# Communication-Efficient Stochastic Zeroth-Order Optimization for Federated Learning

Wenzhi Fang, Ziyi Yu, Yuning Jiang, *Member, IEEE*, Yuanming Shi, *Senior Member, IEEE*,
Colin N. Jones, *Senior Member, IEEE*, and Yong Zhou, *Senior Member, IEEE*

*Abstract*—Federated learning (FL), as an emerging edge artificial intelligence paradigm, enables many edge devices to collaboratively train a global model without sharing their private data. To enhance the training efficiency of FL, various algorithms have been proposed, ranging from first-order to second-order methods. However, these algorithms cannot be applied in scenarios where the gradient information is not available, e.g., federated black-box attack and federated hyperparameter tuning. To address this issue, in this paper we propose a derivative-free federated zeroth-order optimization (FedZO) algorithm featured by performing multiple local updates based on stochastic gradient estimators in each communication round and enabling partial device participation. Under non-convex settings, we derive the convergence performance of the FedZO algorithm on non-independent and identically distributed data and characterize the impact of the numbers of local iterates and participating edge devices on the convergence. To enable communication-efficient FedZO over wireless networks, we further propose an over-the-air computation (AirComp) assisted FedZO algorithm. With an appropriate transceiver design, we show that the convergence of AirComp-assisted FedZO can still be preserved under certain signal-to-noise ratio conditions. Simulation results demonstrate the effectiveness of the FedZO algorithm and validate the theoretical observations.

*Index Terms*—Federated learning, zeroth-order optimization, convergence, over-the-air computation.

## I. INTRODUCTION

With the rapid advancement of the Internet of Things (IoT), a massive amount of data is generated and collected by various edge devices (e.g., sensors, smart phones). Because of the limited radio spectrum resource and increasing privacy concerns, gathering geographically distributed data from a large number of edge devices into a cloud server to enable cloud artificial intelligence (AI) may not be practical. To this end, edge AI has recently been envisioned as a promising AI paradigm [1]. Unlike cloud AI that relies on a cloud server to conduct centralized training, edge AI exploits the computing power of multiple edge devices to perform model training

with their own local data in a distributed manner. Federated learning (FL) [2], as a representative edge AI framework, enables multiple edge devices to collaboratively train a shared model without exchanging their local data, which effectively alleviates the communication burden and privacy concerns. Nowadays, FL has found application in various fields, including autonomous driving [3], recommendation systems [4], healthcare informatics [5], etc.

As a result of the popularity of FL, the federated optimization problem for model training has attracted a growing body of attention from both academia and industry in recent years. Various algorithms have been proposed to attain a fast convergence rate and reduce the communication load, including both first- (e.g., FedAvg [2], FedPD [6], FedNova [7]) and second-order algorithms (e.g., FedDANE [8]). Most existing algorithms rely on gradient and/or Hessian information to solve the federated optimization problem. However, such information cannot be obtained in scenarios where the analytic expressions of the loss functions are unavailable, such as federated hyperparameter tuning [9] or distributed black-box attack of deep neural networks (DNN) [10]. In other words, existing algorithms cannot tackle federated optimization problems when gradient information is not available. This motivates us to develop a communication-efficient federated zeroth-order optimization algorithm that does not require gradient or Hessian information.

Parallel with the research on algorithm design for FL, the implementation of FL over wireless networks is also an emerging research topic. Random channel fading and receiver noise raise unique challenges for the training of FL over wireless networks. Guaranteeing the learning performance with limited radio resource is a challenging task, which requires the joint design of the learning algorithm and communication strategy. Along this line of research, the authors in [11] studied the joint resource allocation and edge device selection to enhance learning performance. Both studies adopted the orthogonal multiple access (OMA) scheme, where the number of edge devices that can participate in each communication round is restricted by the number of available time/frequency resource blocks. The limited radio resource turns out to be the main performance bottleneck of wireless FL. Fortunately, over-the-air computation (AirComp), as a non-orthogonal multiple access scheme, allows concurrent transmissions over the same radio channel to enable low-latency and spectrum-efficient wireless data aggregation [12]–[14], thereby mitigating the communication bottleneck [15]. Motivated by this observation, various AirComp-assisted FL algorithms were proposed in

[16]–[18] to achieve fast model aggregation, wherein all of them adopted first-order optimization algorithm. There still lacks a thorough investigation on AirComp-assisted FL with zeroth-order optimization.

## A. Main Contributions

In this paper, we consider a federated optimization problem, where the gradient of the loss function is not available. We propose a derivative-free federated zeroth-order optimization algorithm, named FedZO, whose key features are performing multiple local updates based on a stochastic gradient estimator in each communication round and enabling partial device participation. We establish a convergence guarantee for the proposed FedZO algorithm and study its implementation over wireless networks with the assistance of AirComp, which is challenging for the following reasons. First, although executing multiple local iterates in each communication round reduces the communication overhead, it also increases the discrepancies among local models due to the data heterogeneity and may even lead to algorithmic divergence. To reduce the communication overhead while preserving the convergence, the relationship between the convergence behavior and the number of local iterates needs to be characterized. Second, the stochastic gradient estimator adopted in FedZO is not an unbiased estimate of the actual gradient, as demonstrated in [19], [20]. These unique features together with multiple local iterates and partial device participation per communication round make the existing convergence analysis framework for FedAvg not applicable to the proposed FedZO algorithm. Third, to characterize the convergence of the AirComp-assisted FedZO algorithm, the impact of the random channel fading and receiver noise in global model aggregation needs to be further taken into account. This not only complicates the convergence analysis but also poses a new problem for the communication strategy design. In this paper, we develop a unified convergence analysis framework to address the aforementioned challenges. The main contributions of this paper are summarized as follows:

- We develop the derivative-free FedZO algorithm, which inherits the framework of the FedAvg algorithm but only queries the values of the objective function, to handle federated optimization problems without using gradient or Hessian information. To cater for the FL system with a large number of edge devices and to reduce the communication overhead, the proposed FedZO algorithm enables partial device participation and performs multiple local iterates in each communication round.
- We establish a convergence guarantee for the proposed FedZO algorithm under non-convex settings and on non-independent and identically distributed (non-i.i.d.) data, and then derive the maximum number of local iterates required for preserving convergence. We demonstrate that the proposed FedZO algorithm can attain linear speedup in the number of local iterates and the number of participating edge devices.
- We study the implementation of the FedZO algorithm over wireless networks with the assistance of AirComp

for the aggregation of local model updates in the uplink. With an appropriate transceiver design that can mitigate the impact of the fading and noise perturbation, we study the convergence behavior of the AirComp-assisted FedZO algorithm and characterize the impact of the signal-to-noise ratio (SNR) on the convergence performance.

We conduct extensive simulations to evaluate the performance of the proposed FedZO and AirComp-assisted FedZO algorithms. Simulation results show that the proposed FedZO algorithm is convergent under various parameter settings and outperforms existing distributed zeroth-order methods. Moreover, simulations illustrate that the performance of the proposed FedZO algorithm is comparable to that of the FedAvg algorithm, which indicates that our proposed algorithm can serve as a satisfactory alternative for FedAvg when first-order information is not available. Results also confirm that, with an appropriate SNR setting, the AirComp-assisted FedZO algorithm preserves convergence.

## B. Related Works

The study of FL started from the seminal work [2], where the authors proposed a communication-efficient federated optimization algorithm known as FedAvg. Subsequently, various articles established convergence guarantees for the FedAvg algorithm [21]–[23]. Following the FedAvg algorithm, many other first-order methods have been proposed, e.g., FedPD [6], FedNova [7], FedProx [24], SCAFFOLD [25], and FedSplit [26]. To further reduce the communication overhead, several second-order optimization algorithms were proposed, such as FedDANE [8] and GIANT [27]. Although the aforementioned first- and second-order algorithms have broad applications, there are still many FL tasks where the gradient and Hessian information are unavailable and thus require zeroth-order optimization.

Recently, several works [10], [28]–[33] focused on studying distributed zeroth-order optimization. Specifically, the authors in [28] developed a so-called ZONE-S algorithm based on the primal-dual technique. In [29], the authors employed the gradient tracking technique to develop a fast distributed zeroth-order algorithm. However, ZONE-S requires $\mathcal{O}(T)$ ($T$ denotes the number of total iterations) sampling complexity per iteration while the algorithm proposed in [29] considered the deterministic setting. More recently, the authors in [10] proposed an algorithm with $\mathcal{O}(1)$ sampling complexity per iteration, which attains linear speedup in the number of edge devices. [30] proposed a decentralized zeroth-order algorithm that allows multiple local updates. However, the theoretic analysis in [30] focused on the strongly convex scenario and relied on the Lipschitzness of local functions, which is relatively restrictive [34]. The authors in [31] proposed and analyzed a distributed zeroth-order Frank-Wolfe algorithm for constrained optimization. Based on single-point and Kiefer-Wolfowitz type gradient estimators, the authors in [32], [33] proposed two distributed zeroth-order algorithms over time-varying graphs. It is worth noting that the aforementioned works mainly consider the peer-to-peer architecture while the studies on the central-server-based architecture are very

limited. Moreover, most of the existing distributed zeroth-order algorithms focused on full device participation, which may not be practical for FL systems with limited radio resources and a large number of edge devices [22].

AirComp has recently been adopted to support the implementation of FL over wireless networks [16], [35]–[39], where channel fading and receiver noise inevitably distort the model aggregation, and in turn introduce a detrimental impact on the learning performance [16]. The convergence behavior of various FL algorithms, e.g., vanilla gradient method [35] and stochastic gradient method [36], showed that the channel fading and noise perturbation typically introduce a non-diminishing optimality gap, which can be mitigated by transmit power control [36], beamforming design [35], and device scheduling [37]. In [38], the authors proposed a joint learning and transmission scheme to ensure global convergence for strongly convex problems. By utilizing the communication strategy in [38], the authors in [39] developed an AirComp-assisted accelerated gradient descent algorithm. Despite the above progress, the existing works focused on the first-order method, while there is no relevant literature studying the AirComp-assisted zeroth-order optimization algorithm.

### C. Organization

The remainder of this paper is organized as follows. We present the problem formulation and propose a federated zeroth-order optimization algorithm in Section II. Section III provides the convergence analysis. Section IV studies the implementation of the proposed FedZO algorithm over wireless networks using AirComp. The simulation results are provided in Section V. Finally, we conclude this paper in Section VI. **Notation:** We denote the $\ell_2$ norm of vectors by $\| \cdot \|$. $[T]$ denotes the set $\{1, 2, \ldots, T-1\}$. $\mathbb{S}^d = \{\boldsymbol{v} \in \mathbb{R}^d \mid \|\boldsymbol{v}\| = 1\}$ denotes a $d$-dimensional unit sphere. $\mathbb{B}^d = \{\boldsymbol{v} \in \mathbb{R}^d \mid \|\boldsymbol{v}\| \leq 1\}$ denotes a $d$-dimensional unit ball. We denote uniform distributions over $\mathbb{S}^d$ and $\mathbb{B}^d$ by $\mathcal{U}(\mathbb{S}^d)$ and $\mathcal{U}(\mathbb{B}^d)$, respectively. We denote $\sim$ as the uniform sampling. For a function $F$, $\nabla F$ and $\widetilde{\nabla} F$ denote the gradient and gradient estimator, respectively.

## II. FEDERATED ZEROTH-ORDER OPTIMIZATION

In this section, we first introduce the federated optimization problem and then propose a federated zeroth-order optimization algorithm.

### A. Problem Formulation

Consider an FL task over a network consisting of a central server and $N$ edge devices indexed by $\{1, 2, \ldots, N\}$. The goal of the central server is to coordinate all edge devices to collaboratively solve the following federated optimization problem

$$\min_{\boldsymbol{x} \in \mathbb{R}^d} f(\boldsymbol{x}) \triangleq \frac{1}{N} \sum_{i=1}^{N} f_i(\boldsymbol{x}), \qquad (1)$$

where $\boldsymbol{x} \in \mathbb{R}^d$ denotes the model parameter of dimension $d$, and $f_i(\boldsymbol{x})$ and $f(\boldsymbol{x})$ denote the local loss function of edge device $i$ and the global loss function at the central server evaluated at model parameter $\boldsymbol{x}$, respectively. We assume that each edge device with a local dataset is equally important for the global model [23]. In (1), $f_i(\boldsymbol{x})$ measures the expected risk over the local data distribution denoted as $\mathcal{D}_i$ at edge device $i$, given by

$$f_i(\boldsymbol{x}) \triangleq \mathbb{E}_{\xi_i \sim \mathcal{D}_i}[F_i(\boldsymbol{x}, \xi_i)],$$

where $F_i(\boldsymbol{x}, \xi_i)$ represents the loss with respect to $\xi_i$ evaluated at model parameter $\boldsymbol{x}$ and $\xi_i$ denotes a random variable uniformly distributed over $\mathcal{D}_i$. In particular, a realization of $\xi_i$ is a single data sample. With sampling information of $\xi_i$ and $\boldsymbol{x}$, edge device $i$ can query the function value of $F_i$, which serves as a stochastic approximation of the expected loss $f_i(\boldsymbol{x})$. Note that the analytic expression and gradient information of $F_i$ are not available.

**Remark 1.** *The scenarios where the gradient information is not available arise in many practical applications [40], including but not limited to federated black-box attacks of DNN [10] and federated hyperparameter tuning in model training [9]. To be specific, in federated black-box attacks, the gradient information cannot be acquired as the deep model is hidden. In the federated hyperparameter tuning task, there does not exist an analytic relationship between the training loss and the hyperparameters.*

### B. Preliminaries on Stochastic Gradient Estimator

We adopt a mini-batch-type stochastic gradient estimator [41]. Specifically, for function $F_i$, the mini-batch-type stochastic gradient estimator is given by

$$\widetilde{\nabla} F_i\left(\boldsymbol{x}, \{\xi_{i,m}\}_{m=1}^{b_1}, \{\boldsymbol{v}_{i,n}\}_{n=1}^{b_2}, \mu\right)$$
$$= \frac{1}{b_1 b_2} \sum_{m=1}^{b_1} \sum_{n=1}^{b_2} \frac{d\boldsymbol{v}_{i,n}}{\mu}\Big(F_i(\boldsymbol{x}+\mu\boldsymbol{v}_{i,n}, \xi_{i,m}) - F_i(\boldsymbol{x}, \xi_{i,m})\Big), \quad (2)$$

where $\{\xi_{i,m}\}_{m=1}^{b_1}$ is a sequence of independent and identically distributed (i.i.d.) random variables with the same distribution as $\xi_i$, $\{\boldsymbol{v}_{i,n}\}_{n=1}^{b_2}$ is a sequence of i.i.d. random vectors with distribution $\mathcal{U}(\mathbb{S}^d)$, and $\mu$ is a positive step size. It has been shown in [20] that

$$\mathbb{E}\left[\frac{d\boldsymbol{v}_{i,n}}{\mu}\Big(F_i(\boldsymbol{x}+\mu\boldsymbol{v}_{i,n}, \xi_{i,m}) - F_i(\boldsymbol{x}, \xi_{i,m})\Big)\right] = \nabla f_i^\mu(\boldsymbol{x}), \quad (3)$$

where the expectation is taken over $\{\xi_{i,m}, \boldsymbol{v}_{i,n}\}$ and $f_i^\mu(\boldsymbol{x}) = \mathbb{E}_{\boldsymbol{u} \sim \mathcal{U}(\mathbb{B}^d)}[f_i(\boldsymbol{x} + \mu\boldsymbol{u})]$ is a locally averaged version of $f_i(\boldsymbol{x})$. Furthermore, we have

$$\mathbb{E}\left[\widetilde{\nabla} F_i\left(\boldsymbol{x}, \{\xi_{i,m}\}_{m=1}^{b_1}, \{\boldsymbol{v}_{i,n}\}_{n=1}^{b_2}, \mu\right)\right] = \nabla f_i^\mu(\boldsymbol{x}). \quad (4)$$

where the expectation is taken over $\{\xi_{i,m}\}_{m=1}^{b_1}$ and $\{\boldsymbol{v}_{i,n}\}_{n=1}^{b_2}$. As $\nabla f_i^\mu(\boldsymbol{x})$ is a biased approximation of $\nabla f_i(\boldsymbol{x})$ [40], (2) is a biased estimate of the actual gradient.

The mini-batch-type stochastic gradient estimator enjoys a low variance than the two-point stochastic gradient estimator [28], [31], [41]. Besides, the mini-batch sizes $b_1$ and $b_2$ are independent of $d$, and hence the computation complexity of (2) does not scale with the dimension of variable $\boldsymbol{x}$.

## C. FedZO Algorithm

Inspired by the FedAvg algorithm, we develop a federated zeroth-order optimization algorithm summarized in Algorithm 1. The main idea of the FedZO algorithm is to get rid of the dependence on the gradient and reduce the frequency of model exchanges, which are achieved by employing a gradient estimator (2) and performing $H$ steps of stochastic zeroth-order updates per communication round, respectively. The FedZO algorithm consists of the following four phases in each round.

- *Global Model Dissemination:* At the beginning of the $t$-th round, the central server uniformly samples $M$ edge devices to participate in the local training. The set of scheduled edge devices in round $t$ is denoted as $\mathcal{M}_t$. Then, the central server disseminates its current global model parameter $x^t$ to the selected edge devices.
- *Local Model Update:* After receiving the model parameter $x^t$ from the central server, each edge device $i \in \mathcal{M}_t$ initializes its local model $x_i^{(t,0)}$ with the received global model from the central server, i.e., $x_i^{(t,0)} = x^t$, and then takes a total of $H$ iterates of stochastic zeroth-order updates. In particular, at the $k$-th iteration of the $t$-th round, edge device $i$ computes a stochastic gradient estimator according to (2). For notational ease, we denote it as

$$e_i^{(t,k)} = \widetilde{\nabla} F_i \left( x_i^{(t,k)}, \{\xi_{i,m}^{(t,k)}\}_{m=1}^{b_1}, \{v_{i,n}^{(t,k)}\}_{n=1}^{b_2}, \mu \right), \quad (5)$$

where $x_i^{(t,k)}$ represents the local model of edge device $i$ at the $k$-th iteration of the $t$-th round. Subsequently, the sampled edge devices update their local models by performing the following stochastic zeroth-order update

$$x_i^{(t,k+1)} = x_i^{(t,k)} - \eta e_i^{(t,k)}, k = 0, 1, \ldots, H-1, \quad (6)$$

where $\eta$ denotes the learning rate. After $H$ local iterates, edge device $i$ obtains an updated local model $x_i^{(t,H)}$.
- *Local Model Uploading:* All edge devices in set $\mathcal{M}_t$ calculate the updates of their local models in this round, i.e., $\Delta_i^t = x_i^{(t,H)} - x_i^{(t,0)}$, $i \in \mathcal{M}_t$, and then upload these updates to the central server.
- *Global Model Update:* After receiving local model updates from the sampled edge devices, the central server aggregates these updates, i.e., $\Delta^t = \frac{1}{M} \sum_{i \in \mathcal{M}_t} \Delta_i^t$, and then updates the global model, i.e., $x^{t+1} = x^t + \Delta^t$.

Although the proposed FedZO algorithm adopts a similar framework as the FedAvg algorithm, the convergence analysis of the latter cannot be directly extended to that of the FedZO algorithm. The key factor hindering the extension is that the gradient estimator does not preserve specific properties of the stochastic gradient. For instance, the gradient estimator (5) is not an unbiased estimate of the true gradient. Besides, the existing theoretical analysis framework for the distributed zeroth-order optimization method cannot be applied to the FedZO algorithm as existing zeroth-order algorithms [10], [28], [29] do not consider multiple steps of local model updates and partial device participation. A larger number of local iterates reduces the communication overhead, but also increases the

---

**Algorithm 1:** FedZO Algorithm

**Input:** Initial model $x^0$, learning rate $\eta$, step size $\mu$,
    mini-batch sizes $b_1$, $b_2$,
number of participating edge devices $M$ **for**
$t \in \{0, 1, \ldots, T-1\}$ **do**
    Uniformly sample a subset $\mathcal{M}_t$ of $M$ edge devices
    Disseminate global model $x^t$ to all edge devices in
    set $\mathcal{M}_t$
    **for** *edge device* $i \in \mathcal{M}_t$ *in parallel* **do**
        Initialize local model $x_i^{(t,0)} = x^t$
        **for** $k = 0, \ldots, H-1$ **do**
            Generate $\{\xi_{i,m}^{(t,k)}\}_{m=1}^{b_1} \sim \mathcal{D}_i$ independently
            Generate $\{v_{i,n}^{(t,k)}\}_{n=1}^{b_2} \sim \mathcal{U}(\mathbb{S}^d)$ independently
            Update $x_i^{(t,k+1)}$ by (6)
        **end**
        Compute local model updates
        $\Delta_i^t = x_i^{(t,H)} - x_i^{(t,0)}$
        Upload local model updates to central server
    **end**
    Aggregate local changes $\Delta^t = \frac{1}{M} \sum_{i \in \mathcal{M}_t} \Delta_i^t$
    Update global model $x^{t+1} = x^t + \Delta^t$
**end**

---

local model discrepancies and may even lead to algorithmic divergence. To preserve convergence for the developed FedZO algorithm, it is necessary to bound these discrepancies by appropriately choosing the number of local updates, i.e., $H$. In Section III, we will provide the convergence analysis for the FedZO algorithm.

## III. CONVERGENCE ANALYSIS FOR FEDZO

In this section, we present the convergence analysis of the FedZO algorithm with full and partial device participation. To make our analysis applicable for more practical scenarios, we focus on the settings of non-convex loss functions and the non-i.i.d. data. We make the following assumptions for the tractability of convergence analysis.

**Assumption 1.** *The global loss in* (1), *i.e.,* $f(x)$, *is lower bounded by* $f_*$, *i.e.,* $f(x) \geq f_* > -\infty$.

**Assumption 2.** $F_i(x, \xi_i)$, $f_i(x)$, *and* $f(x)$ *are* $L$-*smooth. Mathematically, for any* $x \in \mathbb{R}^d$ *and* $y \in \mathbb{R}^d$, *we have*

$$\|\nabla f_i(y) - \nabla f_i(x)\| \leq L\|y - x\|, \ \forall i,$$

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2}\|y - x\|^2.$$

**Assumption 3.** *The second-order moment of stochastic gradient* $\nabla F_i(x, \xi_i)$ *satisfies* $\mathbb{E}_{\xi_i}\|\nabla F_i(x, \xi_i)\|^2 \leq c_g\|\nabla f_i(x)\|^2 + \sigma_g^2$, $\forall x \in \mathbb{R}^d$, $\forall i$, *where* $c_g \geq 1$.

**Assumption 4.** *The gradient dissimilarity between each local loss function and the global loss function is bounded as* $\|\nabla f(x) - \nabla f_i(x)\|^2 \leq c_h\|\nabla f(x)\|^2 + \sigma_h^2$, $\forall x \in \mathbb{R}^d$, $\forall i$, *where* $c_h$ *is a positive constant.*

Assumptions 1-3 are commonly used in stochastic optimization [42]. Assumption 4, also known as the bounded gradient

dissimilarity assumption [10], is adopted to characterize the non-i.i.d. extent of the local data distribution. Similar assumptions have also been made in the literature [7], [21]–[25] for the convergence analysis under the non-i.i.d. setting. Note that these assumptions are only required for convergence analysis which are standard in zeroth-order optimization [10], [20].

In the following, we first present the convergence analysis for full device participation and then extend the analysis to partial device participation.

### A. Full Device Participation

We first characterize the convergence of the FedZO algorithm with full device participation in Theorem 1. We take the squared gradient $\|\nabla f(\boldsymbol{x}^t)\|^2$ to evaluate the suboptimality of the iterates. The speed of approaching a stationary point is an important metric to evaluate the algorithmic effectiveness for non-convex problems [43].

**Theorem 1.** *Suppose Assumptions 1-4 hold and the learning rate satisfies*

$$\eta \leq \min\left\{\frac{N}{72\tilde{c}_g\tilde{c}_h L}, \frac{2}{NH^2L}, \frac{1}{3\sqrt{\tilde{c}_g}HL}\right\}, \quad (7)$$

*the FedZO algorithm with full device participation satisfies*

$$\min_{t\in[T]} \mathbb{E}\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2 \leq 4\frac{f\left(\boldsymbol{x}^0\right) - f_*}{HT\eta} + \eta\frac{24L}{N}\tilde{\sigma}^2$$
$$+ \frac{dL^2\mu^2}{12} + 5L^2\mu^2, \quad (8)$$

*where $\tilde{\sigma}^2 = 3\left(1 + \frac{c_g d}{b_1 b_2}\right)\sigma_h^2 + \frac{d\sigma_g^2}{b_1 b_2}$, $\tilde{c}_g = 1 + \frac{c_g d}{b_1 b_2}$, and $\tilde{c}_h = 1 + c_h$.*

*Proof.* Please refer to Appendix A. □

According to Theorem 1, the upper bound of the minimum squared gradient among the global model sequence is composed of four terms. The first term shows that the optimality gap relies on the initial optimality. The second term shows that the optimality gap depends on the the non-i.i.d. extent of the local data distribution. The rest of the terms are related to step size $\mu$ for computing the gradient estimator that is unique in zeroth-order optimization. As pointed out in [43], we can select an appropriate step size to attain the desired accuracy. The following corollary follows by substituting a suitable learning rate $\eta$ and step size $\mu$ into Theorem 1.

**Corollary 1.** *Suppose Assumptions 1-4 hold and let $b_1 b_2 \leq d$, $\mu = (db_1 b_2 NHT)^{-\frac{1}{4}}$, and $\eta = (Nb_1 b_2)^{\frac{1}{2}}(dHT)^{-\frac{1}{2}}$, which holds for (7) if $T$ is large enough. The FedZO algorithm with full device participation satisfies*

$$\min_{t\in[T]} \mathbb{E}\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2 \leq \mathcal{O}\left(d^{\frac{1}{2}}(NHTb_1 b_2)^{-\frac{1}{2}}\right)$$
$$+ \mathcal{O}\left((db_1 b_2 NHT)^{-\frac{1}{2}}\right), \quad (9)$$

*where the right hand side of (9) is dominated by $\mathcal{O}\left(d^{\frac{1}{2}}(NHTb_1 b_2)^{-\frac{1}{2}}\right)$.*

We consider the case of $b_1 b_2 \leq d$ in Corollary 1 since the dimension $d$ is generally very large in many ML tasks.

TABLE I: Convergence rates of some typical algorithms for stochastic nonconvex unconstrained optimization.

| | Algorithm | Convergence rate | Maximum value of $H$ |
|---|---|---|---|
| FL setting | FedZO | $O\left(\sqrt{d/NHTb_1 b_2}\right)$ | $\min\left\{\mathcal{O}\left((dT)^{\frac{1}{3}}(b_1 b_2)^{-\frac{1}{3}}N^{-1}\right), \mathcal{O}\left(TN^{-1}\right)\right\}$ |
| | FedAvg [21] | $O\left(\sqrt{1/NHT}\right)$ | $O\left(T^{\frac{1}{3}}N^{-1}\right)$ |
| Distributed zeroth-order | ZONE-S [28] | $O(d^3/T)$ | — |
| | DZOPA [10] | $O\left(\sqrt{d/NT}\right)$ | — |
| Centralized zeroth-order | ZO-SGD [44] | $O\left(\sqrt{d/T}\right)$ | — |

Note: $T$ denotes the number of total communication rounds for FedZO and FedAvg, and the number of total iterations for others.

Besides, when $b_1 b_2 \leq d$, the computational consumption of (5) is lower than that of the Kiefer-Wolfowitz type scheme [33]. On the other hand, if $b_1 b_2 > d$, $\eta = N^{\frac{1}{2}}(HT)^{-\frac{1}{2}}$, and $\mu = d^{-\frac{1}{2}}(NHT)^{-\frac{1}{4}}$, according to Theorem 1, we obtain a convergence rate $\mathcal{O}\left((NHT)^{-\frac{1}{2}}\right)$ for the FedZO algorithm which is independent of dimension $d$. Such a convergence rate is the same as that of the FedAvg algorithm. The learning rate of the zeroth-order methods is generally $\sqrt{d}$−times smaller than that of their first-order counterparts [21], [43], [44], as the two-point gradient estimator is less accurate than the gradient. In our work, by adopting the mini-batch-type gradient estimator, we can increase the mini-batch sizes $b_1$ and $b_2$ to enhance the accuracy of the gradient estimator and also balance the effects of $d$ and $T$ on the learning rate.

**Remark 2.** *In Corollary 1, we set the learning rate $\eta = (Nb_1 b_2)^{\frac{1}{2}}(dHT)^{-\frac{1}{2}}$, which decreases as the number of local updates (i.e., $H$) increases. In particular, the progress of one local update shrinks by $1/\sqrt{H}$. However, by performing $H$ steps of local updates, we obtain a $\sqrt{H}$-times speedup per communication round. This accords with the derived convergence rate $\mathcal{O}\left(d^{\frac{1}{2}}(NHTb_1 b_2)^{-\frac{1}{2}}\right)$. According to Corollary 1, to reach an $\epsilon$-stationary solution, the FedZO algorithm takes $\mathcal{O}\left(d(Nb_1 b_2)^{-1}\epsilon^{-2}\right)$ iterations (i.e., $HT = \mathcal{O}\left(d(Nb_1 b_2)^{-1}\epsilon^{-2}\right)$) with learning rate $\eta = (Nb_1 b_2)^{\frac{1}{2}}(dHT)^{-\frac{1}{2}}$. If we increase $H$, then the learning rate (i.e., $\eta$) decreases, and we can attain a higher-accuracy solution with the same number of communication rounds (i.e., $T$). Besides, when the total iteration number (i.e., $HT$) and the learning rate (i.e., $\eta$) are fixed, we can increase the number of local iterations (i.e., $H$) and reduce the number of communication rounds (i.e., $T$) to reach the same accuracy, which enhances the communication efficiency. It is worth noting that the number of local iterations cannot be arbitrarily large. According to (7), to achieve the largest reduction in communication overhead while preserving convergence, the optimal value of $H$ is $\min\left\{\mathcal{O}\left((dT)^{\frac{1}{3}}(b_1 b_2)^{-\frac{1}{3}}N^{-1}\right), \mathcal{O}\left(TN^{-1}\right)\right\}$.*

**Remark 3.** *From Corollary 1, we notice that the proposed FedZO algorithm can attain convergence rate $\mathcal{O}\left(d^{\frac{1}{2}}(NHTb_1 b_2)^{-\frac{1}{2}}\right)$. In particular, FedZO achieves linear speedup in terms of the number of local iterates and the number of participating edge devices compared with the centralized zeroth-order algorithm (i.e., ZO-SGD) that achieves convergence rate $\mathcal{O}\left(d^{\frac{1}{2}}T^{-\frac{1}{2}}\right)$ [44]. For fair comparison, we consider $b_1 = b_2 = 1$. To attain the same accuracy,*

*compared to DZOPA with convergence rate $\mathcal{O}\left(d^{\frac{1}{2}}(NT)^{-\frac{1}{2}}\right)$ [10], the number of communication rounds required by the FedZO algorithm can be reduced by a factor of H. Besides, it is worth noting that the convergence rate of the FedZO algorithm depends on the dimension of the model parameter. In particular, the convergence speed of FedZO is $\sqrt{d}$ times slower than that of its first-order counterpart, i.e., FedAvg. Such a degeneration is the same as its centralized counterpart [44]. In addition, for FedAvg, the maximum value of H is $O\left(T^{\frac{1}{3}}N^{-1}\right)$ [21], which is smaller than that of FedZO mentioned in Remark 2. This is because the learning rate of FedZO is lower than that of FedAvg, which allows more local updates while preserving convergence. The detailed comparison between the proposed algorithm and the related algorithms is summarized in Table I.*

### B. Partial Device Participation

In this subsection, we show the convergence of the FedZO algorithm with partial device participation. By bounding the minimum squared gradient among the global model sequence, we characterize the convergence of the FedZO algorithm in the following theorem.

**Theorem 2.** *Suppose Assumptions 1-4 hold and the learning rate satisfies*

$$\eta \leq \min\left\{\frac{M}{192\tilde{c}_g\tilde{c}_hL}, \frac{M}{72c_hHL}, \frac{2}{MH^2L}, \frac{1}{3\sqrt{\tilde{c}_g}HL},\right.$$
$$\left.\frac{1}{3\sqrt{MH^3L}}\right\}, \tag{10}$$

*the FedZO algorithm with partial device participation satisfies*

$$\min_{t\in[T]}\mathbb{E}\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2 \leq 4\frac{f\left(\boldsymbol{x}^0\right)-f_*}{HT\eta} + \eta\frac{32L}{M}\tilde{\sigma}^2$$
$$+\eta\frac{36HL\sigma_h^2}{M}+\frac{dL^2\mu^2}{24}+13L^2\mu^2, \tag{11}$$

*where $\tilde{c}_g$, $\tilde{c}_h$, and $\tilde{\sigma}^2$ are defined in Theorem 1.*

*Proof.* Please refer to Appendix B. $\qquad\square$

By comparing (11) with (8), we notice that the third term in (11) does not appear in (8), which is induced by the randomness of device sampling, while full device participation eliminates this randomness, thereby reducing the optimality gap.

Similarly, the following corollary follows by substituting suitable learning rate $\eta$ and step size $\mu$ into Theorem 2.

**Corollary 2.** *Suppose Assumptions 1-4 hold and let $b_1b_2 \leq d$, $\mu = (db_1b_2MHT)^{-\frac{1}{4}}$, and $\eta = (Mb_1b_2)^{\frac{1}{2}}(dHT)^{-\frac{1}{2}}$, which holds for (10) if T is large enough. The FedZO algorithm with partial device participation satisfies*

$$\min_{t\in[T]}\mathbb{E}\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2 \leq \mathcal{O}\left(d^{\frac{1}{2}}(MHTb_1b_2)^{-\frac{1}{2}}\right)$$
$$+\mathcal{O}\left((b_1b_2H)^{\frac{1}{2}}(dMT)^{-\frac{1}{2}}\right)+\mathcal{O}\left((db_1b_2MHT)^{-\frac{1}{2}}\right). \tag{12}$$

According to (12), to attain a linear speedup in terms of the number of local iterates and participating edge devices, the number of local iterates cannot exceed $\mathcal{O}(d(b_1b_2)^{-1})$. Combining it with constraint (10), we can derive the largest value of H as $\min\left\{\mathcal{O}\left((dT)^{\frac{1}{3}}(b_1b_2)^{-\frac{1}{3}}M^{-1}\right), \mathcal{O}\left(TM^{-1}\right), \mathcal{O}\left(d(b_1b_2)^{-1}\right)\right\}$.

## IV. AirComp-Assisted FedZO Algorithm

In this section, we study the implementation of the proposed FedZO algorithm over wireless networks using AirComp, where the edge devices communicate with the central server via wireless fading channels.

In each communication round, both the downlink model dissemination phase and the uplink model uploading phase involve wireless transmissions. As the central server generally has a much greater transmit power than the edge devices, the downlink model dissemination is assumed to be error-free as in most of the existing studies [35]–[39] and we focus on the uplink model uploading.

### A. Over-the-Air Aggregation

For the FedZO algorithm, a key observation is that the central server is interested in receiving an average of local model updates of scheduled edge devices rather than each individual one. In particular, at the t-th round, the central server aims to acquire

$$\boldsymbol{\Delta}^t = \frac{1}{|\mathcal{M}_t|}\sum_{i\in\mathcal{M}_t}\boldsymbol{\Delta}_i^t, \tag{13}$$

where $|\mathcal{M}_t|$ denotes cardinality of set $\mathcal{M}_t$. With conventional OMA schemes, the central server in the t-th round first receives the local model update, e.g., $\boldsymbol{\Delta}_i^t$, from each edge device, and then takes an average to obtain the desired global model update, i.e., $\boldsymbol{\Delta}^t$. However, these schemes may not be spectrum-efficient as the number of required resource blocks or the communication latency linearly increases with the number of participating edge devices. AirComp, as a new non-orthogonal multiple access scheme for scalable transmission, allows all edge devices to concurrently transmit their local model updates and exploits the waveform superposition property to achieve spectrum-efficient model aggregation. The communication resource needed for model uploading using AirComp is independent of the number of participating edge devices. Hence, we adopt AirComp for the aggregation of local model updates in this paper.

Consider a wireless FL system where all edge devices and the server are equipped with a single antenna. Over wireless fading channels, the local model updates transmitted by edge devices suffer from detrimental channel distortion, which in turn degenerates the convergence performance of the AirComp-assisted FedZO algorithm. We thus set a threshold $h_{\min}$ and choose a subset of edge devices $\mathcal{M}_t = \{i \mid |h_i^t| \geq h_{\min}\}$ to participate in the training, where $h_i^t \in \mathbb{C}$ represents the channel coefficient between edge device $i$ and the central server in round $t$. We assume that $h_i^t$ are i.i.d. across different edge devices and communication rounds [38], [39]. Note that we can treat the adopted device scheduling strategy as uniform sampling analyzed in Section III-B.

With AirComp, the scheduled edge devices concurrently transmit their precoded model updates, e.g., $\alpha_i^t \boldsymbol{\Delta}_i^t$, to the central server, where $\alpha_i^t$ is the transmit scalar of edge device $i$ at the $t$-th round. Note that synchronization is required among distributed edge devices as in [35]–[39], which can be realized by sharing a reference-clock across the edge devices [45] or utilizing the timing advance technique commonly adopted in 4G long term evolution (LTE) and 5G new radio (NR) [46]. We assume that the model update vector $\boldsymbol{\Delta}_i^t$ of dimension $d$ can be transmitted within one transmission block while the channel coefficient is invariant during one transmission block [35]–[39]. Thus, the aggregated signal received at the central server can be expressed as

$$s^t = \sum_{i \in \mathcal{M}_t} h_i^t \alpha_i^t \boldsymbol{\Delta}_i^t + \boldsymbol{n}_t, \tag{14}$$

where $\boldsymbol{n}_t \sim \mathcal{CN}(0, \sigma_w^2 \boldsymbol{I}_d)$ represents the additive white Gaussian noise (AWGN) vector at the central server.

### B. Transceiver Design

The transmitted signal at each edge device is subject to an energy constraint during one communication round, i.e., $\|\alpha_i^t \boldsymbol{\Delta}_i^t\|^2 \le dP$, where $dP$ is the total energy of each edge device in one communication round. We assume that the channel state information (CSI) is available at both the central server and edge devices as in [35]–[39]. To meet the energy constraint of each edge device, we set the transmit scalar of device $i$ as

$$\alpha_i^t = \frac{h_{\min}}{h_i^t} \sqrt{\frac{dP}{\Delta_{\max}^t}}, \ \forall i, \tag{15}$$

where $\Delta_{\max}^t = \max_{i \in \mathcal{M}_t} \|\boldsymbol{\Delta}_i^t\|^2$. The received signal is thus given by

$$s^t = \sqrt{\frac{dP h_{\min}^2}{\Delta_{\max}^t}} \sum_{i \in \mathcal{M}_t} \boldsymbol{\Delta}_i^t + \boldsymbol{n}_t. \tag{16}$$

To recover the desired global model update $\boldsymbol{\Delta}^t$ in (13) from $s^t$ in (16), the central server scales $s^t$ with a receive scalar $\frac{1}{|\mathcal{M}_t|} \sqrt{\frac{\Delta_{\max}^t}{dP h_{\min}^2}}$, and obtains a noisy version of the global model update as follows

$$\boldsymbol{y}^t = \boldsymbol{\Delta}^t + \tilde{\boldsymbol{n}}_t, \tag{17}$$

where $\tilde{\boldsymbol{n}}_t \sim \mathcal{CN}\left(0, \frac{\sigma_w^2 \Delta_{\max}^t}{|\mathcal{M}_t|^2 dP h_{\min}^2} \boldsymbol{I}_d\right)$. As a result, the global model at the central server is updated as $\boldsymbol{x}^{t+1} = \boldsymbol{x}^t + \boldsymbol{\Delta}^t + \tilde{\boldsymbol{n}}_t$. To facilitate the transceiver design, each edge device needs to know the maximum of squared norm of local model updates among the participating edge devices, i.e., $\Delta_{\max}^t$, and the instantaneous channel coefficient between itself and the central server, i.e., $h_i^t$, which can be obtained via feedback from the central server. Before uplink model aggregation, the central server collects the squared norm of local model update $\|\boldsymbol{\Delta}_i^t\|^2$ from each edge device $i \in \mathcal{M}_t$, and then broadcasts $\Delta_{\max}^t$ to all edge devices. Besides, the central server estimates and feeds back the channel coefficients to these corresponding edge devices. It is worth noting that the communication overhead introduced by the exchange of these scalars is negligible when compared with the transmission of high-dimensional model parameters.

**Remark 4.** *Different from most existing studies [16], [35]–[37] that only focus on compensating for channel fading, the adopted transmitter design, i.e., (15), takes the scale of the model update into account. This ensures that the distortion between the obtained signal and the desired signal, i.e., the scaled receiver noise, is proportional to the maximum of the squared norm of local updates. This distortion diminishes when the local model converges. In other words, with such a transmitter design, the detrimental effect of the noise can be eliminated as the iteration proceeds for a convergent algorithm.*

### C. Convergence Analysis for AirComp-Assisted FedZO

In the following theorem, we characterize the convergence of the AirComp-assisted FedZO algorithm described in the previous two subsections.

**Theorem 3.** *Suppose Assumptions 1-4 hold and the learning rate satisfies*

$$\eta \le \min \left\{ \frac{\tilde{M}}{288 \tilde{c}_g \tilde{c}_h L}, \frac{\tilde{M}}{108 c_h HL}, \frac{3}{2NH^2 L}, \frac{1}{3\sqrt{\tilde{c}_g} HL}, \right.$$
$$\left. \frac{1}{2\sqrt{3}NH^3 L}, \frac{\sqrt{\tilde{M}\gamma}}{L\sqrt{2\tilde{c}_g NH^3}}, \frac{\tilde{M}^2 \gamma}{36 \tilde{c}_g \tilde{c}_h NHL} \right\}, \tag{18}$$

*where $\tilde{M} = \min\{|\mathcal{M}_t|, t \in [T]\}$. The AirComp-assisted FedZO algorithm satisfies*

$$\min_{t \in [T]} \mathbb{E}\|\nabla f(\boldsymbol{x}^t)\|^2 \le 4 \frac{f(\boldsymbol{x}^0) - f_*}{HT\eta} + \eta \frac{32L}{\tilde{M}} \hat{C} \tilde{\sigma}^2$$
$$+ \eta \frac{36HL\sigma_h^2}{\tilde{M}} + \hat{C} \frac{dL^2 \mu^2}{36} + \left(12 + \frac{\hat{C}}{9}\right) L^2 \mu^2, \tag{19}$$

*where $\hat{C} = 1 + \frac{NH}{8\tilde{M}\gamma}$ and $\gamma = \frac{Ph_{\min}^2}{\sigma_w^2}$. $\tilde{c}_g$, $\tilde{c}_h$, and $\tilde{\sigma}^2$ are defined in Theorem 1.*

As can be observed from Theorem 3, the convergence rate depends on $\gamma$, which is the minimum receive SNR. Theorem 3 almost reduces to Theorem 2 when $\gamma$ goes to infinity, i.e., noise-free case. Obviously, a smaller value of SNR leads to a slower convergence speed, which meets our intuition. In the following corollary, we show that a same-order convergence rate as the noise-free case presented in Section III-B can be achieved with appropriate receive SNR $\gamma$, learning rate $\eta$, and step size $\mu$.

**Corollary 3.** *Suppose Assumptions 1-4 hold, $8\tilde{M}\gamma \ge NH$, i.e., the communication quality is good enough, and let $b_1 b_2 \le d$, $\mu = (db_1 b_2 \tilde{M}HT)^{-\frac{1}{4}}$ and $\eta = (\tilde{M}b_1 b_2)^{\frac{1}{2}}(dHT)^{-\frac{1}{2}}$ that holds for (18), we have*

$$\min_{t \in [T]} \mathbb{E}\|\nabla f(\boldsymbol{x}^t)\|^2 \le \mathcal{O}\left(d^{\frac{1}{2}}(\tilde{M}HTb_1 b_2)^{-\frac{1}{2}}\right)$$
$$+ \mathcal{O}\left((b_1 b_2 H)^{\frac{1}{2}}(d\tilde{M}T)^{-\frac{1}{2}}\right) + \mathcal{O}\left((db_1 b_2 \tilde{M}HT)^{-\frac{1}{2}}\right). \tag{20}$$

**Remark 5.** *From Corollary 3, it can be observed that the upper bound of the minimum squared gradient among the global model sequence approaches to zero as $T$ goes to infinity, while that of the existing algorithms with AirComp is only shown to be bounded by a non-diminishing optimality gap [35]–[37]. Moreover, the transceiver design in [35]–[37] is transformed to an optimization problem aiming to minimize this gap, which is computationally expensive. In contrast, our transceiver design follows the principle of COTAF [38] and mitigates the detrimental impact of channel fading and receiver noise perturbation without the need of solving optimization problems. Note that the analysis in [38] concentrates on the first-order algorithm under the strongly convex setup and relies on the assumption that the second-order moment of the stochastic gradient is bounded by a constant, which is restrictive [34] and not required in this paper.*

## V. SIMULATION RESULTS

In this section, we present simulation results to evaluate the effectiveness of the proposed FedZO algorithm for applications of federated black-box attack and softmax regression.

### A. Federated Black-Box Attack

The robustness of machine learning (ML) models is an important performance metric for their practical application. For example, in an image classification model, the prediction results of the ML model are expected to be the same as the decision that humans make. In other words, the same output should be generated by a robust model if the input image is perturbed by a noise imperceptible to human. To evaluate the robustness of ML models, black-box attacks can be adopted, where the adversary acts as a standard user that does not have access to the inner structure of ML models and can only query the outputs (label or confidence score) for different inputs. This situation occurs when attacking ML cloud services where the model only serves as an API. Due to the black-box property, the optimization of black-box attacks falls into the category of zeroth-order optimization.

We consider federated black-box attacks [10] on the image classification DNN models that are well trained on some standard datasets. Federated black-box attacks aim to collaboratively generate a common perturbation such that the perturbed images are visually imperceptible to a human but could mislead the classifier. For image $z_i$, the attack loss [47] is given by

$$\psi_i(\boldsymbol{x}) = \max\left\{ \Phi_{y_i}\left(\frac{1}{2}\tanh\left(\tanh^{-1}2\boldsymbol{z}_i + \boldsymbol{x}\right)\right) \right. $$
$$- \max_{j \neq y_i}\left\{ \Phi_j\left(\frac{1}{2}\tanh\left(\tanh^{-1}2\boldsymbol{z}_i + \boldsymbol{x}\right)\right) \right\}, 0 \right\}$$
$$+ c\left\| \frac{1}{2}\tanh\left(\tanh^{-1}2\boldsymbol{z}_i + \boldsymbol{x}\right) - \boldsymbol{z}_i \right\|^2, \quad (21)$$

where $y_i$ denotes the label of image $\boldsymbol{z}_i$, $\Phi_j(\boldsymbol{z})$ represents the prediction confidence of image $\boldsymbol{z}$ to class $j$, $\frac{1}{2}\tanh\left(\tanh^{-1}2\boldsymbol{z}_i + \boldsymbol{x}\right)$ is the adversarial example of $\boldsymbol{z}_i$, $\frac{1}{2}\tanh\left(\tanh^{-1}2\boldsymbol{z}_i + \boldsymbol{x}\right) - \boldsymbol{z}_i$ is the distortion perturbed by $\boldsymbol{x}$

in the original image space. The first term of $\psi_i(\boldsymbol{x})$ measures the probability of failing to attack. The last term of $\psi_i(\boldsymbol{x})$ represents the distortion induced by $\boldsymbol{x}$ in the original image space. The goal of attack is to find a visually small perturbation to mislead the classifier $\Phi(\cdot)$ that can be realized by minimizing $\psi_i(\boldsymbol{x})$. Parameter $c$ balances the trade-off between the adversarial success and distortion loss. We denote the dataset at edge device $n$ as $\mathcal{D}_n$. The attack loss of device $n$ can be expressed as $f_n(\boldsymbol{x}) = \frac{1}{|\mathcal{D}_n|}\sum_{i \in \mathcal{D}_n}\psi_i(x)$. Federated black-box attacks of a DNN model can be formulated as: $\min_{\boldsymbol{x} \in \mathbb{R}^d}\frac{1}{N}\sum_{n=1}^{N}f_n(\boldsymbol{x})$, which can be tackled by the proposed FedZO algorithm.

In this experiment setting, all edge devices share one well-trained DNN classifier[1] that has a testing accuracy of $82.3\%$ on CIFAR-10 dataset [47]. We pick 4992 correctly classified samples from the training set of image class "deer" (containing 5,000 samples) and then distribute these samples to edge devices without overlapping. Each edge device is assigned a random number of samples.

We set the balancing parameter $c = 1$. The mini-batch sizes are set to $b_1 = 25$ and $b_2 = 20$. The learning rate and step size are set to $\eta = 0.001$ and $\mu = 0.001$, respectively.

In Fig. 1a, we show the impact of the number of local updates on the convergence performance of the proposed FedZO algorithm with full device participation. Specifically, we vary the number of local updates $H \in \{5, 10, 20, 50\}$ and present the attack loss versus the number of communication rounds. It can be observed that the FedZO algorithm can effectively reduce the attack loss for different values of $H$. Besides, as $H$ increases, the convergence speed of the FedZO algorithm tends to increase. This demonstrates the speedup in the number of the local iterates as shown in Section III. We further compare the performance of the proposed FedZO algorithm with DZOPA [10] and ZONE-S [28]. For DZOPA, the learning rate (i.e., $\eta$) and step size (i.e., $\mu$) are set to $0.005$ and $0.001$, respectively. For ZONE-S, the penalty parameter (i.e., $\rho$, defined in [28]) and step size (i.e., $\mu$) are set to $500$ and $0.001$, respectively. Note that DZOPA was proposed for the peer-to-peer architecture which cannot be directly applied to our considered server-client architecture. For comparison, we depict the performance of DZOPA under a fully-connected graph. For fairness, we also upgrade the two-point stochastic gradient estimator of [10] to a mini-batch-type one as in (2). Results show that the FedZO algorithm outperforms the baselines even when $H = 5$. With a larger number of local updates, the attack loss of the FedZO algorithm decreases much faster than that of the baselines.

Fig. 1b shows the convergence performance of the FedZO algorithm versus the number of participating edge devices when $H = 20$. The number of participating edge devices $M$ takes values from set $\{5, 10, 25, 50\}$. It is clear that our proposed algorithm works well in terms of reducing attack loss under the four different values of $M$. As can be observed, increasing the number of edge devices gives rise to a better convergence speed. We can observe the speedup in the number of participating devices, which matches well with our analysis in Section III.

---

[1] https://github.com/carlini/nn_robust_attacks

(a) Impact of number of local updates when $N = 10$ and $M = 10$.

(b) Impact of number of participating edge devices when $N = 50$ and $H = 20$.
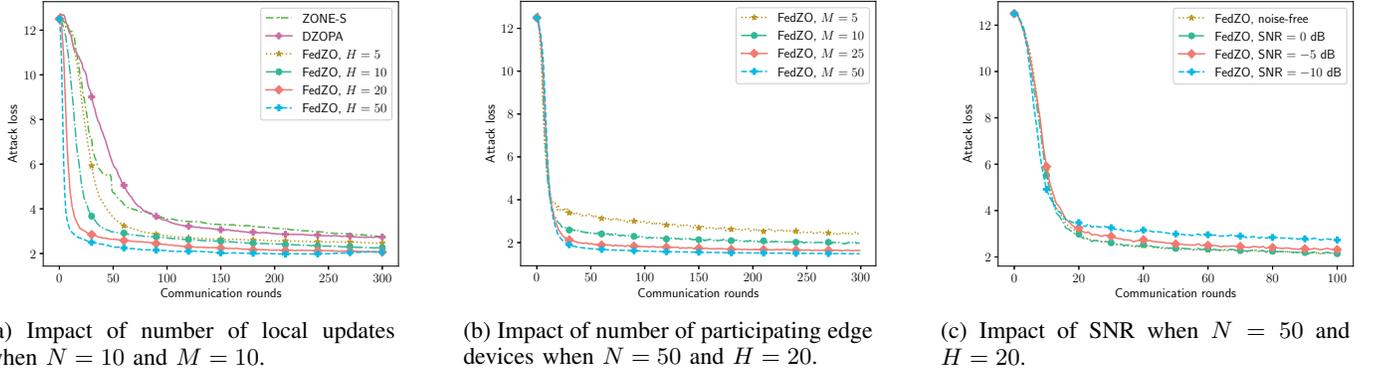
(c) Impact of SNR when $N = 50$ and $H = 20$.

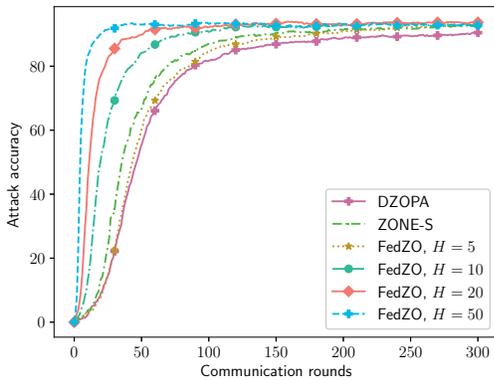Fig. 1: Attack loss of the federated black-box attack on CIFAR-10 dataset.



Fig. 2: Attack accuracy of the federated black-box attack on CIFAR-10 dataset.

In Fig. 1c, we show the performance of the AirComp-assisted FedZO algorithm presented in Section IV. Without loss of generality, we model the channel as $h_i^t \sim \mathcal{CN}(0,1)$, $\forall i, t$, and set the threshold $h_{\min} = 0.8$. Besides, we set the number of local iterates to $H = 20$. We take the FedZO algorithm with noise-free aggregation as benchmark with $H = 20$, where the participating edge devices are the same as that of the case with noise. We plot the attack loss versus the number of communication rounds under different SNR, i.e., $P/\sigma_w^2 \in \{-10 \text{ dB}, -5 \text{ dB}, 0 \text{ dB}\}$. As can be observed, the convergence of the FedZO algorithm can be preserved under such SNR settings. Besides, increasing the SNR accelerates the convergence speed of the FedZO algorithm. Especially, the FedZO algorithm with noise SNR $= 0$ dB attains a comparable performance with the noise-free case. These observations are in line with our theoretical analysis in Section IV-C.

Fig. 2 further demonstrates the superiority of the proposed FedZO algorithm over the baselines in terms of the attack accuracy. As can be observed, the proposed FedZO algorithm achieves a better attack accuracy than ZONE-S and DZOPA when $H$ is greater than 10. When $H = 5$, ZONE-S achieves a higher attack accuracy than FedZO at the cost of incurring a higher attack loss. This is because the perturbation generated by ZONE-S brings a large distortion.
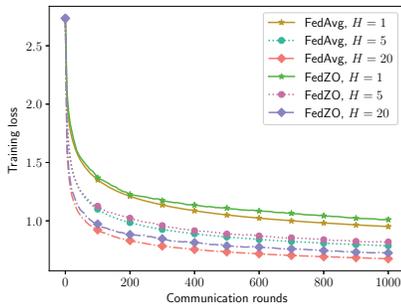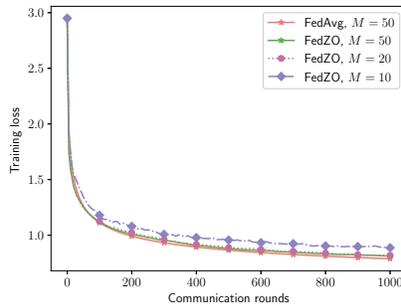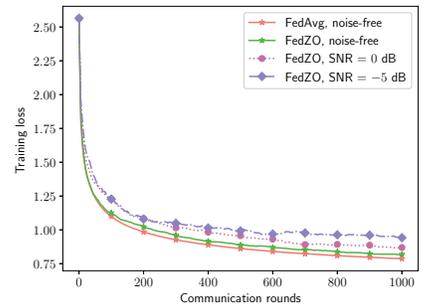
## B. Softmax Regression

We further validate our algorithm on the task of softmax regression, which corresponds to the multinomial classifier. We compare the FedZO algorithm with the FedAvg algorithm [2], which is the most representative first-order method. We set the learning rate (i.e., $\eta$), step size (i.e., $\mu$), and mini-batch sizes (i.e., $b_1$ and $b_2$) for the FedZO algorithm as 0.001, 0.001, 25, and 20, respectively. For the FedAvg algorithm, we set the learning rate (i.e., $\eta$) as 0.001.

We apply the softmax regression model to a 10-class-classification task on Fashion-MNIST [48] dataset. Throughout the experiment, we set the number of devices $N = 50$. Our strategy for constructing the non-i.i.d. data distribution follows the seminal work [2]. In particular, we sort the samples in the training set according to their labels, and then divide the training set into 100 shards of size 600. We then assign two shards to each device, such that each device owns a dataset of $1,200$ samples. Each edge device is assigned with at most four distinctive image labels.
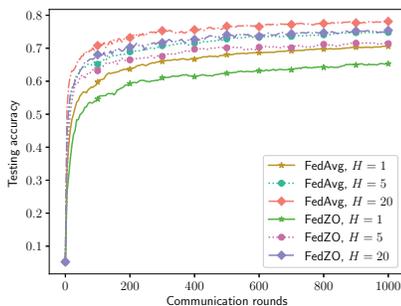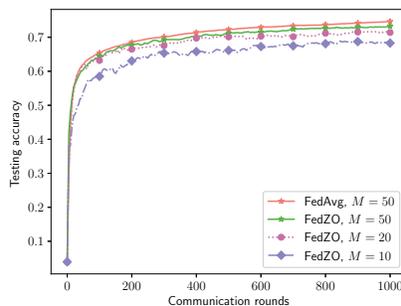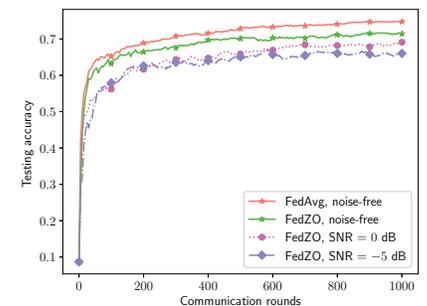
As shown in Fig. 3, the convergence speed of the FedZO algorithm is slightly slower than that of the FedAvg algorithm under the same number of local updates. This gap is brought by the uncertainty of the gradient estimator that FedZO utilizes. Further, we notice that the FedZO algorithm with $H = 20$ achieves comparable performance as the FedAvg algorithm with $H = 5$. As the FedZO algorithm only relies on the zeroth-order information, the slightly decreased performance is reasonable, and demonstrates the effectiveness of the FedZO algorithm. This also shows that the FedZO algorithm can serve as a satisfactory alternative for the FedAvg algorithm when the first-order information is not available.

In Fig. 4, we take the FedAvg algorithm as a benchmark when $H = 5$ and $M = 50$, and further investigate the impact of the number of participating edge devices, i.e., $M$, on the convergence behaviour of the FedZO algorithm. The phenomenon of speedup in $M$ can be witnessed in both the training loss and testing accuracy. We observe that the FedZO algorithm with $H = 5$ and $M = 50$ attains a comparable performance with the FedAvg algorithm.

Fig. 5 shows the performance of the AirComp-assisted FedZO algorithm over wireless networks with the same channel setting as mentioned in Section V-A. It can be observed

(a) Impact of $H$ on the training loss.



(a) Impact of $M$ on the training loss.



(a) Impact of SNR on the training loss.



(b) Impact of $H$ on the testing accuracy.

Fig. 3: The convergence results on the softmax regression problem with Fashion-MNIST dataset when $N = 50$ and $M = 20$.



(b) Impact of $M$ on the testing accuracy.

Fig. 4: The convergence results on the softmax regression problem with Fashion-MNIST dataset when $N = 50$ and $H = 5$.



(b) Impact of SNR on the testing accuracy

Fig. 5: The convergence results on the softmax regression problem with Fashion-MNIST dataset when $N = 50$ and $H = 5$.

that our proposed algorithm converges as the number of communication rounds increases and performs well when the SNR is not very small, e.g., SNR $\in \{-5$ dB, $0$ dB$\}$. Results also show that a greater SNR leads to a higher convergence speed. This result fits well with our analysis.

## VI. CONCLUSION

In this paper, we developed a derivative-free FedZO algorithm to handle federated optimization problems without using the gradient or Hessian information. Under non-convex settings, we characterized its convergence rate on non-i.i.d. data, and demonstrated the linear speedup in terms of the number of participating devices and local iterates. Subsequently, we established the convergence guarantee for the AirComp-assisted FedZO algorithm to support the implementation of the proposed algorithm over wireless networks. Simulation results demonstrated the effectiveness of the proposed FedZO algorithm and showed that the FedZO algorithm could serve as a satisfactory alternative for the FedAvg algorithm. It was also validated that the AirComp-assisted FedZO algorithm could attain a comparable performance with that of the noise-free case under certain SNR conditions.

## APPENDIX

To prove Theorems 1 and 2, we first characterize per round progress by Lemmas 1 and 3, respectively, and then bound the client drift during $H$ local iterates by Lemma 2. To prove

Theorem 3, we further bound the wireless noise by Lemma 4. The proofs of these lemmas are deferred to Appendix D. Before presenting the proofs, we first introduce some notations that are frequently used in this appendix.

Let $\mathcal{F}^{(t,k)}$ be a $\sigma$-field representing all the historical information of the FedZO algorithm up to the start of the $k$-th iteration of the $t$-th round. $\mathbb{E}_t$ and $\mathbb{E}_t^k$ denote expectations conditioning on $\mathcal{F}^{(t,0)}$ and $\mathcal{F}^{(t,k)}$, respectively. Let $\zeta_t = \left\{ \{\xi_{i,m}^{(t,k)}, \boldsymbol{v}_{i,n}^{(t,k)}\}_{i=1,2,...,N;k=0,1,...,H-1}^{m=1,2,...,b_1;n=1,2,...,b_2} \right\}$ and $\zeta_t^k = \left\{ \{\xi_{i,m}^{(t,\tau)}, \boldsymbol{v}_{i,n}^{(t,\tau)}\}_{i=1,2,...,N;\tau=0,1,...,k}^{m=1,2,...,b_1;n=1,2,...,b_2} \right\}$. $\mathbb{E}_{\zeta_t}$, $\mathbb{E}_{\zeta_t^k}$, and $\mathbb{E}_{\mathcal{M}_t}$ denote expectations over $\zeta_t$, $\zeta_t^k$, and $\mathcal{M}_t$, respectively.

### A. Proof of Theorem 1

For notational ease, we denote $\delta_t = \mathbb{E}_{\zeta_t}\left[ \frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1} \|\boldsymbol{x}_i^{(t,k)} - \boldsymbol{x}^t\|^2 \right]$. Before proving Theorem 1, we present the following two lemmas. The first lemma characterizes how the global loss, i.e., $f(\boldsymbol{x}^t)$, evolves as the iteration continues.

**Lemma 1.** *With Assumptions 1-4 and full device participation, by letting $\eta \le \frac{1}{2HL}$, we have*

$$\mathbb{E}_t\left[ f\left(\boldsymbol{x}^{t+1}\right) \right] \le f\left(\boldsymbol{x}^t\right) - \left( \frac{\eta H}{2} - \eta^2 \frac{6\tilde{c}_g \tilde{c}_h HL}{N} \right) \left\| \nabla f\left(\boldsymbol{x}^t\right) \right\|^2$$

$$+ \left( \eta L^2 + \eta^2 \frac{6\tilde{c}_g L^3}{N} \right) \delta_t + \eta^2 \frac{2HL}{N} \tilde{\sigma}^2$$

$$+ \eta^2 \frac{2HL^3\mu^2}{N} + \eta^2 \frac{d^2 HL^3\mu^2}{2Nb_1b_2} + \eta HL^2\mu^2. \tag{22}$$

*Please refer to Appendix D2 for the proof.*

Lemma 1 implies that we need to bound $\delta_t$, which is tackled by the following lemma.

**Lemma 2.** *With Assumptions 1-4 and $\eta \leq \frac{1}{3HL\sqrt{\tilde{c}_g}}$, we have*

$$\delta_t \leq 3\eta^2\tilde{c}_g\tilde{c}_h H^3 \left\|\nabla f(\boldsymbol{x}^t)\right\|^2 + H^3 L^2\eta^2\mu^2 + \eta^2 H^3\tilde{\sigma}^2 + \frac{d^2 H^3 L^2}{4b_1 b_2}\eta^2\mu^2.$$

*Please refer to Appendix D3 for the proof.*

By substituting the upper bound of $\delta_t$ in Lemma 2 into (22), we obtain

$$\mathbb{E}_t\left[f\left(\boldsymbol{x}^{t+1}\right)\right] \leq f\left(\boldsymbol{x}^t\right) - \left(\frac{\eta H}{2} - 3\eta^2\tilde{c}_g\tilde{c}_h Q_1(\eta)\right)\left\|\nabla f(\boldsymbol{x}^t)\right\|^2$$
$$+ \eta^2\tilde{\sigma}^2 Q_1(\eta) + Q_1(\eta)\left(L^2\eta^2\mu^2 + \frac{d^2 L^2}{4b_1 b_2}\eta^2\mu^2\right) + \eta HL^2\mu^2,$$

where $Q_1(\eta) = \frac{2HL}{N} + \left(\eta L^2 + \eta^2\frac{6\tilde{c}_g L^3}{N}\right)H^3$. Under condition (7), we have

$$Q_1(\eta) \leq \frac{6HL}{N}, \quad \eta^2\frac{18\tilde{c}_g\tilde{c}_h HL}{N} \leq \frac{\eta H}{4}, \quad L^2\eta^2\mu^2 \leq \frac{NL\eta\mu^2}{72},$$

and $\frac{d^2 L^2\eta^2\mu^2}{4b_1 b_2} \leq \frac{d^2 NL\eta\mu^2}{288b_1 b_2\tilde{c}_g\tilde{c}_h}$. Recalling the definition of $\tilde{c}_g$ and the fact that $c_g\tilde{c}_h \geq 1$, we have $\frac{d^2}{b_1 b_2\tilde{c}_g\tilde{c}_h} \leq \frac{d}{c_g\tilde{c}_h} \leq d$. With the above inequalities, we have

$$\mathbb{E}_t\left[f\left(\boldsymbol{x}^{t+1}\right)\right] \leq f\left(\boldsymbol{x}^t\right) - \frac{\eta H}{4}\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2 + \eta^2\frac{6HL}{N}\tilde{\sigma}^2$$
$$+ \frac{\eta dHL^2\mu^2}{48} + \frac{13}{12}\eta HL^2\mu^2. \quad (23)$$

By taking expectation on both sides of (23) and telescoping from $t = 0$ to $T - 1$, we obtain

$$\frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2 \leq 4\frac{f\left(\boldsymbol{x}^0\right) - \mathbb{E}\left[f\left(\boldsymbol{x}^T\right)\right]}{HT\eta} + \eta\frac{24L}{N}\tilde{\sigma}^2$$
$$+ \frac{dL^2\mu^2}{12} + 5L^2\mu^2. \quad (24)$$

By Assumption 1, i.e., $f(\boldsymbol{x}) \geq f_*$, we obtain Theorem 1.

### B. Proof of Theorem 2

We first characterize how the global loss, i.e., $f(\boldsymbol{x}^t)$, evolves as the iteration continues in the following lemma.

**Lemma 3.** *With Assumptions 1-4 and partial device participation, by letting the learning rate $\eta \leq \frac{1}{2HL}$, we have*

$$\mathbb{E}_t\left[f\left(\boldsymbol{x}^{t+1}\right)\right]$$
$$\leq f\left(\boldsymbol{x}^t\right) - \left(\frac{\eta H}{2} - \eta^2\frac{6\tilde{c}_g\tilde{c}_h HL}{M} - \eta^2\frac{9H^2 Lc_h}{M}\right)\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2$$
$$+ \left(\eta L^2 + \eta^2\frac{6\tilde{c}_g L^3}{M} + \eta^2 18HL^3\right)\delta_t + \eta^2\frac{2HL}{M}\tilde{\sigma}^2 + \frac{9\eta^2 H^2 L\sigma_h^2}{M}$$
$$+ \eta^2\frac{2HL^3\mu^2}{M} + \eta^2\frac{d^2 HL^3\mu^2}{2Mb_1 b_2} + 6\eta^2 H^2 L^3\mu^2 + \eta HL^2\mu^2.$$

*Please refer to Appendix D1 for the proof.*

By combining Lemmas 2 and 3, we have

$$\mathbb{E}_t\left[f\left(\boldsymbol{x}^{t+1}\right)\right]$$
$$\leq f\left(\boldsymbol{x}^t\right) - \left(\frac{\eta H}{2} - \eta^2\frac{9H^2 Lc_h}{M} - 3\tilde{c}_g\tilde{c}_h\eta^2 Q_2(\eta)\right)\left\|\nabla f(\boldsymbol{x}^t)\right\|^2$$
$$+ \eta^2\tilde{\sigma}^2 Q_2(\eta) + \frac{9\eta^2 H^2 L\sigma_h^2}{M} + Q_2(\eta)\left(L^2\eta^2\mu^2 + \frac{d^2 L^2}{4b_1 b_2}\eta^2\mu^2\right)$$
$$+ 6\eta^2 H^2 L^3\mu^2 + \eta HL^2\mu^2,$$

where $Q_2(\eta) = \frac{2HL}{M} + \left(\eta L^2 + \eta^2\frac{6\tilde{c}_g L^3}{M} + \eta^2 18HL^3\right)H^3$. Under condition (10), we have

$$\eta^2\frac{9H^2 Lc_h}{M} \leq \frac{\eta H}{8}, \quad Q_2(\eta) \leq \frac{8HL}{M}, \quad \frac{24\eta^2\tilde{c}_g\tilde{c}_h HL}{M} \leq \frac{\eta H}{8},$$

$$L^2\eta^2\mu^2 \leq \frac{ML\eta\mu^2}{192}, \quad \frac{d^2 L^2\eta^2\mu^2}{4b_1 b_2} \leq \frac{d^2 ML\eta\mu^2}{768b_1 b_2\tilde{c}_g\tilde{c}_h},$$

and $6\eta^2 H^2 L^3\mu^2 \leq 2\eta HL^2\mu^2$. Recalling the definition of $\tilde{c}_g$ and the fact that $c_g\tilde{c}_h \geq 1$, we have $\frac{d^2}{b_1 b_2\tilde{c}_g\tilde{c}_h} \leq \frac{d}{c_g\tilde{c}_h} \leq d$. With the above inequalities, we have

$$\mathbb{E}_t\left[f\left(\boldsymbol{x}^{t+1}\right)\right] \leq f\left(\boldsymbol{x}^t\right) - \frac{\eta H}{4}\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2 + \eta^2\frac{8HL}{M}\tilde{\sigma}^2$$
$$+ \frac{9\eta^2 H^2 L\sigma_h^2}{M} + \frac{\eta HdL^2\mu^2}{96} + \left(3 + \frac{1}{24}\right)\eta HL^2\mu^2.$$

By following the similar steps as in Appendix A, we obtain Theorem 2.

### C. Proof of Theorem 3

By denoting $\tilde{\boldsymbol{x}}^{t+1} = \boldsymbol{x}^t + \frac{1}{|\mathcal{M}_t|}\sum_{i\in\mathcal{M}_t}\Delta_i^t$ as the noise-free aggregated model, we have $\boldsymbol{x}^{t+1} = \tilde{\boldsymbol{x}}^{t+1} + \tilde{\boldsymbol{n}}_t$. By applying the smoothness of $f(\boldsymbol{x})$ in Assumption 2, we obtain

$$f\left(\boldsymbol{x}^{t+1}\right) \leq f\left(\tilde{\boldsymbol{x}}^{t+1}\right) + \left\langle\nabla f\left(\tilde{\boldsymbol{x}}^{t+1}\right), \tilde{\boldsymbol{n}}_t\right\rangle + \frac{L}{2}\left\|\tilde{\boldsymbol{n}}_t\right\|^2. \quad (25)$$

We denote $s^{(t,H)} = \frac{1}{N}\sum_{i=1}^N\mathbb{E}_{\zeta_t}\|\boldsymbol{x}_i^{(t,H)} - \boldsymbol{x}^t\|^2$. By taking an expectation for (25) conditioning on $\mathcal{F}^{(t,0)}$ and utilizing $\mathbb{E}_{\zeta_t}\left[\max_i\|\boldsymbol{x}_i^{(t,H)} - \boldsymbol{x}^t\|^2\right] \leq Ns^{(t,H)}$, we have

$$\mathbb{E}_t\left[f\left(\boldsymbol{x}^{t+1}\right)\right] \leq \mathbb{E}_t\left[f\left(\tilde{\boldsymbol{x}}^{t+1}\right)\right] + \frac{L}{2\gamma}\frac{N}{|\mathcal{M}_t|^2}s^{(t,H)}. \quad (26)$$

The above result suggests that we need to bound $s^{(t,H)}$, which can be handled by the following lemma.

**Lemma 4.** *With Assumptions 1-4 hold, we have*

$$s^{(t,H)} \leq 6\tilde{c}_g HL^2\eta^2\delta_t + 6\tilde{c}_g\tilde{c}_h H^2\eta^2\left\|\nabla f(\boldsymbol{x}^t)\right\|^2 + 2H^2\eta^2\tilde{\sigma}^2$$
$$+ 2H^2 L^2\eta^2\mu^2 + \frac{d^2 H^2 L^2}{2b_1 b_2}\eta^2\mu^2. \quad (27)$$

*Please refer to Appendix D4 for the proof.*

By combining Lemmas 2, 3, and 4, we obtain

$$\mathbb{E}_t\left[f\left(\boldsymbol{x}^{t+1}\right)\right]$$
$$\leq f\left(\boldsymbol{x}^t\right) - \left(\frac{\eta H}{2} - \eta^2\frac{9H^2 Lc_h}{|\mathcal{M}_t|} - 3\tilde{c}_g\tilde{c}_h\eta^2\tilde{Q}_3(\eta)\right)\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2$$
$$+ \eta^2\tilde{\sigma}^2\tilde{Q}_3(\eta) + \frac{9\eta^2 H^2 L\sigma_h^2}{|\mathcal{M}_t|} + \tilde{Q}_3(\eta)\left(L^2\eta^2\mu^2 + \frac{d^2 L^2}{4b_1 b_2}\eta^2\mu^2\right)$$
$$+ 6\eta^2 H^2 L^3\mu^2 + \eta HL^2\mu^2,$$

where $\tilde{Q}_3(\eta) = \frac{NH^2 L}{|\mathcal{M}_t|^2\gamma} + Q_3(\eta)$ and

$$Q_3(\eta) = \frac{2HL}{|\mathcal{M}_t|} + \left(\eta L^2 + \eta^2\frac{6\tilde{c}_g L^3}{|\mathcal{M}_t|} + \eta^2 18HL^3 + \frac{3\tilde{c}_g NHL^3\eta^2}{|\mathcal{M}_t|^2\gamma}\right)H^3.$$

Under condition (18), we have

$$\eta^2\frac{9H^2 Lc_h}{|\mathcal{M}_t|} \leq \frac{\eta H}{12}, \quad \frac{3\tilde{c}_g\tilde{c}_h NH^2 L\eta^2}{|\mathcal{M}_t|^2\gamma} \leq \frac{\eta H}{12}, \quad Q_3(\eta) \leq \frac{8HL}{|\mathcal{M}_t|},$$

$$\eta^2 \frac{24\tilde{c}_g\tilde{c}_h HL}{|\mathcal{M}_t|} \leq \frac{\eta H}{12}, \ 6\eta^2 H^2 L^3 \mu^2 \leq 2\eta HL^2\mu^2,$$

$$L^2\eta^2\mu^2 \leq \frac{|\mathcal{M}_t|L\eta\mu^2}{288}, \text{and } \frac{d^2L^2\eta^2\mu^2}{4b_1b_2} \leq \frac{d^2|\mathcal{M}_t|L\eta\mu^2}{1152b_1b_2\tilde{c}_g\tilde{c}_h}.$$

Recalling the definition of $\tilde{c}_g$ and the fact that $c_g\tilde{c}_h \geq 1$, we have $\frac{d^2}{b_1b_2\tilde{c}_g\tilde{c}_h} \leq \frac{d}{c_g\tilde{c}_h} \leq d$. With the above inequalities, we have

$$\mathbb{E}_t\left[(\boldsymbol{x}^{t+1})\right] \leq f\left(\boldsymbol{x}^t\right) - \frac{\eta H}{4}\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2 + \eta^2\frac{8HL}{|\mathcal{M}_t|}\hat{C}\tilde{\sigma}^2$$

$$+ \frac{9\eta^2H^2L\sigma_h^2}{|\mathcal{M}_t|} + \hat{C}\frac{\eta dHL^2\mu^2}{144} + \left(3 + \frac{\hat{C}}{36}\right)\eta HL^2\mu^2, \quad (28)$$

where $\hat{C} = 1 + \frac{NH}{8M\gamma}$ and $\tilde{M} \leq |\mathcal{M}_t|, \ \forall t$. After reorganizing (28), we obtain

$$\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2 \leq 4\frac{f\left(\boldsymbol{x}^t\right) - \mathbb{E}_t\left[f\left(\boldsymbol{x}^{t+1}\right)\right]}{\eta H} + \eta\frac{32L}{\tilde{M}}\hat{C}\tilde{\sigma}^2$$

$$+ \frac{36\eta HL\sigma_h^2}{\tilde{M}} + \hat{C}\frac{dL^2\mu^2}{36} + \left(12 + \frac{\hat{C}}{9}\right)L^2\mu^2. \quad (29)$$

By following the similar derivation as in Appendix A, we obtain Theorem 3.

### D. Proof of Lemmas

As Lemma 1 is a simplified version of Lemma 3, we first prove Lemma 3 and then prove Lemma 1.

*1) Proof of Lemma 3:* Based on Assumption 2 and $\eta\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\boldsymbol{e}_i^{(t,k)} = \boldsymbol{x}^{t+1} - \boldsymbol{x}^t$, we have

$$f\left(\boldsymbol{x}^{t+1}\right) \leq f\left(\boldsymbol{x}^t\right) - \eta\left\langle\nabla f\left(\boldsymbol{x}^t\right), \frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\boldsymbol{e}_i^{(t,k)}\right\rangle$$

$$+ \eta^2\frac{L}{2}\left\|\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\boldsymbol{e}_i^{(t,k)}\right\|^2. \quad (30)$$

By taking an expectation for (30) conditioning on $\mathcal{F}^{(t,0)}$, we obtain

$$\mathbb{E}_t\left[f\left(\boldsymbol{x}^{t+1}\right)\right] \leq f\left(\boldsymbol{x}^t\right) \underbrace{-\eta\mathbb{E}_{\zeta_t}\left[\left\langle\nabla f\left(\boldsymbol{x}^t\right), \frac{1}{M}\mathbb{E}_{\mathcal{M}_t}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\boldsymbol{e}_i^{(t,k)}\right\rangle\right]}_{T_1}$$

$$+ \eta^2\frac{L}{2}\underbrace{\mathbb{E}_t\left\|\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\boldsymbol{e}_i^{(t,k)}\right\|^2}_{T_2}. \quad (31)$$

As $\mathcal{M}_t$ is uniformly sampled from $N$ edge devices, by utilizing [22, Lemma 4], we have

$$T_1 = -\eta\mathbb{E}_{\zeta_t}\left[\left\langle\nabla f\left(\boldsymbol{x}^t\right), \frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\boldsymbol{e}_i^{(t,k)}\right\rangle\right].$$

According to (4), we have

$$\mathbb{E}_{\zeta_t}\left[\frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\left(\boldsymbol{e}_i^{(t,k)} - \nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right)\right)\right] = 0. \quad (32)$$

With the above two equalities, we have

$$T_1 = -\eta\mathbb{E}_{\zeta_t}\left[\left\langle\nabla f\left(\boldsymbol{x}^t\right), \frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right)\right\rangle\right].$$

Because of the equality $2\langle\boldsymbol{a}, \boldsymbol{b}\rangle = \|\boldsymbol{a}\|^2 + \|\boldsymbol{b}\|^2 - \|\boldsymbol{a}-\boldsymbol{b}\|^2$, we obtain

$$T_1 = -\frac{\eta H}{2}\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2 - \frac{\eta H}{2}\mathbb{E}_{\zeta_t}\left\|\frac{1}{NH}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right)\right\|^2$$

$$+ \frac{\eta H}{2}\mathbb{E}_{\zeta_t}\underbrace{\left\|\frac{1}{NH}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\left(\nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right) - \nabla f_i\left(\boldsymbol{x}^t\right)\right)\right\|^2}_{T_3}. \quad (33)$$

For $T_3$, we have

$$T_3 \leq \frac{1}{NH}\mathbb{E}_{\zeta_t}\left[\sum_{i=1}^{N}\sum_{k=0}^{H-1}\left\|\nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right) - \nabla f_i\left(\boldsymbol{x}^t\right)\right\|^2\right]$$

$$= \frac{1}{NH}\mathbb{E}_{\zeta_t}\left[\sum_{i=1}^{N}\sum_{k=0}^{H-1}\left\|\nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right) \mp \nabla f_i\left(\boldsymbol{x}_i^{(t,k)}\right) - \nabla f_i\left(\boldsymbol{x}^t\right)\right\|^2\right]$$

$$\leq \frac{2}{NH}\mathbb{E}_{\zeta_t}\left[\sum_{i=1}^{N}\sum_{k=0}^{H-1}\left\|\nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right) - \nabla f_i\left(\boldsymbol{x}_i^{(t,k)}\right)\right\|^2\right]$$

$$+ \frac{2}{NH}\mathbb{E}_{\zeta_t}\left[\sum_{i=1}^{N}\sum_{k=0}^{H-1}\left\|\nabla f_i\left(\boldsymbol{x}_i^{(t,k)}\right) - \nabla f_i\left(\boldsymbol{x}^t\right)\right\|^2\right]$$

$$\leq 2L^2\mu^2 + \frac{2L^2}{NH}\mathbb{E}_{\zeta_t}\left[\sum_{i=1}^{N}\sum_{k=0}^{H-1}\left\|\boldsymbol{x}_i^{(t,k)} - \boldsymbol{x}^t\right\|^2\right], \quad (34)$$

where the first inequality follows by the Jensen's inequality, $a \mp b$ represents $a - b + b$, the second inequality holds because of the Cauchy-Schwartz inequality, and the last inequality follows by [10, Lemma 2] and the smoothness of $f_i(\boldsymbol{x})$ in Assumption 2. By substituting (34) into (33), we have

$$T_1 \leq -\frac{\eta H}{2}\left\|\nabla f\left(\boldsymbol{x}^t\right)\right\|^2 - \frac{\eta H}{2}\mathbb{E}_{\zeta_t}\left\|\frac{1}{NH}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})\right\|^2$$

$$+ \eta HL^2\mu^2 + \eta L^2\mathbb{E}_{\zeta_t}\left[\frac{1}{N}\sum_{i=1}^{N}\sum_{k}^{H}\|\boldsymbol{x}_i^{(t,k)} - \boldsymbol{x}^t\|^2\right]. \quad (35)$$

For $T_2$, according to the Cauchy-Schwartz inequality, we have

$$T_2 \leq 2\mathbb{E}_t\underbrace{\left\|\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\left(\boldsymbol{e}_i^{(t,k)} - \nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right)\right)\right\|^2}_{T_4}$$

$$+ 2\mathbb{E}_t\underbrace{\left\|\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right)\right\|^2}_{T_5}. \quad (36)$$

By denoting $\boldsymbol{h}_i = \sum_{k=0}^{H-1}(\boldsymbol{e}_i^{(t,k)} - \nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)}))$ and utilizing (4), we have $\mathbb{E}_{\zeta_t}[\boldsymbol{h}_i] = 0$. Due to the independence between $\boldsymbol{h}_i$ and $\boldsymbol{h}_j, \forall j \neq i$, we have $\mathbb{E}_{\zeta_t}[\langle\boldsymbol{h}_i, \boldsymbol{h}_j\rangle] = 0$. We thus obtain

$$T_4 = \frac{1}{M^2}\mathbb{E}_{\mathcal{M}_t}\left[\sum_{i\in\mathcal{M}_t}\mathbb{E}_{\zeta_t}\left\|\sum_{k=0}^{H-1}\left(\boldsymbol{e}_i^{(t,k)} - \nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right)\right)\right\|^2\right]$$

$$= \frac{1}{MN}\sum_{i=1}^{N}\mathbb{E}_{\zeta_t}\left\|\sum_{k=0}^{H-1}\left(\boldsymbol{e}_i^{(t,k)} - \nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right)\right)\right\|^2. \quad (37)$$

According to (4) and [7, Lemma 2], it follows that

$$T_4 = \frac{1}{MN}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\mathbb{E}_{\zeta_t^k}\left\|\boldsymbol{e}_i^{(t,k)} - \nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right)\right\|^2. \quad (38)$$

As $\mathbb{E}\|\boldsymbol{z}-\mathbb{E}[\boldsymbol{z}]\|^2 \leq \mathbb{E}\|\boldsymbol{z}\|^2$, we have

$$T_4 \leq \frac{1}{MN}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\mathbb{E}_{\zeta_t^k}\left\|\boldsymbol{e}_i^{(t,k)}\right\|^2. \tag{39}$$

Recalling (2) and (5), we have $\boldsymbol{e}_i^{(t,k)} = \frac{1}{b_1 b_2}\sum_{m=1}^{b_1}\sum_{n=1}^{b_2}\boldsymbol{e}_{i,m,n}^{(t,k)}$ by denoting

$$\boldsymbol{e}_{i,m,n}^{(t,k)} = \frac{d\boldsymbol{v}_{i,n}^{(t,k)}}{\mu}\left(F_i(\boldsymbol{x}_i^{(t,k)}+\mu\boldsymbol{v}_{i,n}^{(t,k)},\xi_{i,m}^{(t,k)})-F_i(\boldsymbol{x}_i^{(t,k)},\xi_{i,m}^{(t,k)})\right).$$

According to (3), we have

$$\mathbb{E}_t^k\left[\boldsymbol{e}_{i,m,n}^{(t,k)}\right] = \nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right),\ \forall m,\ n. \tag{40}$$

Therefore, we can bound $\mathbb{E}_{\zeta_t^k}\left\|\boldsymbol{e}_i^{(t,k)}\right\|^2$ as follows

$$\mathbb{E}_{\zeta_t^k}\left\|\boldsymbol{e}_i^{(t,k)}\right\|^2 = \mathbb{E}_{\zeta_t^{k-1}}\left[\mathbb{E}_t^k\left\|\boldsymbol{e}_i^{(t,k)}\right\|^2\right]$$

$$=\mathbb{E}_{\zeta_t^{k-1}}\left[\mathbb{E}_t^k\left\|\frac{1}{b_1 b_2}\sum_{m=1}^{b_1}\sum_{n=1}^{b_2}\boldsymbol{e}_{i,m,n}^{(t,k)}-\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})\right\|^2+\left\|\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})\right\|^2\right]$$

$$=\frac{1}{b_1 b_2}\mathbb{E}_{\zeta_t^k}\left\|\boldsymbol{e}_{i,1,1}^{(t,k)}-\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})\right\|^2+\mathbb{E}_{\zeta_t^{k-1}}\left\|\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})\right\|^2$$

$$\leq\frac{1}{b_1 b_2}\mathbb{E}_{\zeta_t^k}\left\|\boldsymbol{e}_{i,1,1}^{(t,k)}\right\|^2+\mathbb{E}_{\zeta_t^{k-1}}\left\|\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})\right\|^2, \tag{41}$$

where the second equality follows by $\mathbb{E}\|\boldsymbol{z}\|^2 = \|\mathbb{E}[\boldsymbol{z}]\|^2+\mathbb{E}\|\boldsymbol{z}-\mathbb{E}[\boldsymbol{z}]\|^2$, the third equality follows by the fact that $\boldsymbol{e}_{i,m,n}^{(t,k)}$ and $\boldsymbol{e}_{i,m',n'}^{(t,k)}$ are independent when $m\neq m'$ or $n\neq n'$, and the inequality holds because of $\mathbb{E}\|\boldsymbol{z}-\mathbb{E}[\boldsymbol{z}]\|^2 \leq \mathbb{E}\|\boldsymbol{z}\|^2$. We further bound $\mathbb{E}_{\zeta_t^k}\left\|\boldsymbol{e}_{i,1,1}^{(t,k)}\right\|^2$ as follows

$$\mathbb{E}_{\zeta_t^k}\left\|\boldsymbol{e}_{i,1,1}^{(t,k)}\right\|^2 \leq 2d\mathbb{E}_{\zeta_t^k}\left\|\nabla F_i(\boldsymbol{x}_i^{(t,k)},\xi_{i,1}^{(t,k)})\right\|^2+\frac{1}{2}d^2 L^2\mu^2$$

$$\leq 2c_g d\mathbb{E}_{\zeta_t^{k-1}}\left\|\nabla f_i(\boldsymbol{x}_i^{(t,k)})\right\|^2+2d\sigma_g^2+\frac{1}{2}d^2 L^2\mu^2, \tag{42}$$

where the first inequality follows by [20, Lemma 4.1] and the second inequality follows by Assumption 3. Besides,

$$\mathbb{E}_{\zeta_t^{k-1}}\left\|\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})\right\|^2$$

$$=\mathbb{E}_{\zeta_t^{k-1}}\left\|\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})-\nabla f_i(\boldsymbol{x}_i^{(t,k)})+\nabla f_i(\boldsymbol{x}_i^{(t,k)})\right\|^2$$

$$\leq 2\mathbb{E}_{\zeta_t^{k-1}}\left\|\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})-\nabla f_i(\boldsymbol{x}_i^{(t,k)})\right\|^2+2\mathbb{E}_{\zeta_t^{k-1}}\left\|\nabla f_i(\boldsymbol{x}_i^{(t,k)})\right\|^2$$

$$\leq 2\mu^2 L^2+2\mathbb{E}_{\zeta_t^{k-1}}\left\|\nabla f_i(\boldsymbol{x}_i^{(t,k)})\right\|^2, \tag{43}$$

where the last inequality follows by [10, Lemma 2]. By combining (41), (42), and (43), we have

$$\mathbb{E}_{\zeta_t^k}\left\|\boldsymbol{e}_i^{(t,k)}\right\|^2 \leq \left(2+\frac{2c_g d}{b_1 b_2}\right)\mathbb{E}_{\zeta_t^{k-1}}\left\|\nabla f_i(\boldsymbol{x}_i^{(t,k)})\right\|^2$$

$$+\frac{2d\sigma_g^2}{b_1 b_2}+\frac{d^2 L^2\mu^2}{2b_1 b_2}+2L^2\mu^2. \tag{44}$$

We next bound $\mathbb{E}_{\zeta_t^{k-1}}\left\|\nabla f_i(\boldsymbol{x}_i^{(t,k)})\right\|^2$ as below

$$\mathbb{E}_{\zeta_t^{k-1}}\left\|\nabla f_i(\boldsymbol{x}_i^{(t,k)})\right\|^2$$

$$=\mathbb{E}_{\zeta_t^{k-1}}\left\|\nabla f_i(\boldsymbol{x}_i^{(t,k)})\mp\nabla f_i(\boldsymbol{x}^t)\mp\nabla f(\boldsymbol{x}^t)\right\|^2$$

$$\leq 3L^2\mathbb{E}_{\zeta_t^{k-1}}\left\|\boldsymbol{x}_i^{(t,k)}-\boldsymbol{x}^t\right\|^2+3\sigma_h^2+(3c_h+3)\left\|\nabla f(\boldsymbol{x}^t)\right\|^2, \tag{45}$$

where the inequality follows by the Cauchy-Schwartz inequality, Assumption 2, and Assumption 4. According to (39), (44), and (45), we obtain

$$T_4 \leq \frac{6\tilde{c}_g L^2}{M}\mathbb{E}_{\zeta_t}\left[\frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\left\|\boldsymbol{x}_i^{(t,k)}-\boldsymbol{x}^t\right\|^2\right]+\frac{6\tilde{c}_g\tilde{c}_h H}{M}\left\|\nabla f(\boldsymbol{x}^t)\right\|^2$$

$$+\frac{2H}{M}\tilde{\sigma}^2+\frac{2HL^2\mu^2}{M}+\frac{d^2 HL^2\mu^2}{2Mb_1 b_2}, \tag{46}$$

where $\tilde{\sigma}^2$, $\tilde{c}_g$, and $\tilde{c}_h$ are defined in Theorem 1.

Next, we split $T_5$ as follows

$$T_5 = \mathbb{E}_{\zeta_t}\left[\mathbb{E}_{\mathcal{M}_t}\left\|\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right)\right\|^2\right] =$$

$$\mathbb{E}_{\zeta_t}\left[\mathbb{E}_{\mathcal{M}_t}\underbrace{\left\|\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right)-\frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right)\right\|^2}_{T_6}\right.$$

$$\left.+\left\|\frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\nabla f_i^\mu\left(\boldsymbol{x}_i^{(t,k)}\right)\right\|^2\right], \tag{47}$$

where the equality follows by $\mathbb{E}\|\boldsymbol{z}\|^2 = \|\mathbb{E}[\boldsymbol{z}]\|^2+\mathbb{E}\|\boldsymbol{z}-\mathbb{E}[\boldsymbol{z}]\|^2$.

For $T_6$, we provide the following upper bounds

$$T_6 = \mathbb{E}_t\left\|\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})\mp\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\nabla f_i(\boldsymbol{x}_i^{(t,k)})\right.$$

$$\left.\mp\frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\nabla f_i(\boldsymbol{x}_i^{(t,k)})-\frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})\right\|^2$$

$$\leq 3\mathbb{E}_t\left\|\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\left(\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})-\nabla f_i(\boldsymbol{x}_i^{(t,k)})\right)\right\|^2$$

$$+3\mathbb{E}_t\left\|\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\nabla f_i(\boldsymbol{x}_i^{(t,k)})-\frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\nabla f_i(\boldsymbol{x}_i^{(t,k)})\right\|^2$$

$$+3\mathbb{E}_t\left\|\frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\left(\nabla f_i(\boldsymbol{x}_i^{(t,k)})-\nabla f_i^\mu(\boldsymbol{x}_i^{(t,k)})\right)\right\|^2$$

$$\leq 3\mathbb{E}_t\underbrace{\left\|\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\nabla f_i(\boldsymbol{x}_i^{(t,k)})-\frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\nabla f_i(\boldsymbol{x}_i^{(t,k)})\right\|^2}_{T_7}$$

$$+6H^2 L^2\mu^2, \tag{48}$$

where the first inequality follows from the Cauchy-Schwartz inequality and the second inequality follows by the Jensen's inequality and [10, Lemma 2].

By substituting $\mp\frac{1}{M}\sum_{i\in\mathcal{M}_t}\sum_{k=0}^{H-1}\nabla f_i(\boldsymbol{x}^t)$ and $\mp\frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\nabla f_i(\boldsymbol{x}^t)$ into $T_7$, and then following a similar derivation for bounding $T_6$, we can bound $T_7$ as follows

$$T_7 \leq 18HL^2\mathbb{E}_{\zeta_t}\left[\frac{1}{N}\sum_{i=1}^{N}\sum_{k=0}^{H-1}\left\|\boldsymbol{x}_i^{(t,k)}-\boldsymbol{x}^t\right\|^2\right]$$

$$+9H^2\,\mathbb{E}_{\mathcal{M}_t}\underbrace{\left\|\frac{1}{M}\sum_{i\in\mathcal{M}_t}\nabla f_i\left(\boldsymbol{x}^t\right)-\nabla f\left(\boldsymbol{x}^t\right)\right\|^2}_{T_8}. \tag{49}$$

We continue by bounding $T_8$ as below

$$T_8 = \mathbb{E}_{\mathcal{M}_t} \left\| \frac{1}{M} \sum_{i \in \mathcal{M}_t} \left( \nabla f_i \left( \boldsymbol{x}^t \right) - \nabla f \left( \boldsymbol{x}^t \right) \right) \right\|^2$$

$$= \frac{1}{M^2} \mathbb{E}_{\mathcal{M}_t} \sum_{i \in \mathcal{M}_t} \left\| \nabla f_i \left( \boldsymbol{x}^t \right) - \nabla f \left( \boldsymbol{x}^t \right) \right\|^2$$

$$\leq \frac{\sigma_h^2}{M} + \frac{c_h}{M} \left\| \nabla f \left( \boldsymbol{x}^t \right) \right\|^2, \tag{50}$$

where the second identity follows by the fact that devices in $\mathcal{M}_t$ are independently sampled from device set $\{1, 2, \ldots, N\}$ and $\mathbb{E}_{i \sim \{1,2,\ldots,N\}} \left[ \nabla f_i \left( \boldsymbol{x}^t \right) \right] = \nabla f \left( \boldsymbol{x}^t \right)$, and the last inequality comes from Assumption 4.

By combining (47), (48), (49), and (50), we bound $T_5$ as

$$T_5 \leq 18 H L^2 \mathbb{E}_{\zeta_t} \left[ \frac{1}{N} \sum_{i=1}^{N} \sum_{k=0}^{H-1} \left\| \boldsymbol{x}_i^{(t,k)} - \boldsymbol{x}^t \right\|^2 \right] + \frac{9 c_h H^2}{M} \left\| \nabla f \left( \boldsymbol{x}^t \right) \right\|^2$$

$$+ \frac{9 H^2 \sigma_h^2}{M} + 6 H^2 L^2 \mu^2 + \mathbb{E}_{\zeta_t} \left\| \frac{1}{N} \sum_{i=1}^{N} \sum_{k=0}^{H-1} \nabla f_i^{\mu} \left( \boldsymbol{x}_i^{(t,k)} \right) \right\|^2. \tag{51}$$

By combining (31), (35), (36), (46), and (51), we obtain Lemma 3.

*2) Proof of Lemma 1:* Most of the intermediary results for proving Lemma 3 in Appendix D1 can be directly applied here by replacing $M$ with $N$. The main difference is that there is no randomness in $\mathcal{M}_t$, and we thus do not need to construct an upper bound for $T_5$ in (36). By combining (31), (35), (36), and (46), we obtain Lemma 1.

*3) Proof of Lemma 2:* By denoting $\frac{1}{N} \sum_{i=1}^{N} \mathbb{E}_{\zeta_t^{k-1}} \left\| \boldsymbol{x}_i^{(t,k)} - \boldsymbol{x}^t \right\|^2$ as $s^{(t,k)}$, we have

$$s^{(t,\tau)} = \eta^2 \frac{1}{N} \sum_{i=1}^{N} \mathbb{E}_{\zeta_t^{\tau-1}} \left\| \sum_{k=0}^{\tau-1} \boldsymbol{e}_i^{(t,k)} \right\|^2$$

$$\leq \tau \eta^2 \sum_{k=0}^{\tau-1} \frac{1}{N} \sum_{i=1}^{N} \mathbb{E}_{\zeta_t^k} \left\| \boldsymbol{e}_i^{(t,k)} \right\|^2, \tag{52}$$

where the inequality follows by the Cauchy-Schwartz inequality. By combining (44), (45), and (52), we have

$$s^{(t,\tau)} \leq 6 \tilde{c}_g L^2 \tau \eta^2 \sum_{k=0}^{\tau-1} s^{(t,k)} + 6 \tau^2 \eta^2 \tilde{c}_g \tilde{c}_h \left\| \nabla f(\boldsymbol{x}^t) \right\|^2$$

$$+ 2 \tau^2 \eta^2 \tilde{\sigma}^2 + 2 L^2 \tau^2 \eta^2 \mu^2 + \frac{d^2 L^2 \tau^2}{2 b_1 b_2} \eta^2 \mu^2. \tag{53}$$

By taking summation over $\tau$ from 1 to $H-1$, we obtain

$$\sum_{\tau=1}^{H-1} s^{(t,\tau)} \leq 6 \tilde{c}_g L^2 \eta^2 \sum_{\tau=1}^{H-1} \tau \sum_{k=0}^{\tau-1} s^{(t,k)} + C_0$$

$$\leq 3 \tilde{c}_g H^2 L^2 \eta^2 \sum_{k=0}^{H-1} s^{(t,k)} + C_0, \tag{54}$$

where we utilize the property of arithmetic sequence and

$$C_0 = 2 H^3 \eta^2 \tilde{c}_g \tilde{c}_h \left\| \nabla f(\boldsymbol{x}^t) \right\|^2 + \frac{2}{3} H^3 \eta^2 \tilde{\sigma}^2 + \frac{2}{3} L^2 H^3 \eta^2 \mu^2$$

$$+ \frac{d^2 L^2 H^3}{6 b_1 b_2} \eta^2 \mu^2.$$

As $s^{(t,0)} = 0$ and by rearranging (54), we have

$$\left( 1 - 3 \tilde{c}_g H^2 L^2 \eta^2 \right) \sum_{\tau=0}^{H-1} s^{(t,\tau)} \leq 2 H^3 \eta^2 \tilde{c}_g \tilde{c}_h \left\| \nabla f(\boldsymbol{x}^t) \right\|^2$$

$$+ \frac{2}{3} L^2 H^3 \eta^2 \mu^2 + \frac{2}{3} H^3 \eta^2 \tilde{\sigma}^2 + \frac{d^2 L^2 H^3}{6 b_1 b_2} \eta^2 \mu^2. \tag{55}$$

As $\eta \leq \frac{1}{3 H L \sqrt{\tilde{c}_g}}$, we have $3 \left( 1 - 3 \tilde{c}_g H^2 L^2 \eta^2 \right) \geq 2$. Note that $\sum_{\tau=0}^{H-1} s^{(t,\tau)} = \delta_t$. We thus obtain Lemma 2.

*4) Proof of Lemma 4:* Lemma 4 is a byproduct of the derivation of Lemma 2. It follows from (53) by setting $\tau = H$.

## REFERENCES

[1] K. B. Letaief, Y. Shi, J. Lu, and J. Lu, "Edge artificial intelligence for 6G: Vision, enabling technologies, and applications," *IEEE J. Sel. Areas Commun.*, pp. 1–1, 2021, doi:10.1109/JSAC.2021.3126076.

[2] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, "Communication-Efficient Learning of Deep Networks from Decentralized Data," in *Proc. Int. Conf. Artif. Intell. Stat. (AISTATS)*, 2017.

[3] Y. Shi, K. Yang, Z. Yang, and Y. Zhou, *Mobile Edge Artificial Intelligence: Opportunities and Challenges*. Elsevier, 2021.

[4] Y. Qiang, "Federated recommendation systems," in *Proc. IEEE Int. Conf. Big Data*, 2019.

[5] K. Yang, Y. Shi, Y. Zhou, Z. Yang, L. Fu, and W. Chen, "Federated machine learning for intelligent IoT via reconfigurable intelligent surface," *IEEE Netw.*, vol. 34, no. 5, pp. 16–22, Sept. 2020.

[6] X. Zhang, M. Hong, S. Dhople, W. Yin, and Y. Liu, "FedPD: A federated learning framework with adaptivity to non-IID data," *IEEE Trans. Signal Process.*, vol. 69, pp. 6055–6070, Oct. 2021.

[7] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, "A novel framework for the analysis and design of heterogeneous federated learning," *IEEE Trans. Signal Process.*, vol. 69, pp. 5234–5249, 2021.

[8] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smithy, "FedDANE: A federated newton-type method," in *Proc. IEEE Asilomar Conf. Signals, Systems, Computers (ACSSC)*, 2019.

[9] Z. Dai, B. K. H. Low, and P. Jaillet, "Federated bayesian optimization via thompson sampling," in *Proc. Neural Inf. Process. Syst. (NeurIPS)*, 2020.

[10] X. Yi, S. Zhang, T. Yang, T. Chai, and K. H. Johansson, "Zeroth-order algorithms for stochastic distributed nonconvex optimization," *arXiv preprint arXiv:2106.02958*, 2021. [Online]. Available: https://arxiv.org/pdf/2106.02958.pdf

[11] J. Ren, Y. He, D. Wen, G. Yu, K. Huang, and D. Guo, "Scheduling for cellular federated edge learning with importance and channel awareness," *IEEE Trans. Wireless Commun.*, vol. 19, no. 11, pp. 7690–7703, Nov. 2020.

[12] Z. Wang, Y. Shi, Y. Zhou, H. Zhou, and N. Zhang, "Wireless-powered over-the-air computation in intelligent reflecting surface-aided IoT networks," *IEEE Internet Things J.*, vol. 8, no. 3, pp. 1585–1598, 2021.

[13] G. Zhu, J. Xu, K. Huang, and S. Cui, "Over-the-air computing for wireless data aggregation in massive IoT," *IEEE Wirel. Commun.*, vol. 28, no. 4, pp. 57–65, Aug. 2021.

[14] M. Fu, Y. Zhou, Y. Shi, W. Chen, and R. Zhang, "UAV aided over-the-air computation," *IEEE Trans. Commun.*, pp. 1–1, 2021, doi:10.1109/TWC.2021.3134327.

[15] W. Fang, Y. Jiang, Y. Shi, Y. Zhou, W. Chen, and K. B. Letaief, "Over-the-air computation via reconfigurable intelligent surface," *IEEE Trans. Commun.*, vol. 69, no. 12, pp. 8612–8626, Dec. 2021.

[16] K. Yang, T. Jiang, Y. Shi, and Z. Ding, "Federated learning via over-the-air computation," *IEEE Trans. Wireless Commun.*, vol. 19, no. 3, pp. 2022–2035, Mar. 2020.

[17] M. M. Amiri and D. Gündüz, "Machine learning at the wireless edge: Distributed stochastic gradient descent over-the-air," *IEEE Trans. Signal Process.*, vol. 68, pp. 2155–2169, Mar. 2020.

[18] 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, Jan. 2020.

[19] S. Liu, B. Kailkhura, P.-Y. Chen, P. Ting, S. Chang, and L. Amini, "Zeroth-order stochastic variance reduction for nonconvex optimization," *Proc. Neural Inf. Process. Syst. (NeurIPS)*, 2018.

[20] X. Gao, B. Jiang, and S. Zhang, "On the information-adaptive variants of the ADMM: An iteration complexity perspective," *J. Sci. Comput.*, vol. 76, no. 1, pp. 327–363, 2018.

[21] H. Yu, S. Yang, and S. Zhu, "Parallel restarted SGD with faster convergence and less communication: Demystifying why model averaging works for deep learning," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2019.

[22] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of FedAvg on Non-IID data," in *Proc. Int. Conf. Learn. Representations (ICLR)*, 2020.

[23] J. Wang, Z. Charles, Z. Xu, G. Joshi, H. B. McMahan, M. Al-Shedivat, G. Andrew, S. Avestimehr, K. Daly, D. Data *et al.*, "A field guide to federated optimization," *arXiv preprint arXiv:2107.06917*, 2021. [Online]. Available: https://arxiv.org/pdf/2107.06917.pdf

[24] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," in *Proc. Mach. Learn. Syst. (MLSys)*, 2020.

[25] S. P. Karimireddy, S. Kale, M. Mohri, S. J. Reddi, S. U. Stich, and A. T. Suresh, "SCAFFOLD: stochastic controlled averaging for federated learning," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2020.

[26] R. Pathak and M. J. Wainwright, "Fedsplit: an algorithmic framework for fast federated optimization," in *Proc. Neural Inf. Process. Syst. (NeurIPS)*, 2020.

[27] S. Wang, F. Roosta-Khorasani, P. Xu, and M. W. Mahoney, "GIANT: Globally improved approximate newton method for distributed optimization," in *Proc. Neural Inf. Process. Syst. (NeurIPS)*, 2018.

[28] D. Hajinezhad, M. Hong, and A. Garcia, "ZONE: Zeroth-order nonconvex multiagent optimization over networks," *IEEE Trans. Automat. Contr.*, vol. 64, no. 10, pp. 3995–4010, Oct. 2019.

[29] Y. Tang, J. Zhang, and N. Li, "Distributed zero-order algorithms for nonconvex multiagent optimization," *IEEE Trans. Control. Netw. Syst.*, vol. 8, no. 1, pp. 269–281, Mar. 2021.

[30] Z. Li and L. Chen, "Communication-efficient decentralized zeroth-order method on heterogeneous data," in *Int. Conf. Wirel. Commun. Signal Process. (WCSP)*, Changsha, China, 2021, pp. 1–6.

[31] A. K. Sahu and S. Kar, "Decentralized zeroth-order constrained stochastic optimization algorithms: Frank–wolfe and variants with applications to black-box adversarial attacks," *Proc. IEEE*, vol. 108, no. 11, pp. 1890–1905, 2020.

[32] W. Li and M. Assaad, "Distributed zeroth-order stochastic optimization in time-varying networks," *arXiv preprint arXiv: 2105.12597*, 2021. [Online]. Available: https://arxiv.org/abs/2105.12597

[33] A. K. Sahu, D. Jakovetic, D. Bajovic, and S. Kar, "Distributed zeroth order optimization over random networks: A kiefer-wolfowitz stochastic approximation approach," in *2018 IEEE Conference on Decision and Control (CDC)*. IEEE, 2018, pp. 4951–4958.

[34] A. Khaled, K. Mishchenko, and P. Richtárik, "Tighter theory for local SGD on identical and heterogeneous data," in *Proc. Int. Conf. Artif. Intell. Stat. (AISTATS)*, 2020.

[35] H. Liu, X. Yuan, and Y.-J. A. Zhang, "Reconfigurable intelligent surface enabled federated learning: A unified communication-learning design approach," *IEEE Trans. Wireless Commun.*, vol. 20, no. 11, pp. 7595–7609, Nov. 2021.

[36] X. Cao, G. Zhu, J. Xu, Z. Wang, and S. Cui, "Optimized power control design for over-the-air federated edge learning," *IEEE J. Sel. Areas Commun.*, pp. 1–1, 2021, doi:10.1109/JSAC.2021.3126060.

[37] Z. Wang, J. Qiu, Y. Zhou, Y. Shi, L. Fu, W. Chen, and K. B. Letaief, "Federated learning via intelligent reflecting surface," *IEEE Trans. Wireless Commun.*, pp. 1–1, 2021, doi:10.1109/TWC.2021.3099505.

[38] T. Sery, N. Shlezinger, K. Cohen, and Y. C. Eldar, "Over-the-air federated learning from heterogeneous data," *IEEE Trans. Signal Process.*, vol. 69, pp. 3796–3811, Jun. 2021.

[39] R. Paul, Y. Friedman, and K. Cohen, "Accelerated gradient descent learning over multiple access fading channels," *IEEE J. Sel. Areas Commun.*, pp. 1–1, 2021, doi:10.1109/JSAC.2021.3118410.

[40] S. Liu, P.-Y. Chen, B. Kailkhura, G. Zhang, A. O. Hero III, and P. K. Varshney, "A primer on zeroth-order optimization in signal processing and machine learning: Principals, recent advances, and applications," *IEEE Signal Process. Mag.*, vol. 37, no. 5, pp. 43–54, Sept. 2020.

[41] S. Liu, J. Chen, P.-Y. Chen, and A. Hero, "Zeroth-order online alternating direction method of multipliers: Convergence analysis and applications," in *Proc. Int. Conf. Artif. Intell. Stat. (AISTATS)*, 2018.

[42] L. Bottou, F. E. Curtis, and J. Nocedal, "Optimization methods for large-scale machine learning," *SIAM Review*, vol. 60, no. 2, pp. 223–311, 2018.

[43] Y. Nesterov and V. Spokoiny, "Random gradient-free minimization of convex functions," *Foundations of Computational Mathematics*, vol. 17, no. 2, pp. 527–566, 2017.

[44] S. Ghadimi and G. Lan, "Stochastic first- and zeroth-order methods for nonconvex stochastic programming," *SIAM J. Optim.*, vol. 23, no. 4, pp. 2341–2368, 2013.

[45] O. Abari, H. Rahul, D. Katabi, and M. Pant, "Airshare: Distributed coherent transmission made seamless," in *Proc. IEEE Conf. Comput. Commun. (INFOCOM)*, 2015, pp. 1742–1750.

[46] A. Mahmood, M. I. Ashraf, M. Gidlund, J. Torsner, and J. Sachs, "Time synchronization in 5G wireless edge: Requirements and solutions for critical-mtc," *IEEE Commun. Mag.*, vol. 57, no. 12, pp. 45–51, 2019.

[47] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *Proc. Symp. Security Privacy (SP)*, 2017.

[48] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017. [Online]. Available: https://arxiv.org/abs/1708.07747