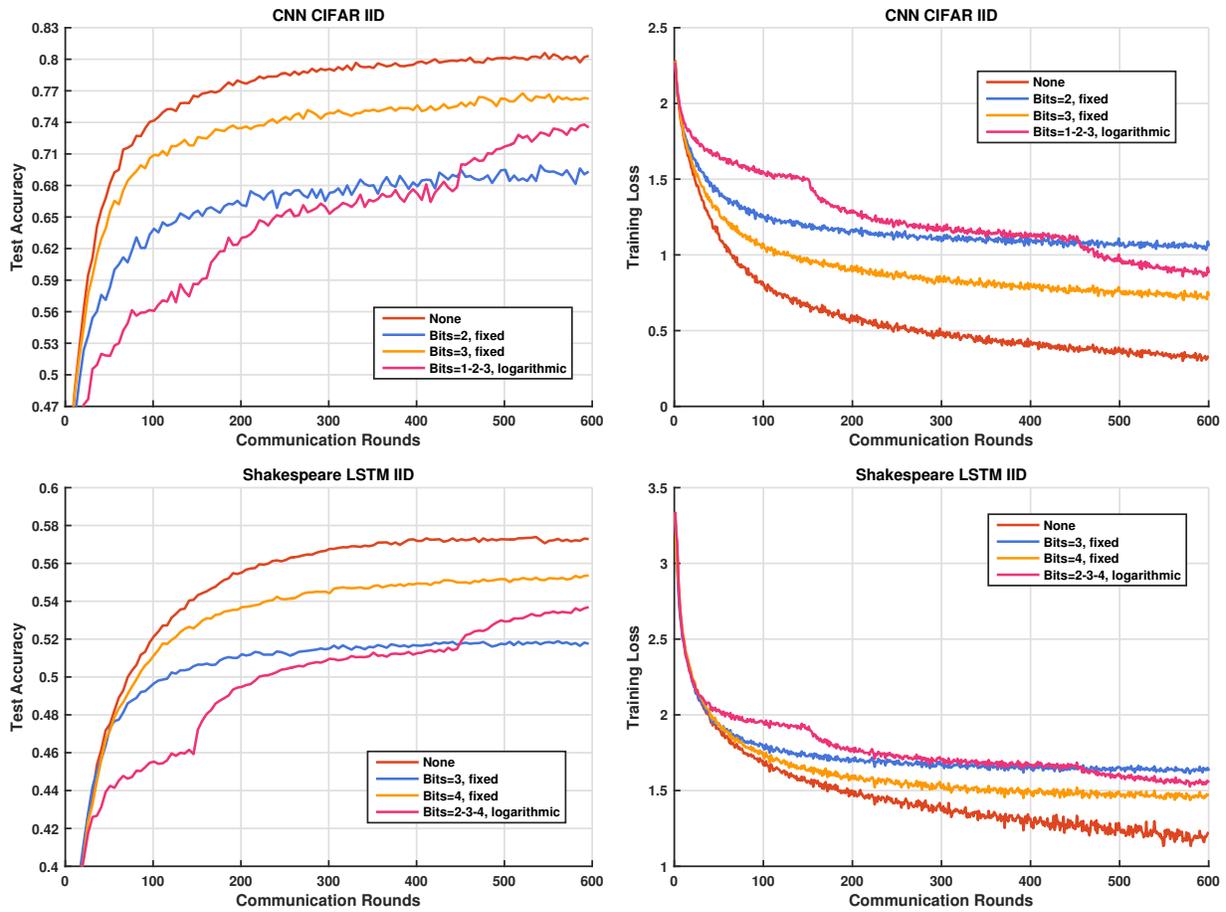


Fig. 6. Comparing the performance of fixed and increasing quantization level on CIFAR-10 (top two subplots) with  $f = 2, p = 75$  and Shakespeare (bottom two subplots) with  $f = 4, p = 37.5$ .

less than  $9e-5$  while at round 500, this proportion achieves 90%. This is another useful observation, as it indicates that we may be able to decrease the quantization bit-width at the late stage of training.

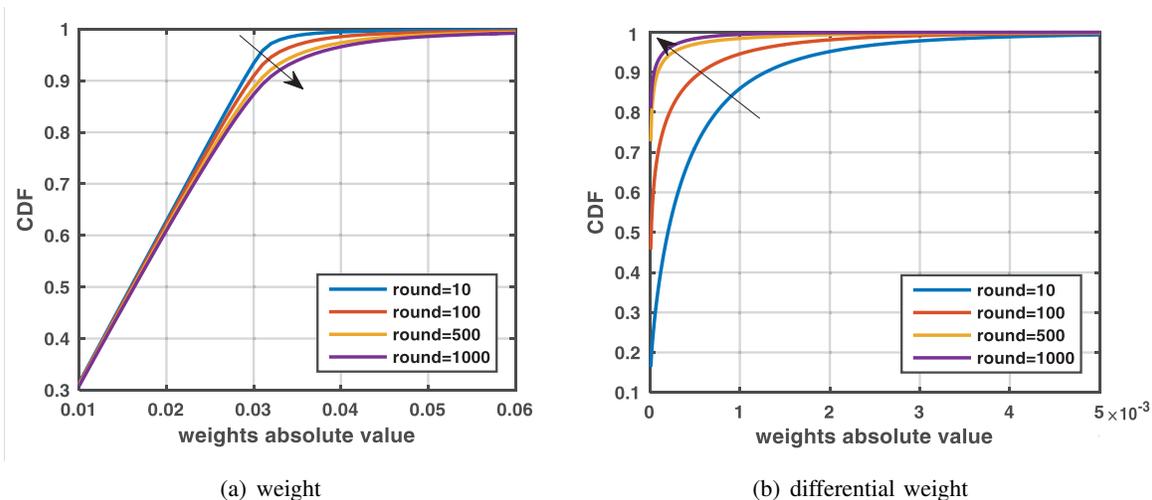


Fig. 7. Comparing the distribution and value range of the weight itself and differential weight (MNIST, i.i.d.).  $K$  is set to 20 and at specified rounds we make statistics of the weights of the 20 selected clients.

**Putting all techniques together.** Finally, we report an experiment where TQ, SR and DT are combined, which represents the best quantization design for uplink communication in our paper. We are interested in evaluating how well this design performs, especially comparing against the floating-point baseline (no quantization). Fig. 8 shows that, for both i.i.d. and non-i.i.d. cases, we are able to quantize the floating-point weight differential to 1-bit representations with almost negligible performance loss:

- **i.i.d.** 99.08% accuracy (99.83% of the baseline accuracy) for 1-bit (3.13% of the baseline bandwidth); 99.18% accuracy (99.93% of the baseline accuracy) for 2-bit (6.25% of the baseline bandwidth);
- **non-i.i.d.** 98.59% accuracy (99.41% of the baseline accuracy) for 1-bit (3.13% of the baseline bandwidth); and 98.99% accuracy (99.81% of the baseline accuracy) for 2-bit (6.25% of the baseline bandwidth).

These results suggest that the proposed design achieves the best communication efficiency in this FL task, to the best of the authors' knowledge.

**Other datasets.** To further evaluate the performance of our uplink design, we also run experiments on CIFAR-10, Shakespeare and F-EMNIST, in which the training tasks are harder than classification on MNIST. We report the results with the best quantization method (combining TQ, SR and DT) of these three datasets in Fig. 9, Fig. 10 and Fig. 11, respectively. The results suggest that, with a well-designed uplink quantization, using 3 bits or fewer allows federated learning to achieve sufficiently good performance, for both i.i.d. and non-i.i.d. dataset.

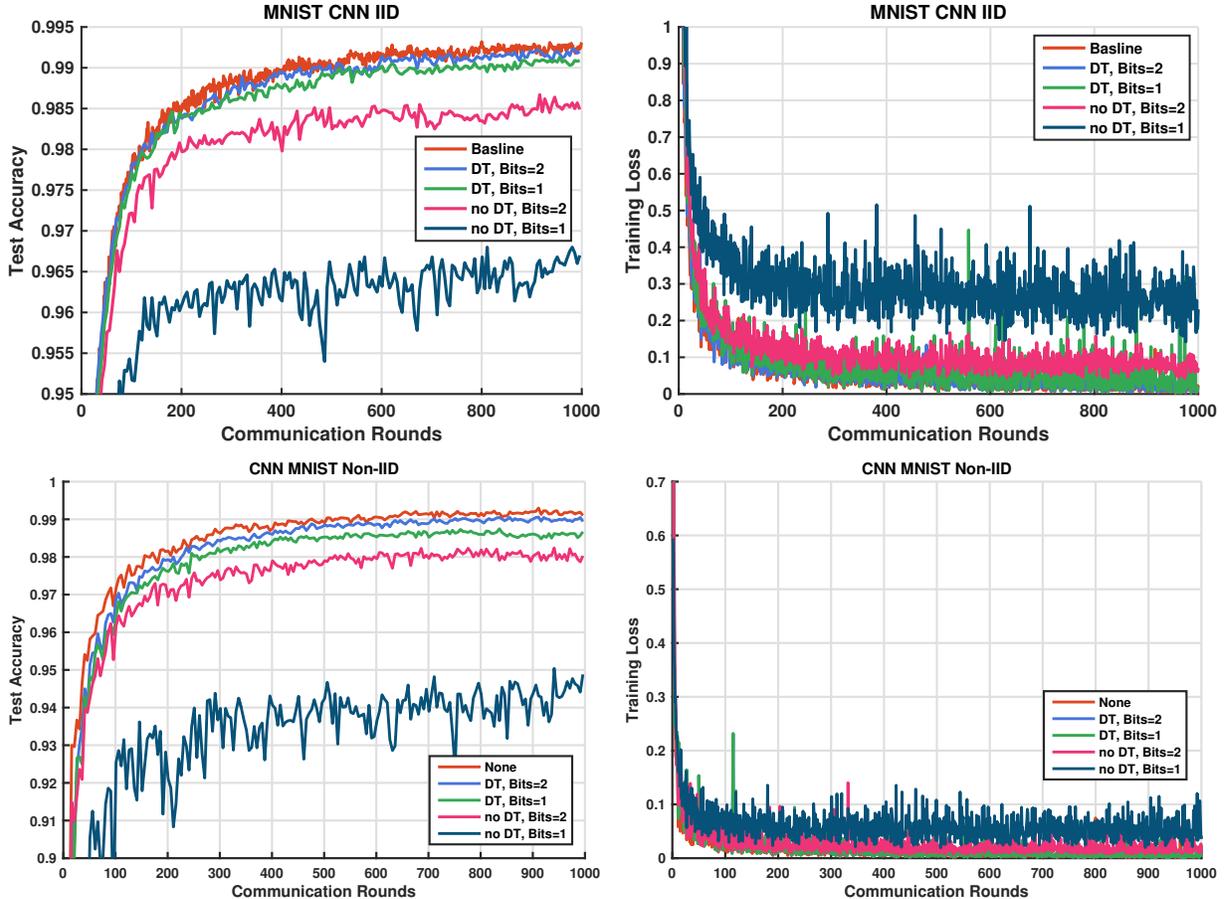


Fig. 8. Comparing the performance and transmission with and without DT on i.i.d. (top two subplots) and non-i.i.d. (bottom two subplots) MNIST dataset. Both are quantized with TQ and SR. For the 1-bit DT of the non-i.i.d. case, the learning rate is reduced to 0.03.

### C. Results for Downlink Communication

**Quantization has a bigger impact on downlink communication.** We evaluate the impact of low-precision quantization on downlink communication in this subsection. Our experimental results in Fig. 12 suggest that a poorly designed quantization scheme (e.g., NQ with NR) for downlink can significantly degrade the performance of the overall FL – for the same quantization level  $B$  and the same quantization method, quantization in downlink has worse performance than quantization in uplink. This can be intuitively understood since the downloaded model is used by many clients, and hence the inaccuracy can manifest, resulting in a broader impact than upload inaccuracy [2].

**Performance of quantization in downlink.** We have mentioned in Section V that the quantization scheme designed for uplink communication can be adapted in downlink except DT. We now report the results of using quantization with TQ and SR on the CIFAR-10 dataset in Fig. 13. The results show that

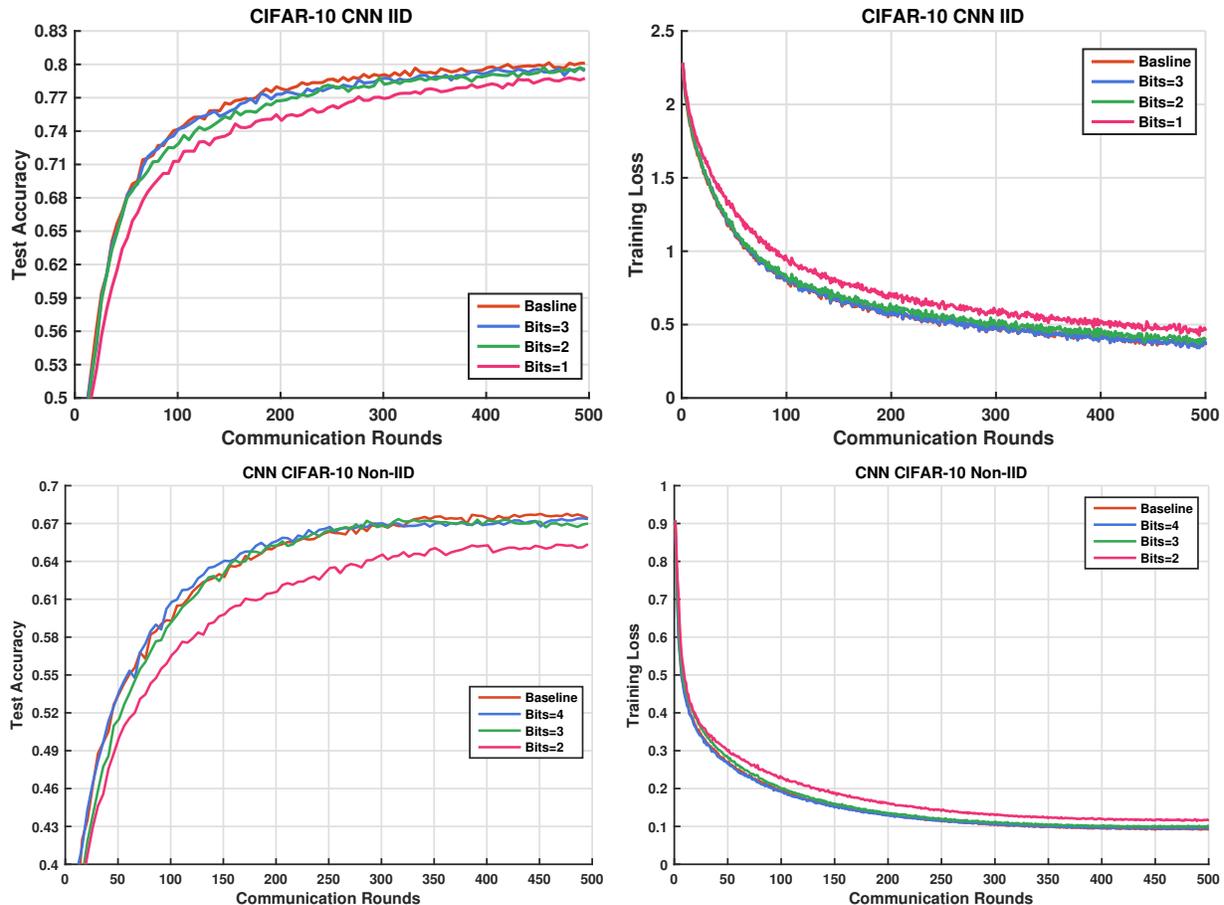


Fig. 9. The performance of uplink quantization on i.i.d. (top two subplots) and non-i.i.d. (bottom two subplots) CIFAR-10 dataset.

with a well-design scheme, the downlink communication can also be made efficient and effective. Unlike the results in Fig. 12, the accuracy of a 6-bit quantization can achieve 78.33% accuracy (98% of the baseline accuracy). However, without the support of DT, we see that there is a noticeable performance reduction when the bit-width falls below 3.

We now evaluate layered quantization (LQ) and see if it can improve the performance. We use the method described in Section V-A to carefully set an appropriate quantization gain for each layer. The results reported in Fig. 14 suggest that LQ is effective for both i.i.d. and non-i.i.d. cases. For CIFAR-10 (F-EMNIST), it improves the performance of the 3-bit communication from 74.48% to 77.15% (75.23% to 78.29%) and the 4-bit communication from 76.04% to 78.42% (78.48% to 80.46%), respectively.

**Results on other datasets.** We now have identified the combination of TQ, SR and layered quantization as the best design options for downlink communication, and we now validate this combination on other datasets. The results from Fig. 15 further confirm that, at least for the three datasets we have evaluated,

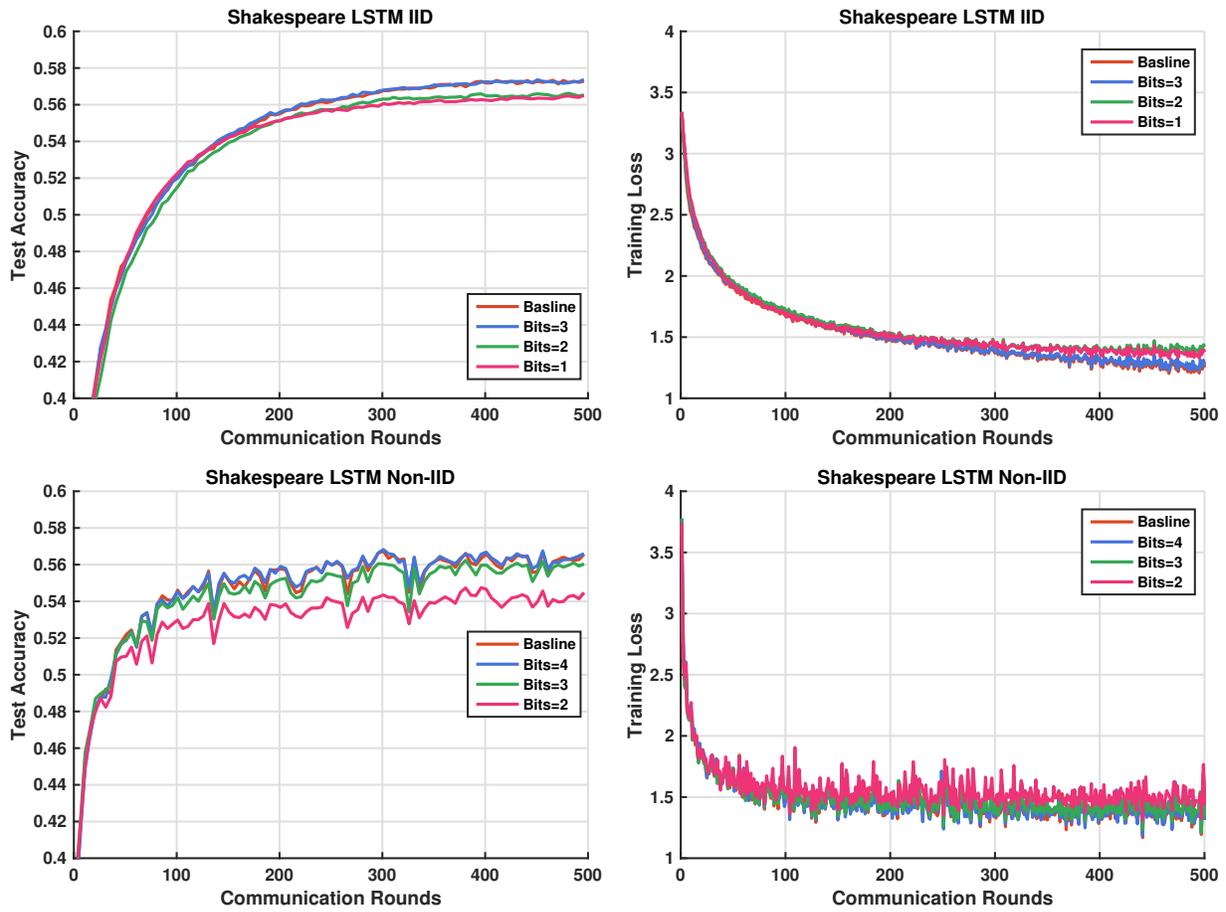


Fig. 10. The performance of uplink quantization on i.i.d. (top two subplots) and non-i.i.d. (bottom two subplots) Shakespeare dataset.

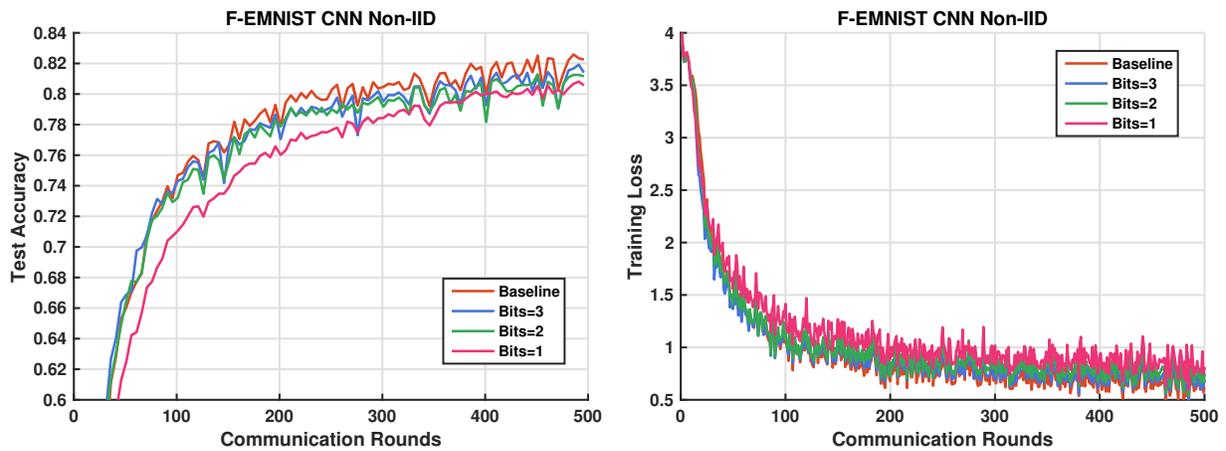


Fig. 11. The performance of uplink quantization on the naturally non-i.i.d. F-EMNIST dataset.

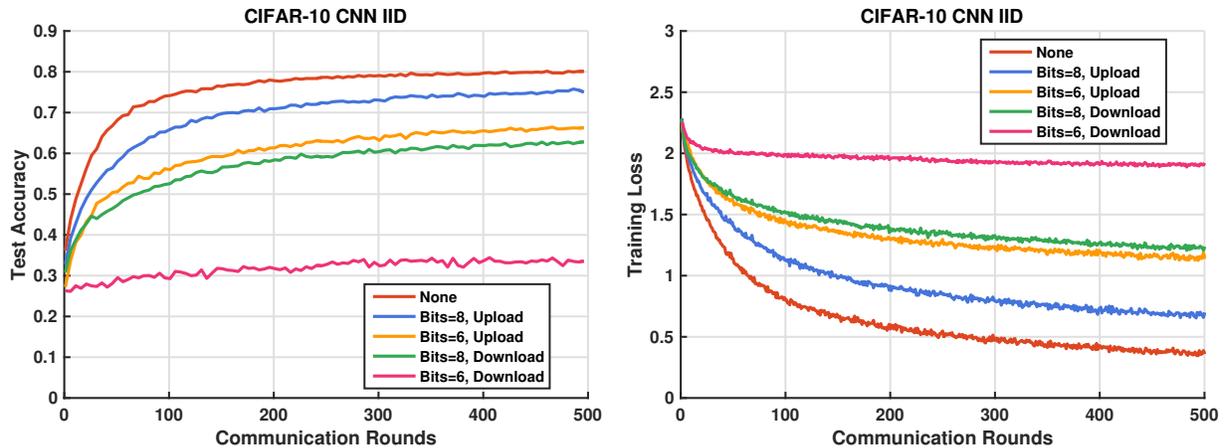


Fig. 12. Comparison of the impact of quantization (NQ and NR) on uplink and downlink communications.

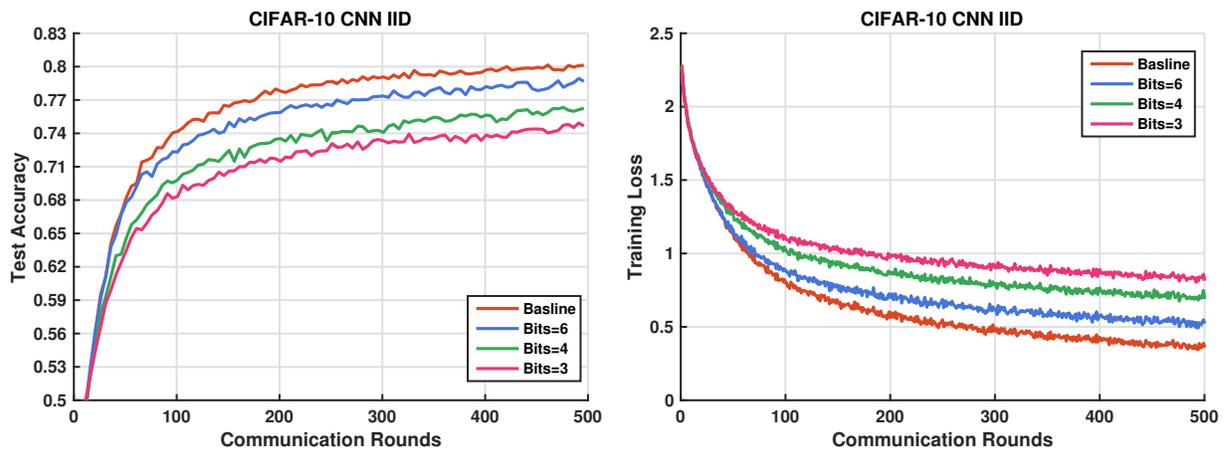


Fig. 13. Performance of quantization with TQ and SR in downlink communication.

the proposed design can reduce the quantization bit-width to 4 (12.5% of the baseline bandwidth) while achieving an accuracy degradation within 2% of the baseline accuracy.

#### D. Results of Quantization on Both Uplink and Downlink

Lastly, we carry out experiment with simultaneous quantization on both uplink and downlink communications. The experimental results on different datasets are reported in Table I. We run 1000 rounds for MNIST and average the final 100 rounds as the final (convergence) accuracy (the fourth column). As for CIFAR-10, Shakespeare and F-EMNIST, we run 500 rounds and average the final 50 rounds. The last column shows the percentage of the baseline (using 32-bit float) can be achieved by the learning with quantized communications in both uplink and downlink. For all experiments, layered quantization with TQ and SR is used for downlink while DT with TQ and SR is used for uplink.

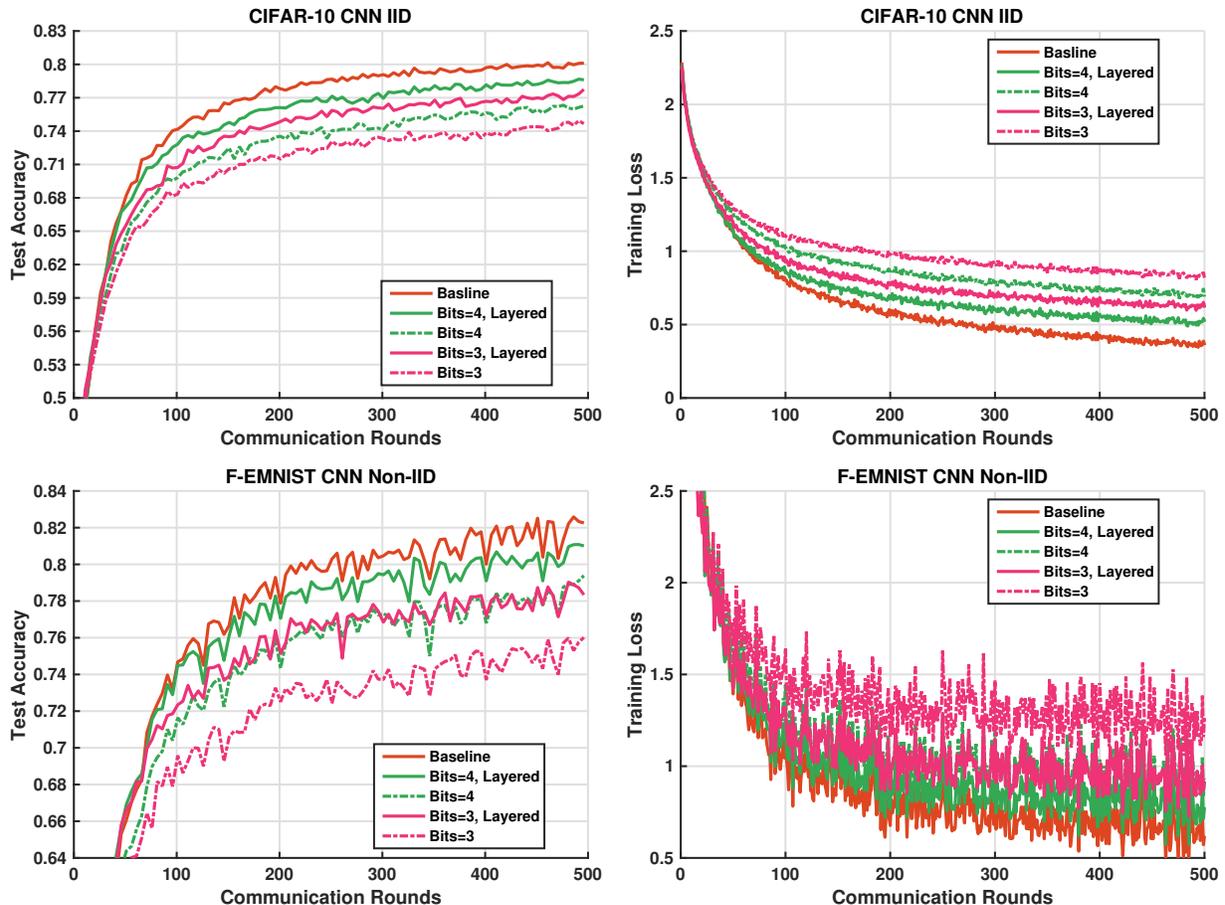


Fig. 14. Comparing the performance of quantization with and without layered quantization on CIFAR-10 (top two subplots) and F-EMNIST (bottom two subplots) datasets.

We evaluate how much communication payload can be reduced while maintaining a small accuracy loss (defined as less than 2%). The results in Table I show that well designed quantization schemes are important to improve the communication efficiency. Take MNIST (i.i.d.) as an example, 2-bit for both downlink or uplink are sufficiently good, which can reduce (from the baseline) 93.75% in communications for both uplink and downlink. Even for the more complex cases such as CIFAR-10 (non-i.i.d.), 6-bit for downlink and 4-bit for uplink have very good performance, reducing 81.25% and 87.5% of the communication bandwidth for each client on downlink and uplink respectively. Overall, we conclude that the proposed designs are effective in addressing the communication bottleneck of federated learning.

### E. Impact of hyperparameters

There are several hyperparameters that impact the training of FL, some of which have been discussed in [4]. Fig. 16 shows the relationship between the quantization and these hyperparameters. Local batch size

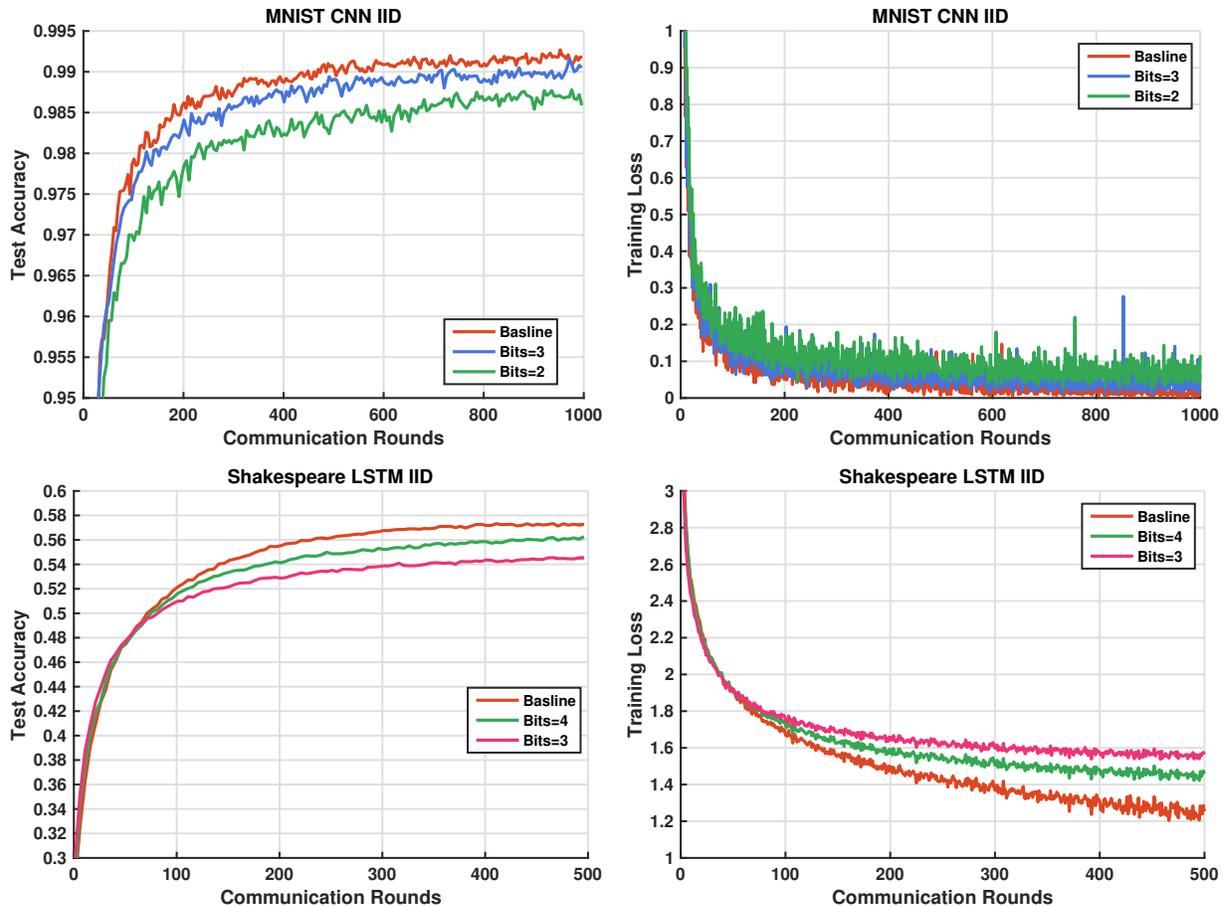


Fig. 15. The performance of the quantization scheme designed for downlink on MNIST (top two subplots) and Shakespeare (bottom two subplots) datasets.

TABLE I  
PERFORMANCE OF SIMULTANEOUS QUANTIZATION ON BOTH UPLINK AND DOWNLINK.

Dataset	Downlink	Uplink	Accuracy (Baseline)	Percentage*
i.i.d.				
MNIST	2-bit	2-bit	98.46% (99.11%)	99.34%
CIFAR-10	5-bit	2-bit	78.43% (79.93%)	98.12%
Shakespeare	5-bit	2-bit	56.17% (57.25%)	98.11%
Non-i.i.d.				
MNIST	2-bit	2-bit	97.41% (99.10%)	98.29%
CIFAR-10	6-bit	4-bit	61.67% (62.61%)	98.50%
Shakespeare	5-bit	3-bit	55.16% (56.16%)	98.22%
F-EMNIST	5-bit	3-bit	80.24% (81.82%)	98.07%

\*: The last column represents the percentage of FL accuracy against the baseline accuracy.

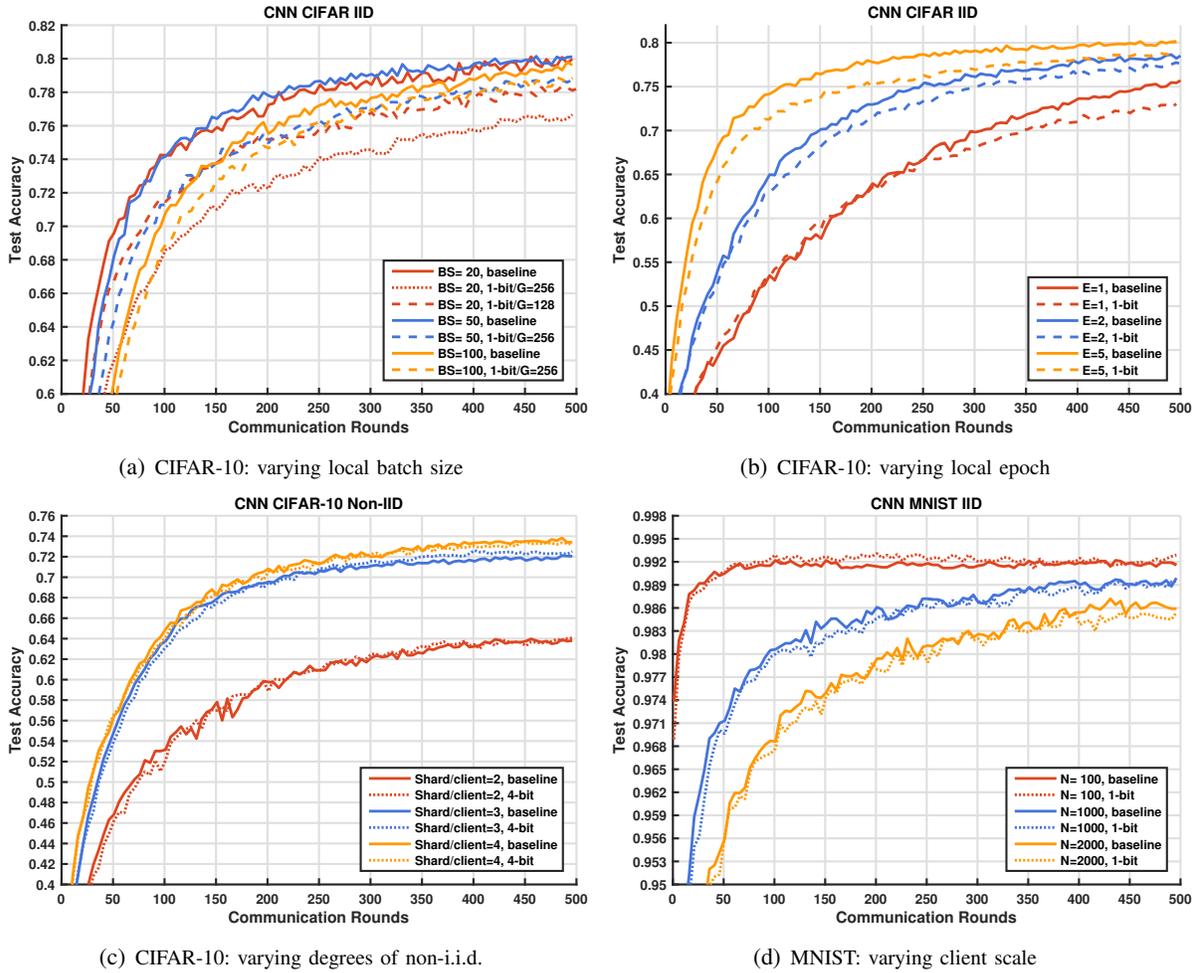


Fig. 16. Impact of different hyperparameters with quantized uplink transmission in FL. Differential transmission is used for all the experiments. The baseline curves mean the results without any quantization.

( $BS$ ) is suggested to be small in most cases, and we can see in Fig. 16(a) that the baseline performance of  $BS = 20$  and  $50$  are better. However, we notice that the accuracy loss between the baseline and 1-bit quantization is the smallest for  $BS = 100$ . For the setting of  $BS = 20$ , smaller quantization gain  $G$  should be used; otherwise the performance is severely degraded. The possible reason is that a smaller  $BS$  brings more local iterations on the clients at each round and then increases the dynamic range of weight differentials, which indicates that there is a tradeoff between increasing computation per client and better quantization performance. As for the local epochs, although a larger  $E$  might also increase the dynamic range, this becomes less important, compared to the benefit in convergence speed, as shown in Fig. 16(b). Therefore, even with quantization, a relative large  $E$  can be adopted, especially for the i.i.d. cases. Fig. 16(c) and Fig. 16(d) imply that the proposed quantization scheme is not very sensitive to the degree of non-i.i.d. or client scale, since the theoretical results have shown that the impact of quantization

is decoupled with  $\Gamma$  or  $N$ .

## VII. CONCLUSIONS

We have studied the design and analysis of physical layer quantization and transmission methods for wireless federated learning. If nothing else, this paper showed that the communication design must tailor to the characteristics of FL. In particular, the choice of *what* to transmit and *how* to transmit them has a profound impact on the performance of federated learning in a wireless system, and we established this conclusion both theoretically, via convergence analysis of various quantization and transmission options in the well-known FEDAVG, and experimentally, via comprehensive evaluation on real-world datasets. An important theoretical convergence result was established, which states that in order to achieve an  $\mathcal{O}(1/T)$  convergence rate with quantization, transmitting the weight requires increasing the quantization level at a *logarithmic* rate, while transmitting the weight differential can keep a constant quantization level. As a crown jewel of the experimental study, we were able to achieve a significant milestone: 1-bit quantization (3.1% of the floating-point baseline bandwidth) achieves 99.8% of the floating-point baseline accuracy at almost the same convergence rate on MNIST, representing the best known bandwidth-accuracy tradeoff to the best of the authors' knowledge.

In addition to enabling efficient communication design for FL, we have noticed that quantization can also be combined with communication resource allocation. For example, the theoretical result of Theorem 1 naturally leads to a resource (bit) allocation problem where one is given a total budget of uplink bandwidth and asked to allocate the bits over communication rounds to optimize the learning performance. Another interesting future research direction is the combination of quantization and client selection. For example, for a given total uplink bandwidth budget, how to balance the increased number of clients and reduced quantization precision.

## APPENDIX A

### PROOF OF THEOREM 1

#### A. Notations

In our analysis, there are three sources of randomness: stochastic gradients, random sampling of clients, and stochastic rounding. To distinguish them, we respectively use the notation  $\mathbb{E}_{SG}[\cdot]$ ,  $\mathbb{E}_{S_t}[\cdot]$  and  $\mathbb{E}_{SR}[\cdot]$  and use  $\mathbb{E}[\cdot]$  for expectation over all three of them. With a slight abuse of notation, we change the timeline to be with respect to the SGD iteration time instead of the communication round. Let  $\mathbf{w}_t^k$  be the model weights on the  $k$ th client at the  $t$ th iteration and  $\mathbf{w}_t$  be the global model at the  $t$ th iteration. In FEDAVG,

clients perform  $E$  local iterations before global aggregation. Hence,  $\mathbf{w}_t$  is only accessible for specific  $t \in \mathcal{I}_E$ , where  $\mathcal{I}_E = \{nE \mid n = 1, 2, \dots\}$ , i.e. the time for communication.

For client  $k$ , it trains the model locally with

$$\mathbf{v}_{t+1}^k = \mathbf{w}_t^k - \eta_t \nabla F_k(\mathbf{w}_t^k, \xi_t^k). \quad (8)$$

If  $t+1 \notin \mathcal{I}_E$ , the next-step result is  $\mathbf{w}_{t+1}^k = \mathbf{v}_{t+1}^k$  since no global aggregation takes place. If  $t+1 \in \mathcal{I}_E$ , all client  $k \in \mathcal{S}_{t+1}$  upload their quantized weights  $Q(\mathbf{v}_{t+1}^k)$ . The global model is updated with  $\mathbf{w}_{t+1} = \frac{1}{K} \sum_{k \in \mathcal{S}_{t+1}} Q(\mathbf{v}_{t+1}^k)$ . Since we do not model downlink quantization, selected clients update their local weights as  $\mathbf{w}_{t+1}^k = \mathbf{w}_{t+1}$  and start the next local training period. We define the following three variables to summarize the aforementioned steps:

$$\begin{aligned} \mathbf{v}_{t+1}^k &\triangleq \mathbf{w}_t^k - \eta_t \nabla F_k(\mathbf{w}_t^k, \xi_t^k); \\ \mathbf{u}_{t+1}^k &\triangleq \begin{cases} \mathbf{v}_{t+1}^k & \text{if } t+1 \notin \mathcal{I}_E, \\ \frac{1}{K} \sum_{i \in \mathcal{S}_{t+1}} \mathbf{v}_{t+1}^i & \text{if } t+1 \in \mathcal{I}_E; \end{cases} \\ \mathbf{w}_{t+1}^k &\triangleq \begin{cases} \mathbf{v}_{t+1}^k & \text{if } t+1 \notin \mathcal{I}_E, \\ \frac{1}{K} \sum_{i \in \mathcal{S}_{t+1}} Q(\mathbf{v}_{t+1}^i) & \text{if } t+1 \in \mathcal{I}_E. \end{cases} \end{aligned}$$

We define three **virtual sequences**  $\bar{\mathbf{v}}_t = \frac{1}{N} \sum_{k=1}^N \mathbf{v}_t^k$ ,  $\bar{\mathbf{w}}_t = \frac{1}{N} \sum_{k=1}^N \mathbf{w}_t^k$  and  $\bar{\mathbf{u}}_t = \frac{1}{N} \sum_{k=1}^N \mathbf{u}_t^k$  to facilitate the analysis. For convenience, we define  $\bar{\mathbf{g}}_t = \frac{1}{N} \sum_{k=1}^N \nabla F_k(\mathbf{w}_t^k)$  and  $\mathbf{g}_t = \frac{1}{N} \sum_{k=1}^N \nabla F_k(\mathbf{w}_t^k, \xi_t^k)$ . Therefore,  $\bar{\mathbf{v}}_t = \bar{\mathbf{w}}_t - \eta_t \mathbf{g}_t$  and  $\mathbb{E}_{SG}[\mathbf{g}_t] = \bar{\mathbf{g}}_t$ . Notice that we take average over all  $N$  instead of  $K$  clients, which is because for  $t+1 \in \mathcal{I}_E$ , we have

$$\bar{\mathbf{w}}_{t+1} = \frac{1}{N} \sum_{k=1}^N \mathbf{w}_t^k = \frac{1}{K} \sum_{k \in \mathcal{S}_{t+1}} Q(\mathbf{v}_{t+1}^k) \quad (9)$$

and the global model is meaningful only at  $t+1 \in \mathcal{I}_E$ .

## B. Lemmas

We present some necessary lemmas that are useful in the proof of Theorem 1. Lemmas 1, 2 and 3 have been established in [11] for floating-point weights. Because (1) Lemma 1 is derived based on the smoothness and convexity of  $F_k(\cdot)$ ; (2) Lemma 2 is derived based on the bounded variance of SGD; and (3) Lemma 3 is derived based on the bounded gradient of  $F_k(\cdot)$ , these lemmas still hold under Assumption 1 for quantized FEDAVG.

**Lemma 1** (Result of one step SGD). *Let Assumption 1-1) and 2) hold. If  $\eta_t \leq \frac{1}{4L}$ , we have*

$$\mathbb{E}_{SG} \|\bar{\mathbf{v}}_{t+1} - \mathbf{w}^*\|^2 \leq (1-\eta_t\mu)\mathbb{E}_{SG} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + \eta_t^2 \mathbb{E}_{SG} \|\mathbf{g}_t - \bar{\mathbf{g}}_t\|^2 + 6L\eta_t^2\Gamma + 2\mathbb{E}_{SG} \left[ \frac{1}{N} \sum_{k=1}^N \|\bar{\mathbf{w}}_t - \mathbf{w}_t^k\|^2 \right].$$

**Lemma 2** (Bounding the variance). *Let Assumption 1-3) hold, It follows that*

$$\mathbb{E}_{SG} \|\mathbf{g}_t - \bar{\mathbf{g}}_t\|^2 \leq \sum_{k=1}^N \frac{\sigma_k^2}{N^2}.$$

**Lemma 3** (Bonding the divergence of  $\mathbf{w}_t^k$ ). *Let Assumption 1-4) hold,  $\eta_t$  is non-increasing and  $\eta_t \leq 2\eta_{t+E}$  for all  $t \geq 0$ . It follows that*

$$\mathbb{E}_{SG} \left[ \frac{1}{N} \sum_{k=1}^N \|\bar{\mathbf{w}}_t - \mathbf{w}_t^k\|^2 \right] \leq 4\eta_t^2(E-1)^2 H^2.$$

Lemmas 4 to 6 are specific for uplink quantization of FEDAVG, whose proofs are deferred to Appendix A-D.

**Lemma 4** (Unbiased and variance bounded sampling). *Let Assumption 1-4) hold. For  $t+1 \in \mathcal{I}_E$ , assume that  $\eta_t \leq 2\eta_{t+E}$  for all  $t \geq 0$ . We have*

$$\begin{aligned} \mathbb{E}_{\mathcal{S}_t} [\bar{\mathbf{u}}_{t+1}] &= \bar{\mathbf{v}}_{t+1}, \\ \mathbb{E} \|\bar{\mathbf{v}}_{t+1} - \bar{\mathbf{u}}_{t+1}\|^2 &\leq \frac{N-K}{N-1} \frac{4}{K} \eta_t^2 E^2 H^2. \end{aligned} \tag{10}$$

**Lemma 5** (Properties of stochastic rounding). *For a vector  $\mathbf{w} \in \mathbb{R}^d$  satisfying  $\|\mathbf{w}\|_\infty \leq M$ , let  $Q(\mathbf{w})$  be the quantization of  $\mathbf{w}$  with stochastic rounding, quantization level  $B$  and quantization gain  $G = \frac{2^{B_t-1}}{M}$ . Then we have:*

$$\begin{aligned} \mathbb{E}_{SR} [Q(\mathbf{w})] &= \mathbf{w}, \\ \mathbb{E}_{SR} \left[ \|Q(\mathbf{w}) - \mathbf{w}\|^2 \right] &\leq d \left( \frac{M}{2^B - 1} \right)^2. \end{aligned}$$

**Lemma 6** (Unbiased and variance bounded quantization). *Let Assumption 2 hold. With stochastic rounding and the quantization level set to  $B_{t+1}$ , we have*

$$\begin{aligned} \mathbb{E}_{SR} [\bar{\mathbf{w}}_{t+1}] &= \bar{\mathbf{u}}_{t+1}, \\ \mathbb{E}_{SR} \|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{u}}_{t+1}\|^2 &\leq q_{t+1}^2 \cdot \frac{dM^2}{K} \end{aligned} \tag{11}$$

for  $t+1 \in \mathcal{I}_E$ , where  $q_{t+1} = 1/(2^{B_{t+1}} - 1)$ .

### C. Proof of Theorem 1

If  $t+1 \notin \mathcal{I}_E$ ,  $\bar{\mathbf{w}}_{t+1} = \bar{\mathbf{v}}_{t+1}$ , then using Lemmas 1 to 3, we first take expectation over the randomness of stochastic gradient and get

$$\begin{aligned} \mathbb{E}_{SG} \|\bar{\mathbf{w}}_{t+1} - \mathbf{w}^*\|^2 &= \mathbb{E}_{SG} \|\bar{\mathbf{v}}_{t+1} - \mathbf{w}^*\|^2 \\ &\leq (1 - \eta_t \mu) \mathbb{E}_{SG} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + \eta_t^2 \mathbb{E}_{SG} \|\mathbf{g}_t - \bar{\mathbf{g}}_t\|^2 + 6L\eta_t^2 \Gamma + 2\mathbb{E}_{SG} \left[ \frac{1}{N} \sum_{k=1}^N \|\bar{\mathbf{w}}_t - \mathbf{w}_t^k\|^2 \right] \\ &\leq (1 - \eta_t \mu) \mathbb{E}_{SG} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + \eta_t^2 \left[ \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 \right]. \end{aligned}$$

We then take expectation over the randomness of  $\mathcal{S}_t$  and stochastic rounding to have

$$\mathbb{E} \|\bar{\mathbf{w}}_{t+1} - \mathbf{w}^*\|^2 \leq (1 - \eta_t \mu) \mathbb{E} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + \eta_t^2 \left[ \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 \right]. \quad (12)$$

If  $t+1 \in \mathcal{I}_E$ , note that

$$\begin{aligned} \|\bar{\mathbf{w}}_{t+1} - \mathbf{w}^*\|^2 &= \|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{u}}_{t+1} + \bar{\mathbf{u}}_{t+1} - \mathbf{w}^*\|^2 \\ &= \underbrace{\|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{u}}_{t+1}\|^2}_{A_1} + \underbrace{\|\bar{\mathbf{u}}_{t+1} - \mathbf{w}^*\|^2}_{A_2} + 2 \underbrace{\langle \bar{\mathbf{w}}_{t+1} - \bar{\mathbf{u}}_{t+1}, \bar{\mathbf{u}}_{t+1} - \mathbf{w}^* \rangle}_{A_3}. \end{aligned} \quad (13)$$

When the expectation is taken over the randomness of stochastic rounding, the last term  $A_3$  vanishes since we have  $\mathbb{E}_{SR} [\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{u}}_{t+1}] = \mathbf{0}$  (from Eqn. (11)).  $A_1$  can be bounded using Lemma 6. As for  $A_2$ , we have

$$\begin{aligned} \|\bar{\mathbf{u}}_{t+1} - \mathbf{w}^*\|^2 &= \|\bar{\mathbf{u}}_{t+1} - \bar{\mathbf{v}}_{t+1} + \bar{\mathbf{v}}_{t+1} - \mathbf{w}^*\|^2 \\ &= \underbrace{\|\bar{\mathbf{u}}_{t+1} - \bar{\mathbf{v}}_{t+1}\|^2}_{B_1} + \underbrace{\|\bar{\mathbf{v}}_{t+1} - \mathbf{w}^*\|^2}_{B_2} + 2 \underbrace{\langle \bar{\mathbf{u}}_{t+1} - \bar{\mathbf{v}}_{t+1}, \bar{\mathbf{v}}_{t+1} - \mathbf{w}^* \rangle}_{B_3}. \end{aligned}$$

When expectation is taken over the randomness of  $\mathcal{S}_t$ , the last term  $B_3$  vanishes because  $\mathbb{E}_{\mathcal{S}_t} [\bar{\mathbf{u}}_{t+1} - \bar{\mathbf{v}}_{t+1}] = \mathbf{0}$  (from Eqn. (10)).  $B_1$  can be bounded using Lemma 4, and  $B_2$  can be bounded using Lemmas 1 to 3.

In summary, by taking expectation over all the three randomnesses on Eqn. (13), we finally have

$$\begin{aligned} \mathbb{E} \|\bar{\mathbf{w}}_{t+1} - \mathbf{w}^*\|^2 &\leq (1 - \eta_t \mu) \mathbb{E} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + q_{t+1}^2 \frac{dM^2}{K} \\ &\quad + \eta_t^2 \left[ \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 + \frac{N-K}{N-1} \frac{4}{K} E^2 H^2 \right]. \end{aligned} \quad (14)$$

If we increase the quantization level  $B_{t+1}$  following

$$B_{t+1} = \log_2(1/\eta_t + 1),$$

then

$$q_{t+1} = 1/(2^{B_{t+1}} - 1) = \eta_t.$$

Let  $\Delta_t = \mathbb{E} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2$ . From Eqn. (12) and Eqn. (14), it is clear that no matter whether  $t+1 \in \mathcal{I}_E$  or  $t+1 \notin \mathcal{I}_E$ , we always have

$$\Delta_{t+1} \leq (1 - \eta_t \mu) \Delta_t + \eta_t^2 D$$

where

$$D = \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 + \frac{N-K}{N-1} \frac{4}{K} E^2 H^2 + \frac{dM^2}{K}.$$

We decay the learning rate with  $\eta_t = \frac{\beta}{t+\gamma}$  for some  $\beta \geq \frac{1}{\mu}$  and  $\gamma \geq 0$  such that  $\eta_1 \leq \min\{\frac{1}{\mu}, \frac{1}{4L}\} = \frac{1}{4L}$  and  $\eta_t \leq 2\eta_{t+E}$ . Now we prove that  $\Delta_t \leq \frac{v}{\gamma+t}$  where

$$v = \max\left\{\frac{\beta^2 D}{\beta\mu - 1}, (\gamma + 1)\Delta_0\right\}$$

by induction. First, the definition of  $v$  ensures that it holds for  $t = 0$ . Assume the conclusion holds for some  $t > 0$ , it follows that

$$\begin{aligned} \Delta_{t+1} &\leq (1 - \eta_t \mu) \Delta_t + \eta_t^2 D \\ &= \left(1 - \frac{\beta\mu}{t+\gamma}\right) \frac{v}{t+\gamma} + \frac{\beta^2 D}{(t+\gamma)^2} \\ &= \frac{t+\gamma-1}{(t+\gamma)^2} v + \left[\frac{\beta^2 D}{(t+\gamma)^2} - \frac{\mu\beta-1}{(t+\gamma)^2} v\right] \\ &\leq \frac{v}{t+\gamma+1}. \end{aligned}$$

Then by the strong convexity of  $F(\cdot)$ ,

$$\mathbb{E}[F(\bar{\mathbf{w}}_t)] - F^* \leq \frac{L}{2} \Delta_t \leq \frac{L}{2} \frac{v}{\gamma+t}.$$

Specially, if we choose  $\beta = \frac{2}{\mu}$ ,  $\gamma = \max\{8\frac{L}{\mu} - 1, E\}$  and denote  $\kappa = \frac{L}{\mu}$ , then  $\eta_t = \frac{2}{\mu} \frac{1}{\gamma+t}$ . Using  $\max\{a, b\} \leq a + b$ , we have

$$\begin{aligned} v &\leq \frac{\beta^2 D}{\beta\mu - 1} + (\gamma + 1)\Delta_0 \\ &= 4 \frac{D}{\mu^2} + (\gamma + 1)\Delta_0 \\ &\leq 4 \frac{D}{\mu^2} + \left(8 \frac{L}{\mu} - 1 + E + 1\right) \Delta_0 \\ &= 4 \frac{D}{\mu^2} + \left(8 \frac{L}{\mu} + E\right) \|\mathbf{w}_0 - \mathbf{w}^*\|^2. \end{aligned}$$

Therefore,

$$\begin{aligned}\mathbb{E}[F(\bar{\mathbf{w}}_t)] - F^* &\leq \frac{L}{2(\gamma+t)} \left[ 4\frac{D}{\mu^2} + \left(8\frac{L}{\mu} + E\right) \|\mathbf{w}_0 - \mathbf{w}^*\|^2 \right] \\ &= \frac{2\kappa}{\gamma+t} \left[ \frac{D}{\mu} + \left(2L + \frac{E\mu}{4}\right) \|\mathbf{w}_0 - \mathbf{w}^*\|^2 \right].\end{aligned}$$

#### D. Deferred proofs of lemmas

**Proof of Lemma 4.** Let  $S_{t+1}$  denote the set of chosen indexes. Note that the number of possible  $S_{t+1}$  is  $C_N^K$  and we denote the  $l$ th possible result as  $S_{t+1}^l = \{i_1^l, \dots, i_K^l\}$ , where  $l = 1, \dots, C_N^K$ . Therefore,

$$\sum_{j=1}^{C_N^K} \sum_{k=1}^K \mathbf{v}_{t+1}^{i_k^j} = \frac{K \cdot C_N^K}{N} \sum_{i=1}^N \mathbf{v}_{t+1}^i = C_{N-1}^{K-1} \sum_{i=1}^N \mathbf{v}_{t+1}^i.$$

Since when  $t+1 \in \mathcal{I}_E$ ,  $\mathbf{u}_{t+1}^k = \frac{1}{K} \sum_{k \in S_{t+1}} \mathbf{v}_{t+1}^k$  for all  $k$ , we have

$$\bar{\mathbf{u}}_{t+1} = \sum_{k=1}^N \mathbf{u}_{t+1}^k = \frac{1}{K} \sum_{k \in S_{t+1}} \mathbf{v}_{t+1}^k.$$

Then

$$\mathbb{E}_{S_t}[\bar{\mathbf{u}}_{t+1}] = \sum_{l=1}^{C_N^K} \mathbb{P}(S_{t+1} = S_{t+1}^l) \frac{1}{K} \sum_{k \in S_{t+1}^l} \mathbf{v}_{t+1}^k = \frac{1}{C_N^K} \frac{1}{K} \sum_{j=1}^{C_N^K} \sum_{k=1}^K \mathbf{v}_{t+1}^{i_k^j} = \frac{C_{N-1}^{K-1}}{C_N^K} \frac{1}{K} \sum_{k=1}^N \mathbf{v}_{t+1}^k = \bar{\mathbf{v}}_{t+1}.$$

As for the variance, we have [11]

$$\begin{aligned}\mathbb{E}_{S_t} \|\bar{\mathbf{u}}_{t+1} - \bar{\mathbf{v}}_{t+1}\|^2 &= \mathbb{E}_{S_t} \left\| \frac{1}{K} \sum_{i \in S_{t+1}} \mathbf{v}_{t+1}^i - \bar{\mathbf{v}}_{t+1} \right\|^2 = \frac{1}{K^2} \mathbb{E}_{S_t} \left\| \sum_{i=1}^N \mathbb{I}\{i \in S_t\} (\mathbf{v}_{t+1}^i - \bar{\mathbf{v}}_{t+1}) \right\|^2 \\ &= \frac{1}{K^2} \left[ \sum_{i \in [N]} \mathbb{P}(i \in S_{t+1}) \|\mathbf{v}_{t+1}^i - \bar{\mathbf{v}}_{t+1}\|^2 + \sum_{i \neq j} \mathbb{P}(i, j \in S_{t+1}) \langle \mathbf{v}_{t+1}^i - \bar{\mathbf{v}}_{t+1}, \mathbf{v}_{t+1}^j - \bar{\mathbf{v}}_{t+1} \rangle \right] \quad (15) \\ &= \frac{1}{KN} \sum_{i=1}^N \|\mathbf{v}_{t+1}^i - \bar{\mathbf{v}}_{t+1}\|^2 + \sum_{i \neq j} \frac{K-1}{KN(N-1)} \langle \mathbf{v}_{t+1}^i - \bar{\mathbf{v}}_{t+1}, \mathbf{v}_{t+1}^j - \bar{\mathbf{v}}_{t+1} \rangle \\ &= \frac{1 - \frac{K}{N}}{K(N-1)} \sum_{i=1}^N \|\mathbf{v}_{t+1}^i - \bar{\mathbf{v}}_{t+1}\|^2\end{aligned}$$

where we use the following results:

$$\mathbb{P}(i \in S_{t+1}) = \frac{K}{N}$$

and

$$\mathbb{P}(i, j \in S_{t+1}) = \frac{K(K-1)}{N(N-1)}$$

for all  $i \neq j$ , and

$$\sum_{i \in [N]} \|\mathbf{v}_{t+1}^i - \bar{\mathbf{v}}_{t+1}\|^2 + \sum_{i \neq j} \left\langle \mathbf{v}_{t+1}^i - \bar{\mathbf{v}}_{t+1}, \mathbf{v}_{t+1}^j - \bar{\mathbf{v}}_{t+1} \right\rangle = 0.$$

Since  $t+1 \in \mathcal{I}_E$ , we know that  $t_0 = t - E + 1 \in \mathcal{I}_E$  is the communication time, implying that  $\{\mathbf{u}_{t_0}^k\}_{k=1}^N$  are identical. Then

$$\begin{aligned} & \sum_{i=1}^N \|\mathbf{v}_{t+1}^i - \bar{\mathbf{v}}_{t+1}\|^2 \\ &= \sum_{i=1}^N \|(\mathbf{v}_{t+1}^i - \bar{\mathbf{u}}_{t_0}) - (\bar{\mathbf{v}}_{t+1} - \bar{\mathbf{u}}_{t_0})\|^2 \\ &= \sum_{i=1}^N \|\mathbf{v}_{t+1}^i - \bar{\mathbf{u}}_{t_0}\|^2 - 2 \left\langle \sum_{i=1}^N \mathbf{v}_{t+1}^i - \bar{\mathbf{u}}_{t_0}, \bar{\mathbf{v}}_{t+1} - \bar{\mathbf{u}}_{t_0} \right\rangle + \sum_{i=1}^N \|\bar{\mathbf{v}}_{t+1} - \bar{\mathbf{u}}_{t_0}\|^2 \\ &= \sum_{i=1}^N \|\mathbf{v}_{t+1}^i - \bar{\mathbf{u}}_{t_0}\|^2 - \sum_{i=1}^N \|\bar{\mathbf{v}}_{t+1} - \bar{\mathbf{u}}_{t_0}\|^2 \\ &\leq \sum_{i=1}^N \|\mathbf{v}_{t+1}^i - \bar{\mathbf{u}}_{t_0}\|^2 \end{aligned}$$

Taking expectation over the randomness of stochastic gradient on Eqn. (15), we have

$$\begin{aligned} & \mathbb{E} \left[ \frac{1}{K(N-1)} \left(1 - \frac{K}{N}\right) \sum_{k=1}^N \|\mathbf{v}_{t+1}^i - \bar{\mathbf{v}}_{t+1}\|^2 \right] \\ & \leq \frac{N-K}{K(N-1)} \frac{1}{N} \sum_{k=1}^N \mathbb{E} \|\mathbf{v}_{t+1}^i - \bar{\mathbf{u}}_{t_0}\|^2 \\ & \leq \frac{N-K}{K(N-1)} \frac{1}{N} \sum_{k=1}^N E \sum_{i=t_0}^t \mathbb{E} \|\eta_i \nabla F_k(\mathbf{u}_i^k, \xi_i^k)\|^2 \\ & \leq \frac{N-K}{K(N-1)} E^2 \eta_{t_0}^2 H^2 \\ & \leq \frac{N-K}{N-1} \frac{4}{K} E^2 \eta_t^2 H^2 \end{aligned}$$

where the last line is because  $\eta_t$  is non-increasing and  $\eta_{t_0} \leq 2\eta_t$ .

**Proof of Lemma 5.** Let  $w$  be an arbitrary element of  $\mathbf{w}$ . Then  $|w| \leq M$ . With  $B$ -bit quantization we can divide  $[-M, +M]$  into  $\zeta$  smaller intervals  $I_1 = [s_1, s_2], I_2 = [s_2, s_3], \dots, I_\zeta = [s_\zeta, s_{\zeta+1}]$ , with  $\zeta = 2^{B-1}$ . Suppose  $w$  is located at the  $i$ th interval, i.e

$$s_i \leq w \leq s_{i+1}.$$

Using stochastic rounding, we get the quantized result as

$$Q(w) = \begin{cases} s_i, & \text{w.p. } \frac{s_{i+1}-w}{s_{i+1}-s_i}, \\ s_{i+1}, & \text{w.p. } \frac{w-s_i}{s_{i+1}-s_i}. \end{cases}$$

Then

$$\mathbb{E}_{SR}[Q(w)] = s_i \frac{s_{i+1}-w}{s_{i+1}-s_i} + s_{i+1} \frac{w-s_i}{s_{i+1}-s_i} = \frac{w(s_{i+1}-s_i)}{s_{i+1}-s_i} = w,$$

and

$$\begin{aligned} \mathbb{E}_{SR}[(Q(w)-w)^2] &= (s_i-w)^2 \frac{s_{i+1}-w}{s_{i+1}-s_i} + (s_{i+1}-w)^2 \frac{w-s_i}{s_{i+1}-s_i} \\ &= (w-s_i)(s_{i+1}-w) \\ &\leq \left(\frac{s_{i+1}-s_i}{2}\right)^2 = \left(\frac{M}{2^B-1}\right)^2. \end{aligned}$$

Hence, for  $\mathbf{w} = [w_1, w_2, \dots, w_d]$ , we have

$$\mathbb{E}_{SR}[Q(\mathbf{w})] = [\mathbb{E}_{SR}[Q(w_1)], \mathbb{E}_{SR}[Q(w_2)], \dots, \mathbb{E}_{SR}[Q(w_d)]] = \mathbf{w},$$

and

$$\mathbb{E}_{SR}\|Q(\mathbf{w}) - \mathbf{w}\|^2 = \sum_{i=1}^d \mathbb{E}_{SR}[(Q(w_i) - w_i)^2] \leq d \left(\frac{M}{2^B-1}\right)^2.$$

**Proof of Lemma 6.** According to Eqn. (9) and Lemma 5, for  $t+1 \in \mathcal{I}_E$ , we have

$$\begin{aligned} \mathbb{E}_{SR}[\bar{\mathbf{w}}_{t+1}] &= \mathbb{E}_{SR}\left[\frac{1}{K} \sum_{k \in \mathcal{S}_{t+1}} Q(\mathbf{v}_{t+1}^k)\right] \\ &= \frac{1}{K} \sum_{k \in \mathcal{S}_{t+1}} \mathbb{E}_{SR}[Q(\mathbf{v}_{t+1}^k)] \\ &= \frac{1}{K} \sum_{k \in \mathcal{S}_{t+1}} \mathbf{v}_{t+1}^k = \bar{\mathbf{u}}_{t+1}. \end{aligned}$$

As the quantization level is set to  $B_{t+1}$ , with Lemma 5, we know that for all  $k \in [K]$ ,

$$\mathbb{E}_{SR}\|Q(\mathbf{v}_{t+1}^k) - \mathbf{v}_{t+1}^k\|^2 \leq q_{t+1}^2 dM^2 \tag{16}$$

where  $q_{t+1} = 1/(2^{B_{t+1}} - 1)$ . Then

$$\begin{aligned}\mathbb{E}_{SR} \|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{u}}_{t+1}\|^2 &= \mathbb{E}_{SR} \left\| \frac{1}{K} \sum_{k \in \mathcal{S}_{t+1}} Q(\mathbf{v}_{t+1}^k) - \frac{1}{K} \sum_{k \in \mathcal{S}_{t+1}} \mathbf{v}_{t+1}^k \right\|^2 \\ &= \frac{1}{K^2} \mathbb{E}_{SR} \left\| \sum_{k \in \mathcal{S}_{t+1}} (Q(\mathbf{v}_{t+1}^k) - \mathbf{v}_{t+1}^k) \right\|^2.\end{aligned}$$

Let  $\mathbf{e}_{t+1}^k = Q(\mathbf{v}_{t+1}^k) - \mathbf{v}_{t+1}^k$ , then

$$\mathbb{E}_{SR} \|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{u}}_{t+1}\|^2 = \frac{1}{K^2} \sum_{k \in \mathcal{S}_{t+1}} \mathbb{E}_{SR} \|\mathbf{e}_{t+1}^k\|^2 + \frac{1}{K^2} \mathbb{E}_{SR} \left[ \sum_{i,j \in \mathcal{S}_{t+1}, i \neq j} \langle \mathbf{e}_{t+1}^i, \mathbf{e}_{t+1}^j \rangle \right].$$

We know  $\mathbb{E}_{SR} [\mathbf{e}_{t+1}^k] = \mathbf{0}$  from Lemma 5, and  $\mathbf{e}_{t+1}^i$  and  $\mathbf{e}_{t+1}^j$  are independent if  $i \neq j$ . Therefore,

$$\mathbb{E}_{SR} \left[ \sum_{i \neq j} \langle \mathbf{e}_{t+1}^i, \mathbf{e}_{t+1}^j \rangle \right] = \sum_{i \neq j} \mathbb{E}_{SR} [\langle \mathbf{e}_{t+1}^i, \mathbf{e}_{t+1}^j \rangle] = \sum_{i \neq j} \langle \mathbb{E}_{SR}[\mathbf{e}_{t+1}^i], \mathbb{E}_{SR}[\mathbf{e}_{t+1}^j] \rangle = 0. \quad (17)$$

With Eqn. (16), we have

$$\begin{aligned}\mathbb{E}_{SR} \|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{u}}_{t+1}\|^2 &= \frac{1}{K^2} \mathbb{E}_{SR} \sum_{k \in \mathcal{S}_{t+1}} \|\mathbf{e}_{t+1}^k\|^2 \\ &= \frac{1}{K^2} \sum_{k \in \mathcal{S}_{t+1}} \mathbb{E}_{SR} \|Q(\mathbf{v}_{t+1}^k) - \mathbf{v}_{t+1}^k\|^2 \leq q_{t+1}^2 \frac{dM^2}{K}.\end{aligned}$$

## APPENDIX B

### PROOF OF THEOREM 2

#### A. Notations

All of the notations in Appendix A can be extended for DT unless  $\mathbf{w}_{t+1}^k$  is slightly different. For quantized differential transmission, if  $t+1 \in \mathcal{I}_E$ , each client in  $\mathcal{S}_{t+1}$  uploads the quantized differential weights  $Q(\mathbf{d}_{t+1}^k)$  where  $\mathbf{d}_{t+1}^k = \mathbf{v}_{t+1}^k - \mathbf{w}_{t+1-E}$  and  $\mathbf{w}_{t+1-E}$  means the most recent global model it downloaded from the server. And the global aggregation is  $\mathbf{w}_{t+1} = \mathbf{w}_{t+1-E} + \frac{1}{K} \sum_{k \in \mathcal{S}_{t+1}} Q(\mathbf{d}_{t+1}^k)$ . Hence, we can redefine  $\mathbf{w}_{t+1}^k$  as

$$\mathbf{w}_{t+1}^k = \begin{cases} \mathbf{v}_{t+1}^k & \text{if } t+1 \notin \mathcal{I}_E, \\ \mathbf{w}_{t+1-E} + \frac{1}{K} \sum_{k \in \mathcal{S}_{t+1}} Q(\mathbf{d}_{t+1}^k) & \text{if } t+1 \in \mathcal{I}_E. \end{cases}$$

#### B. Lemma

**Lemma 7** (Unbiased and variance bounded quantization). *With stochastic rounding and quantization level  $B$  and assuming the quantization gain for  $\mathbf{d}_{t+1}^k$  is  $G = 2^{B-1} / \|\mathbf{d}_{t+1}^k\|_\infty$ , for all  $t+1 \in \mathcal{I}_E, k \in \mathcal{S}_{t+1}$ ,*

we have

$$\mathbb{E}_{SR} [\bar{\mathbf{w}}_{t+1}] = \bar{\mathbf{u}}_{t+1},$$

and

$$\mathbb{E} \|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{u}}_{t+1}\|^2 \leq \frac{4d}{K(2^B - 1)^2} \eta_t^2 E^2 H^2.$$

**Proof of Lemma 7.** Considering the special case of Lemma 5, say  $M = \|\mathbf{w}\|_\infty$  and the corresponding  $G = 2^{B-1}/\|\mathbf{w}\|_\infty$ , we have

$$\mathbb{E}_{SR} \left[ \|Q(\mathbf{w}) - \mathbf{w}\|^2 \right] \leq d \left( \frac{M}{2^B - 1} \right)^2 = d \frac{\|\mathbf{w}\|_\infty^2}{(2^B - 1)^2} \leq d \frac{\|\mathbf{w}\|^2}{(2^B - 1)^2} \quad (18)$$

Then, for  $t + 1 \in \mathcal{I}_E$ ,

$$\begin{aligned} \bar{\mathbf{w}}_{t+1} &= \frac{1}{N} \sum_{k=1}^N \mathbf{w}_{t+1}^k = \mathbf{w}_{t+1-E} + \frac{1}{K} \sum_{k \in S_{t+1}} Q(\mathbf{d}_{t+1}^k) \\ \bar{\mathbf{u}}_{t+1} &= \frac{1}{N} \sum_{k=1}^N \mathbf{u}_{t+1}^k = \frac{1}{K} \sum_{k \in S_{t+1}} \mathbf{v}_{t+1}^k \end{aligned}$$

Therefore, we get

$$\begin{aligned} \mathbb{E}_{SR} [\bar{\mathbf{w}}_{t+1}] &= \mathbf{w}_{t+1-E} + \frac{1}{K} \sum_{k \in S_{t+1}} \mathbb{E}_{SR} [Q(\mathbf{d}_{t+1}^k)] \\ &= \mathbf{w}_{t+1-E} + \frac{1}{K} \sum_{k \in S_{t+1}} \mathbf{d}_{t+1}^k \\ &= \mathbf{w}_{t+1-E} + \frac{1}{K} \sum_{k \in S_{t+1}} (\mathbf{v}_{t+1}^k - \mathbf{w}_{t+1-E}) \\ &= \frac{1}{K} \sum_{k \in S_{t+1}} \mathbf{v}_{t+1}^k = \bar{\mathbf{u}}_{t+1} \end{aligned}$$

As for the variance, we have

$$\begin{aligned} \mathbb{E}_{SR} \|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{u}}_{t+1}\|^2 &= \mathbb{E}_{SR} \left\| \mathbf{w}_{t+1-E} + \frac{1}{K} \sum_{k \in S_{t+1}} Q(\mathbf{d}_{t+1}^k) - \frac{1}{K} \sum_{k \in S_{t+1}} \mathbf{v}_{t+1}^k \right\|^2 \\ &= \frac{1}{K^2} \mathbb{E}_{SR} \left\| \sum_{k \in S_{t+1}} Q(\mathbf{d}_{t+1}^k) - \sum_{k \in S_{t+1}} (\mathbf{v}_{t+1}^k - \mathbf{w}_{t+1-E}) \right\|^2 \\ &= \frac{1}{K^2} \mathbb{E}_{SR} \left\| \sum_{k \in S_{t+1}} (Q(\mathbf{d}_{t+1}^k) - \mathbf{d}_{t+1}^k) \right\|^2 = \frac{1}{K^2} \sum_{k \in S_{t+1}} \mathbb{E}_{SR} \|Q(\mathbf{d}_{t+1}^k) - \mathbf{d}_{t+1}^k\|^2 \end{aligned}$$

where the last equality is due to  $\mathbb{E}_{SR}[Q(\mathbf{d}_{t+1}^k) - \bar{\mathbf{d}}_{t+1}^k] = \mathbf{0}$  (see the proof of Eqn. (17)). Since we set  $G = 1/\|\mathbf{d}_{t+1}^k\|$  for all  $\mathbf{d}_{t+1}^k$ , with Eqn. (18), we get

$$\begin{aligned} \mathbb{E}_{SR} \|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{u}}_{t+1}\|^2 &= \frac{1}{K^2} \sum_{k \in S_{t+1}} \mathbb{E}_{SR} \left\| Q(\mathbf{d}_{t+1}^k) - \mathbf{d}_{t+1}^k \right\|^2 \\ &\leq \frac{1}{K^2} \sum_{k \in S_{t+1}} \frac{d}{(2^B - 1)^2} \left\| \mathbf{d}_{t+1}^k \right\|^2 \\ &= \frac{d}{K^2 (2^B - 1)^2} \sum_{k \in S_{t+1}} \left\| \sum_{\tau=t+1-E}^t \eta_\tau \nabla F_k(\mathbf{w}_\tau^k, \xi_\tau^k) \right\|^2 \\ &\leq \frac{dE}{K^2 (2^B - 1)^2} \sum_{k \in S_{t+1}} \sum_{\tau=t+1-E}^t \eta_\tau^2 \left\| \nabla F_k(\mathbf{w}_\tau^k, \xi_\tau^k) \right\|^2 \end{aligned}$$

By further taking expectation over the randomness of stochastic gradient, we get

$$\begin{aligned} \mathbb{E} \|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{u}}_{t+1}\|^2 &\leq \frac{dE}{K^2 (2^B - 1)^2} \sum_{k \in S_{t+1}} \sum_{\tau=t+1-E}^t \eta_\tau^2 \mathbb{E}_{SG} \left\| \nabla F_k(\mathbf{w}_\tau^k, \xi_\tau^k) \right\|^2 \\ &\leq \frac{dE}{K^2 (2^B - 1)^2} \sum_{k \in S_{t+1}} \sum_{\tau=t+1-E}^t \eta_{t+1-E}^2 H^2 \\ &= \frac{dE^2}{K (2^B - 1)^2} \eta_{t+1-E}^2 H^2 \leq \frac{4d}{K (2^B - 1)^2} \eta_t^2 E^2 H^2 \end{aligned}$$

where we use the fact that  $\eta_t$  is non-increasing and  $2\eta_{t+1-E} \leq 2\eta_t$ .

### C. Proof of Theorem 2

We use Lemma 7 to update Eqn. (14) to

$$\begin{aligned} \mathbb{E} \|\bar{\mathbf{w}}_{t+1} - \mathbf{w}^*\|^2 &\leq (1 - \eta_t \mu) \mathbb{E} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 \\ &+ \eta_t^2 \left[ \frac{4d}{K(2^B - 1)^2} E^2 H^2 + \frac{\sigma_k^2}{N} + 6L\Gamma + 8(E-1)^2 H^2 + \frac{N-K}{N-1} \frac{4}{K} E^2 H^2 \right]. \end{aligned}$$

Let  $\Delta_t = \mathbb{E} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2$ , therefore for  $t+1 \in \mathcal{I}_E$  or  $t+1 \notin \mathcal{I}_E$ , we have

$$\Delta_{t+1} \leq (1 - \eta_t \mu) \Delta_t + \eta_t^2 D,$$

with

$$D = \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 + \frac{N-K}{N-1} \frac{4}{K} E^2 H^2 + \frac{4d}{K(2^B - 1)^2} E^2 H^2.$$

We can then apply the same induction as in Appendix A-C to get the final result.

## APPENDIX C

## PROOF OF THEOREM 3

**Notations.** Again we extend the notations in Appendix A to downlink quantization. The global model aggregation is  $\mathbf{w}_{t+1} = \frac{1}{K} \sum_{k \in \mathcal{S}_{t+1}} \mathbf{v}_{t+1}^k$  and its quantized version  $Q(\mathbf{w}_{t+1})$  is broadcast to  $K$  randomly selected clients for the next round. All notations are similarly defined. We further note that the analysis of convergence should be on  $\|\bar{\mathbf{u}}_{t+1} - \mathbf{w}^*\|^2$  instead of  $\|\bar{\mathbf{w}}_{t+1} - \mathbf{w}^*\|^2$ , since the server has access to unquantized global model aggregation.

**Proof of Theorem 3.** Under Assumption 1, Lemmas 1 to 3 still hold. We need to consider four cases.

1)  $t+1 \notin \mathcal{I}_E$  and  $t \notin \mathcal{I}_E$ . By taking expectation over all the three randomness, we can get

$$\mathbb{E} \|\bar{\mathbf{v}}_{t+1} - \mathbf{w}^*\|^2 \leq (1 - \eta_t \mu) \mathbb{E} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 + \eta_t^2 \left[ \frac{\sigma_k^2}{N} + 6L\Gamma + 8(E-1)^2 H^2 \right]. \quad (19)$$

since  $t \notin \mathcal{I}_E$ , we have  $\bar{\mathbf{w}}_t = \bar{\mathbf{u}}_t = \bar{\mathbf{v}}_t = \frac{1}{N} \sum_{k=1}^N \mathbf{v}_t^k$ . Hence, we can transform Eqn. (19) into

$$\mathbb{E} \|\bar{\mathbf{u}}_{t+1} - \mathbf{w}^*\|^2 \leq (1 - \eta_t \mu) \mathbb{E} \|\bar{\mathbf{u}}_t - \mathbf{w}^*\|^2 + \eta_t^2 \left[ \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 \right].$$

2)  $t+1 \notin \mathcal{I}_E$  and  $t \in \mathcal{I}_E$ . We still have  $\bar{\mathbf{u}}_{t+1} = \bar{\mathbf{v}}_{t+1}$  and  $\mathbb{E} \|\bar{\mathbf{v}}_{t+1} - \mathbf{w}^*\|^2 = \mathbb{E} \|\bar{\mathbf{u}}_{t+1} - \mathbf{w}^*\|^2$ . But now  $\bar{\mathbf{w}}_t = Q(\frac{1}{K} \sum_{k \in \mathcal{S}_t} \mathbf{v}_t^k)$  and  $\bar{\mathbf{u}}_t = \frac{1}{K} \sum_{k \in \mathcal{S}_t} \mathbf{v}_t^k$ . Under Assumption 2, we have that  $\|\mathbf{v}_t^k\|_\infty \leq M$ , which suggests  $\|\bar{\mathbf{u}}_t\|_\infty \leq M$ . Using Lemma 5, we have

$$\mathbb{E}_{SR} [\bar{\mathbf{w}}_t] = \mathbb{E}_{SR} [Q(\bar{\mathbf{u}}_t)] = \bar{\mathbf{u}}_t \quad (20)$$

$$\mathbb{E}_{SR} \left[ \|\bar{\mathbf{w}}_t - \bar{\mathbf{u}}_t\|^2 \right] = \mathbb{E}_{SR} \left[ \|Q(\bar{\mathbf{u}}_t) - \bar{\mathbf{u}}_t\|^2 \right] \leq d \cdot q_t^2 M^2 \quad (21)$$

where  $q_t = 1/(2^{B_t} - 1)$  and  $B_t$  is the quantization level for the  $t$ th iteration. Therefore,

$$\begin{aligned} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 &= \|\bar{\mathbf{w}}_t - \bar{\mathbf{u}}_t + \bar{\mathbf{u}}_t - \mathbf{w}^*\|^2 \\ &= \underbrace{\|\bar{\mathbf{w}}_t - \bar{\mathbf{u}}_t\|^2}_{A_1} + \underbrace{\|\bar{\mathbf{u}}_t - \mathbf{w}^*\|^2}_{A_2} + 2 \underbrace{\langle \bar{\mathbf{w}}_t - \bar{\mathbf{u}}_t, \bar{\mathbf{u}}_t - \mathbf{w}^* \rangle}_{A_3}. \end{aligned} \quad (22)$$

When expectation is taken over the randomness of stochastic rounding, the last term  $A_3$  vanishes because of Eqn. (20) and  $A_3$  can be bounded using Eqn. (21). We further have

$$\mathbb{E} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 \leq \mathbb{E} \|\bar{\mathbf{u}}_t - \mathbf{w}^*\|^2 + d \cdot q_t^2 M^2, \quad (23)$$

which transforms Eqn. (19) into

$$\begin{aligned} \mathbb{E} \|\bar{\mathbf{u}}_{t+1} - \mathbf{w}^*\|^2 &\leq (1 - \eta_t \mu) \mathbb{E} \|\bar{\mathbf{u}}_t - \mathbf{w}^*\|^2 \\ &+ (1 - \eta_t \mu) d q_t^2 M^2 + \eta_t^2 \left[ \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 \right]. \end{aligned}$$

3)  $t+1 \in \mathcal{I}_E$  and  $t \notin \mathcal{I}_E$ . We still have  $\bar{\mathbf{u}}_{t+1} = \bar{\mathbf{w}}_{t+1}$  and  $\mathbb{E} \|\bar{\mathbf{u}}_{t+1} - \mathbf{w}^*\|^2 = \mathbb{E} \|\bar{\mathbf{w}}_{t+1} - \mathbf{w}^*\|^2$ . But now  $\bar{\mathbf{v}}_{t+1} = \frac{1}{N} \sum_{k=1}^N \mathbf{v}_{t+1}^k$  and  $\bar{\mathbf{u}}_{t+1} = \frac{1}{K} \sum_{k \in \mathcal{S}_{t+1}} \mathbf{v}_{t+1}^k$ , and

$$\begin{aligned} \|\bar{\mathbf{u}}_{t+1} - \mathbf{w}^*\|^2 &= \|\bar{\mathbf{u}}_{t+1} - \bar{\mathbf{v}}_{t+1} + \bar{\mathbf{v}}_{t+1} - \mathbf{w}^*\|^2 \\ &= \underbrace{\|\bar{\mathbf{u}}_{t+1} - \bar{\mathbf{v}}_{t+1}\|^2}_{B_1} + \underbrace{\|\bar{\mathbf{v}}_{t+1} - \mathbf{w}^*\|^2}_{B_2} + 2 \underbrace{\langle \bar{\mathbf{u}}_{t+1} - \bar{\mathbf{v}}_{t+1}, \bar{\mathbf{v}}_{t+1} - \mathbf{w}^* \rangle}_{B_3}. \end{aligned} \quad (24)$$

Lemma 4 indicates  $\mathbb{E}_{\mathcal{S}_t} [\bar{\mathbf{u}}_{t+1} - \bar{\mathbf{v}}_{t+1}] = \mathbf{0}$ , so when expectation is taken over the randomness of  $\mathcal{S}_t$ , the last term  $B_3$  vanishes.  $B_1$  can be bounded by Eqn. (10). We finally have

$$\mathbb{E} \|\bar{\mathbf{u}}_{t+1} - \mathbf{w}^*\|^2 \leq \mathbb{E} \|\bar{\mathbf{v}}_{t+1} - \mathbf{w}^*\|^2 + \frac{N-K}{N-1} \frac{4}{K} \eta_t^2 E^2 H^2, \quad (25)$$

With Eqn. (19), and  $\mathbb{E} \|\bar{\mathbf{w}}_t - \mathbf{w}^*\|^2 = \mathbb{E} \|\bar{\mathbf{u}}_t - \mathbf{w}^*\|^2$  since  $t \notin \mathcal{I}_E$ , we can further have

$$\begin{aligned} \mathbb{E} \|\bar{\mathbf{u}}_{t+1} - \mathbf{w}^*\|^2 &\leq (1 - \eta_t \mu) \mathbb{E} \|\bar{\mathbf{u}}_t - \mathbf{w}^*\|^2 \\ &+ \eta_t^2 \left[ \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 + \frac{N-K}{N-1} \frac{4}{K} E^2 H^2 \right]. \end{aligned}$$

4)  $t+1 \in \mathcal{I}_E$  and  $t \in \mathcal{I}_E$ . This case is only possible for  $E = 1$ . In this case,  $\bar{\mathbf{v}}_{t+1} \neq \bar{\mathbf{u}}_{t+1}$  and  $\bar{\mathbf{u}}_{t+1} \neq \bar{\mathbf{w}}_{t+1}$ . We use both Eqn. (23) and Eqn. (25) to transform Eqn. (19) into

$$\begin{aligned} \mathbb{E} \|\bar{\mathbf{u}}_{t+1} - \mathbf{w}^*\|^2 &\leq (1 - \eta_t \mu) \mathbb{E} \|\bar{\mathbf{u}}_t - \mathbf{w}^*\|^2 + (1 - \eta_t \mu) d \cdot q_t^2 M^2 \\ &+ \eta_t^2 \left[ \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 + \frac{N-K}{N-1} \frac{4}{K} E^2 H^2 \right]. \end{aligned} \quad (26)$$

In summary, Eqn. (26) holds for all cases. Let  $\Delta_t = \mathbb{E} \|\bar{\mathbf{u}}_t - \mathbf{w}^*\|^2$ . If we increase the quantization level  $B_t$  following

$$B_t = \log_2 \left( 1 + \frac{\sqrt{1 - \eta_t \mu}}{\eta_t} \right)$$

to make

$$q_t = 1/(2^{B_t} - 1) = \frac{\eta_t}{\sqrt{1 - \eta_t \mu}}.$$

Then we have  $(1 - \eta_t \mu)q_t^2 = \eta_t^2$  and we also have

$$\Delta_{t+1} \leq (1 - \eta_t \mu)\Delta_t + \eta_t^2 D$$

where

$$D = \sum_{k=1}^N \frac{\sigma_k^2}{N^2} + 6L\Gamma + 8(E-1)^2 H^2 + \frac{N-K}{N-1} \frac{4}{K} E^2 H^2 + dM^2.$$

Applying the same induction method in Appendix A-C proves the theorem.

## REFERENCES

- [1] S. Niknam, H. S. Dhillon, and J. H. Reed, “Federated learning for wireless communications: Motivation, opportunities and challenges,” *arXiv preprint arXiv:1908.06847*, 2019.
- [2] W. Y. B. Lim, N. C. Luong, D. T. Hoang, Y. Jiao, Y.-C. Liang, Q. Yang, D. Niyato, and C. Miao, “Federated learning in mobile edge networks: A comprehensive survey,” *IEEE Commun. Surveys Tuts.*, 2020.
- [3] G. Zhu, D. Liu, Y. Du, C. You, J. Zhang, and K. Huang, “Toward an intelligent edge: Wireless communication meets machine learning,” *IEEE Commun. Mag.*, vol. 58, no. 1, pp. 19–25, 2020.
- [4] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, Fort Lauderdale, FL, USA, Apr. 2017, pp. 1273–1282.
- [5] J. Konecny, H. B. McMahan, F. X. Yu, P. Richtarik, A. T. Suresh, and D. Bacon, “Federated learning: Strategies for improving communication efficiency,” in *NIPS Workshop on Private Multi-Party Machine Learning*, 2016.
- [6] 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, 2020.
- [7] P. Kairouz *et al.*, “Advances and open problems in federated learning,” *arXiv preprint arXiv:1912.04977*, 2019.
- [8] G. Zhu, Y. Wang, and K. Huang, “Broadband analog aggregation for low-latency federated edge learning,” *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 491–506, 2019.
- [9] Y. Du, S. Yang, and K. Huang, “High-dimensional stochastic gradient quantization for communication-efficient edge learning,” *IEEE Trans. Signal Processing*, vol. 68, pp. 2128–2142, 2020.
- [10] N. Shlezinger, M. Chen, Y. C. Eldar, H. V. Poor, and S. Cui, “Federated learning with quantization constraints,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 8851–8855.
- [11] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, “On the convergence of FedAvg on non-IID data,” in *International Conference on Learning Representations*, 2020.
- [12] J. Verbraeken, M. Wolting, J. Katzy, J. Kloppenburg, T. Verbelen, and J. S. Rellermeyer, “A survey on distributed machine learning,” *ACM Comput. Surv.*, vol. 53, no. 2, March 2020.
- [13] S. U. Stich, “Local SGD converges fast and communicates little,” in *International Conference on Learning Representations*, 2018.
- [14] J. Wang and G. Joshi, “Cooperative SGD: A unified framework for the design and analysis of communication-efficient SGD algorithms,” in *ICML Workshop on Coding Theory for Machine Learning*, 2019.
- [15] F. Haddadpour, M. M. Kamani, M. Mahdavi, and V. Cadambe, “Local SGD with periodic averaging: Tighter analysis and adaptive synchronization,” in *Advances in Neural Information Processing Systems*, 2019, pp. 11 080–11 092.

- [16] A. Bhowmick, J. Duchi, J. Freudiger, G. Kapoor, and R. Rogers, "Protection against reconstruction and its applications in private federated learning," *arXiv preprint arXiv:1812.00984*, 2018.
- [17] S. Truex, N. Baracaldo, A. Anwar, T. Steinke, H. Ludwig, R. Zhang, and Y. Zhou, "A hybrid approach to privacy-preserving federated learning," in *Proceedings of the 12th ACM Workshop on Artificial Intelligence and Security*, 2019, pp. 1–11.
- [18] C. Niu, F. Wu, S. Tang, L. Hua, R. Jia, C. Lv, Z. Wu, and G. Chen, "Secure federated submodel learning," *arXiv preprint arXiv:1911.02254*, 2019.
- [19] C. Xie, "Zeno++: robust asynchronous SGD with arbitrary number of byzantine workers," *arXiv preprint arXiv:1903.07020*, 2019.
- [20] C. Xie, S. Koyejo, and I. Gupta, "Practical distributed learning: Secure machine learning with communication-efficient local updates," *arXiv preprint arXiv:1903.06996*, 2019.
- [21] K. Bonawitz *et al.*, "Towards federated learning at scale: System design," in *Proceedings of the 2nd SysML Conference*, 2019, pp. 1–15.
- [22] T. Li, M. Sanjabi, and V. Smith, "Fair resource allocation in federated learning," *arXiv preprint arXiv:1905.10497*, 2019.
- [23] X. Mo and J. Xu, "Energy-efficient federated edge learning with joint communication and computation design," *arXiv preprint arXiv:2003.00199*, 2020.
- [24] Q. Zeng, Y. Du, K. K. Leung, and K. Huang, "Energy-efficient radio resource allocation for federated edge learning," *arXiv preprint arXiv:1907.06040*, 2019.
- [25] M. Chen, Z. Yang, W. Saad, C. Yin, H. V. Poor, and S. Cui, "A joint learning and communications framework for federated learning over wireless networks," *arXiv preprint arXiv:1909.07972*, 2019.
- [26] W. Shi, S. Zhou, and Z. Niu, "Device scheduling with fast convergence for wireless federated learning," *arXiv preprint arXiv:1911.00856*, 2019.
- [27] Z. Yang, M. Chen, W. Saad, C. S. Hong, and M. Shikh-Bahaei, "Energy efficient federated learning over wireless communication networks," *arXiv preprint arXiv:1911.02417*, 2019.
- [28] M. Chen, H. V. Poor, W. Saad, and S. Cui, "Convergence time optimization for federated learning over wireless networks," *arXiv preprint arXiv:2001.07845*, 2020.
- [29] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," in *Proceedings of the 3rd MLSys Conference*, 2020.
- [30] 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. Select. Areas Commun.*, vol. 37, no. 6, pp. 1205–1221, 2019.
- [31] 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*, 2017, pp. 1709–1720.
- [32] J. Konečný, H. B. McMahan, D. Ramage, and P. Richtárik, "Federated optimization: Distributed machine learning for on-device intelligence," *arXiv preprint arXiv:1610.02527*, 2016.
- [33] S. Caldas, J. Konečný, H. B. McMahan, and A. Talwalkar, "Expanding the reach of federated learning by reducing client resource requirements," *arXiv preprint arXiv:1812.07210*, 2018.
- [34] J. Bernstein, Y.-X. Wang, K. Azizzadenesheli, and A. Anandkumar, "signSGD: Compressed optimisation for non-convex problems," in *International Conference on Machine Learning*, 2018, pp. 560–569.
- [35] A. Reiszadeh, A. Mokhtari, H. Hassani, and R. Pedarsani, "An exact quantized decentralized gradient descent algorithm," *IEEE Trans. Signal Processing*, vol. 67, no. 19, pp. 4934–4947, 2019.
- [36] G. Zhu, Y. Du, D. Gunduz, and K. Huang, "One-bit over-the-air aggregation for communication-efficient federated edge learning: Design and convergence analysis," *arXiv preprint arXiv:2001.05713*, 2020.

- [37] A. Reisizadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani, “FedPAQ: A communication-efficient federated learning method with periodic averaging and quantization,” in *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics (AISTATS)*, Palermo, Sicily, Italy, 2020.
- [38] M. M. Amiri and D. Gündüz, “Federated learning over wireless fading channels,” *IEEE Trans. Wireless Commun.*, vol. 19, no. 5, pp. 3546–3557, 2020.
- [39] M. M. Amiri, D. Gunduz, S. R. Kulkarni, and H. V. Poor, “Federated learning with quantized global model updates,” *arXiv preprint arXiv:2006.10672*, 2020.
- [40] J. Xu and H. Wang, “Client selection and bandwidth allocation in wireless federated learning networks: A long-term perspective,” *arXiv preprint arXiv:2004.04314*, 2020.
- [41] V. Smith, S. Forte, C. Ma, M. Takáč, M. I. Jordan, and M. Jaggi, “CoCoA: A general framework for communication-efficient distributed optimization,” *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 8590–8638, 2017.
- [42] D. Tse and P. Viswanath, *Fundamentals of Wireless Communication*. Cambridge University Press, 2005.
- [43] S. Gupta, A. Agrawal, K. Gopalakrishnan, and P. Narayanan, “Deep learning with limited numerical precision,” in *International Conference on Machine Learning*, 2015, pp. 1737–1746.
- [44] D. Lin, S. Talathi, and S. Annapureddy, “Fixed point quantization of deep convolutional networks,” in *International Conference on Machine Learning*, 2016, pp. 2849–2858.
- [45] S. Han, H. Mao, and W. J. Dally, “Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding,” in *International Conference on Learning Representations*, 2016.
- [46] G. Montorsi and S. Benedetto, “Design of fixed-point iterative decoders for concatenated codes with interleavers,” *IEEE Journal on Selected Areas in Communications*, vol. 19, no. 5, pp. 871–882, 2001.
- [47] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, “Federated learning with non-IID data,” *arXiv preprint arXiv:1806.00582*, 2018.
- [48] K. Bonawitz, F. Salehi, J. Konečný, B. McMahan, and M. Gruteser, “Federated learning with autotuned communication-efficient secure aggregation,” in *2019 53rd Asilomar Conference on Signals, Systems, and Computers*. IEEE, 2019, pp. 1222–1226.
- [49] S. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, and H. B. McMahan, “Adaptive federated optimization,” 2020.
- [50] 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, 1998.
- [51] A. Krizhevsky, “Learning multiple layers of features from tiny images,” University of Toronto, Tech. Rep., April 2009.
- [52] W. Shakespeare, “The complete works of William Shakespeare,” [EB/OL], <http://www.gutenberg.org/ebooks/100>, Accessed June 23, 2020.
- [53] S. Caldas *et al.*, “LEAF: A benchmark for federated settings,” *arXiv preprint arXiv:1812.01097*, 2018.