Distributed Minimax Fair Optimization over Hierarchical Networks

Wen Xu*, Juncheng Wang[†], Ben Liang^{*}, Gary Boudreau[‡], Hamza Sokun[‡] *University of Toronto, Canada, [†]Hong Kong Baptist University, Hong Kong, [‡]Ericsson, Canada

ABSTRACT

In modern applications, the underlying computation and communication networks are often hierarchical, which is typified by the three-layer client-edge-cloud system that has become prominent in recent times. We study minimax fairness in distributed optimization over such systems, to provide robust performance guarantee for the worst-case mixture of loss functions. We propose HIERMINIMAX, a communication efficient distributed algorithm to solve the minimax optimization problem. We provide convergence analysis for both convex and non-convex loss functions, leading to performance bounds that enable tuning the tradeoff between the communication complexity and the optimization convergence rate. Our experiments on classification problems with canonical datasets show that HIERMINIMAX substantially improves the fairness in learning accuracy and reduces the communication overhead compared with the current best alternatives.

CCS CONCEPTS

Computing methodologies → Distributed algorithms; Machine learning;
 Networks;

KEYWORDS

Minimax optimization, distributed learning, hierarchical networks

ACM Reference Format:

Wen Xu, Juncheng Wang, Ben Liang, Gary Boudreau, and Hamza Sokun. 2024. Distributed Minimax Fair Optimization over Hierarchical Networks. In *The 53rd International Conference on Parallel Processing (ICPP '24), August* 12–15, 2024, Gotland, Sweden. ACM, New York, NY, USA, 10 pages. https: //doi.org/10.1145/3673038.3673137

1 INTRODUCTION

Our work is motivated by federated learning (FL) [23] as a prime example of distributed optimization. In FL, multiple local clients collaboratively train a machine learning (ML) model with the assistance of a server. The standard objective in FL is to find a best global model w that minimizes a weighted sum of the local loss functions at individual clients:

$$\min_{\mathbf{w}} \sum_{n \in \mathcal{N}} q_n f_n(\mathbf{w}), \tag{1}$$

where N is the set of clients, $f_n(w)$ is the local loss function at the *n*th client, and $q_n \ge 0$ is some given local weight. It is common to

ICPP '24, August 12-15, 2024, Gotland, Sweden

© 2024 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-1793-2/24/08

https://doi.org/10.1145/3673038.3673137

choose q_n such that it is proportional to the amount of data at the *n*th client.

FL is typically considered in a *two-layer client-server* architecture [14, 18, 38]. In edge computing, the clients are mobile or IoT devices and the server can be either an edge server or a cloud server [3]. Exemplary algorithms under this architecture include Federated Averaging (FEDAVG) [23] and Local Stochastic Gradient Descent (LOCAL-SGD) [34]. For each training round, each client updates its local model via SGD from a common global model using its local data, and then the server aggregates all the updated local models to generate a new global model.

However, in many practical applications, the underlying communication networks are multi-layer hierarchical. An example is to perform FL as a cloud-based distributed service, where the clients are small remote devices and the server is in some far away cloud center [3]. Therefore, the communication between the clients and the server must go through some intermediate edge server. Such hierarchical communication architecture was ignored in [3, 14, 18, 23, 34, 38]. In contrast, a *three-layer client-edge-cloud* FL system was considered and an optimization algorithm termed HI-ERFAVG was proposed in [21]. Furthermore, an extension to enable multi-step client-edge model aggregation and quantization was studied in [22].

A persistent challenge in FL and distributed optimization in general is that the heterogeneity of local data distributions can significantly harm the optimization performance [14, 18, 38]. Specifically, the computed global model that solves (1) can perform poorly for individual clients if the local data distributions differ significantly. It can also perform poorly when the data ratios of clients in training do not match that of the unseen data in reality, e.g., when the clients process different amounts of data during training due to their differing processing capabilities. To remedy this fairness issue, *minimax* optimization was proposed in [25] as a special case of distributionally robust optimization [11, 30]. In minimax optimization, the weight vector $\boldsymbol{q} = [q_1, \ldots, q_N]$ becomes another set of optimization variables. This new optimization problem is given by

$$\min_{\mathbf{w}} \max_{\mathbf{q} \in Q} F(\mathbf{w}, \mathbf{q}), \text{ where } F(\mathbf{w}, \mathbf{q}) = \sum_{n \in \mathcal{N}} q_n f_n(\mathbf{w}), \qquad (2)$$

and Q is a subset of the (N-1)-dimensional probability simplex. The global model w that solves (2) is guaranteed to be robust against heterogeneity of data distributions. In particular, when constraint Q is inactive, the learned model w performs well even for the worstoff client that has the least favorable local data distribution. The Stochastic Agnostic Federated Learning (STOCHASTIC-AFL) algorithm [25] was proposed to solve the minimax problem (2). It was later extended to accommodate multi-step local updates in Distributionally Robust Federated Averaging (DRFA) [10]. In these methods, for each training round, both the global model w and the weight vector q are updated. However, these works are limited to

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

the conventional two-layer client-server architecture and cannot be applied to the client-edge-cloud network.

This motivates us to design a distributed algorithm to solve the minimax optimization problem (2) over a multi-layer hierarchical network. In particular, we are interest in a solution that has provable performance guarantees. To achieve this objective, we must address several challenges in multiple fronts: 1) The intermediate edge servers have unique impacts on minimax optimization, such as client-edge model aggregation, that has not been explored in existing literature. 2) Minimax optimization requires joint consideration of the mutual impacts of hierarchical model updates and weight updates on algorithm convergence performance. 3) The multi-step local model updates at the clients and multi-step client-edge aggregations can lead to asymmetric synchronization gaps, as each update of the weight vector can only be performed once for multiple steps of model update and aggregation. In this context, the main contributions of this paper are as follows:

- We formulate a constrained minimax optimization problem over the client-edge-cloud architecture. The model parameters and the weight vector for edge areas are jointly optimized to achieve minimax fairness across all edge areas under a general convex constraint. To the best of our knowledge, distributed minimax fair optimization over multi-layer hierarchical networks has not been considered before.
- We propose a distributed algorithm called HIERMINIMAX, which utilizes the client-edge-cloud architecture to reduce the communication overhead via both multi-step local model updates at the clients and multi-step client-edge aggregation. We adopt a checkpoint mechanism, which enables periodic updates of the weight vector for efficient communication. Furthermore, HIERMINIMAX allows partial participation of edge areas, local model updates via local SGD, and most importantly, flexibility of communication frequencies between the clients and the edge servers, as well as between the edge servers and the cloud server.
- We analyze the convergence of HIERMINIMAX to show that it provides performance guarantees for both convex and non-convex loss functions. We explore the communicationconvergence trade-off enabled by HIERMINIMAX, showing how we can reduce the communication complexity without overly degrading the convergence rate in distribution minimax optimization. Specifically, for any $\alpha \in [0, 1)$, we can achieve $O(T^{1-\alpha})$ edge-cloud communication complexity with $O(\frac{1}{T^{(1-\alpha)/2}})$ convergence rate for convex loss functions and $O(T^{1-\alpha})$ edge-cloud communication complexity with $O(\frac{1}{T^{(1-\alpha)/4}})$ convergence rate for non-convex loss functions, where *T* is the number of training time slots.
- For numerical evaluation, we experiment with standard classification datasets for both convex logistic regression and non-convex neural network training. Our experimental results demonstrate that HIERMINIMAX substantially improves the worst-case learning performance and reduces the communication overhead over the current best alternatives.

The rest of this paper is structured as follows. In Section 2, we present the related work. In Section 3, we describe the system model of multi-layer hierarchical networks. We present HIERMINIMAX in

Table 1: Summary of related works on distributed minimax optimization (Hier.: hierarchical, c.c.: communication complexity, c.r.: convergence rate).

		Convex loss		Non-convex loss	
Reference	Hier	. c.c.	c.r.	c.c.	c.r.
[25]	X	O(T)	$O(rac{1}{T^{1/2}})$	N/A	N/A
[10]	X	$O(T^{3/4})$	$O(\frac{1}{T^{3/8}})$	$O(T^{3/4})$	$O(\frac{1}{T^{1/8}})$
Ours	1	$O(T^{1-\alpha})$	$O(\frac{1}{T^{(1-\alpha)/2}})$	$O(T^{1-\alpha})$	$O(\frac{1}{T^{(1-\alpha)/4}})$

Section 4 and its convergence analysis in Section 5. Numerical experiments are provided in Section 6. We conclude the paper in Section 7.

2 RELATED WORK

We provide a literature review on hierarchical FL and minimax optimization. The differences between existing works on distributed minimax optimization and our work are summarized in Table 1.

2.1 Hierarchical FL

The standard FL algorithms were designed for two-layer clientserver systems [14, 18, 38]. To take advantages of both the connectivity to a large number of clients at the cloud and the availability of low-latency communication at the edge, a three-layer clientedge-cloud FL system was considered in [21, 22]. The HIERFAVG algorithm in [21] naturally extends FEDAVG to multi-layer networks, and the HIER-LOCAL-QSGD algorithm in [22] further extends HIER-FAVG with model quantization. A similar study in [5] considered heterogeneous operating rates and showed the dependence of the convergence on the average operating rate, the network topology, and the number of iterations. A convergence bound of hierarchical SGD over non-independent data was derived based on upward and downward divergences in [36]. However, none of these works considered optimizing the model performance over the worst mixture of the data distributions.

2.2 Minimax Optimization

We are interested in learning over multiple distributions via minimax optimization [13, 32] to achieve distributionally robust learning [11, 30]. Minimax optimization has been studied in the game theory and optimization literature [1, 27, 35]. First-order methods for minimax optimization include Gradient Descent Ascent (GDA) [9, 20], Extra-gradient (EG) [16, 24], and Optimistic Gradient Descent Ascent (OGDA) [7]. With the advent and popularity of ML, many ML problems are found to have the minimax structure, such as GANs [12], online learning [31], and adversarial robust learning [33]. However, all of these minimax optimization algorithms require *centralized* implementation.

The authors in [25] were the first to give the formulation of (2) in the context of FL and proposed STOCHASTIC-AFL, which is a stochastic, distributed, and privacy-preserving version of GDA. More recently, in [10], DRFA was proposed to further accommodate a more general communication pattern of FL through periodic model averaging and weight vector updating. Both [10] and [25] are only applicable to conventional two-layer client-server networks. In this work, we study minimax optimization over the client-edge-cloud Distirbuted Minimax Fair Optimization over Hierarchical Networks



Figure 1: System model of hierarchical networks.

architecture. In particular, HIERMINIMAX differs from DRFA [10] in several key aspects: 1) HIERMINIMAX allows both flexible multi-step local model updates and multi-step client-edge aggregation (see Section 4), 2) HIERMINIMAX is guided by theoretical analysis of the convergence rate and communication complexity caused by hierarchical update and aggregation (see Section 5), and 3) HIER-MINIMAX achieves substantial empirical advantages over DRFA (see Section 6).

3 MINIMAX OPTIMIZATION OVER **CLIENT-EDGE-CLOUD ARCHITECTURE**

We consider a multi-layer hub-and-spoke-type network topology. Since the three-layer client-edge-cloud network architecture is common in practical systems, we use it as an representative example as shown in Fig. 1. All edge servers can directly communicate with the cloud server. The cloud server can indirectly communicate with any client through an edge server.

Let N be the set of all clients, whose cardinality is |N| = N. Let \mathcal{E} be the set of all edge servers, whose cardinality is $|\mathcal{E}| = N_{\mathcal{E}}$. The set of clients that is associated with edge server $e \in \mathcal{E}$ is denoted by N_e . We define an edge area as an edge server and all clients associated with it. For convenience of notation, we assume the edge servers have the same number of clients, i.e., $|\mathcal{N}_e| = N_0, \forall e \in \mathcal{E}$. Hence, we have $N = N_0 N_{\mathcal{E}}$. Our work can be easily generalized to the case where different edge servers have different numbers of clients.

The clients, edge servers, and cloud server collaborate to compute a global model w. We assume the clients within each edge area $e \in \mathcal{E}$ share the same local loss function, denoted by $f_e(w)$. In the FL example, this means that the data samples of clients in each edge area are generated from the same distribution [14, 18, 38]. This is without loss of generality, since if there are multiple distributions in an edge area, we can group clients of the same distribution and the edge server as a virtual edge area. Our goal is to jointly optimize *w* and the edge weights, denoted by $p = [p_1, \ldots, p_{N_E}]$, by minimizing a worst-case global loss, i.e., for the worst mixture of local loss functions (e.g., data distributions in FL) at the edge. Thus, problem (2) takes the following new form for a three-layer client-edge-cloud network:

$$\min_{\boldsymbol{w}\in\mathcal{W}}\max_{\boldsymbol{p}\in\mathcal{P}}F(\boldsymbol{w},\boldsymbol{p}), \text{ where } F(\boldsymbol{w},\boldsymbol{p}) = \sum_{e\in\mathcal{E}}p_ef_e(\boldsymbol{w}), \qquad (3)$$

Algorithm 1 Hierarchical Distributed Minimax Optimization (HIERMINIMAX)

- 1: Cloud initializes $\boldsymbol{w}^{(0)}$ and $\boldsymbol{p}^{(0)} = [1/N_{\mathcal{E}}, \dots, 1/N_{\mathcal{E}}]$.
- 2: for each round $k = 0, \ldots, \overline{K} 1$ do // Phase 1
- Cloud samples edge servers $\mathcal{E}^{(k)} \subseteq \mathcal{E}$ by $\boldsymbol{p}^{(k)}$. 3:
- Cloud samples (c_1, c_2) uniformly from $[\tau_1] \times [\tau_2]$. 4:
- Cloud broadcasts $w^{(k)}$ and (c_1, c_2) to $\mathcal{E}^{(k)}$. 5:
- 6:
- for each edge server $e \in \mathcal{E}^{(k)}$ do $w_e^{(k,\tau_2)}, w_e^{(k,c_2,c_1)} = \text{MODELUPDATE}(w^{(k)}, c_1, c_2).$ 7:
- Cloud aggregates the global model via (5). 8:
- Cloud computes the checkpoint model via (6). 9 // Phase 2
- Cloud samples edges $\mathcal{U}^{(k)} \subseteq \mathcal{E}$ uniformly. 10:
- Cloud broadcasts $\mathbf{w}^{(k,c_2,c_1)}$ to $\mathcal{U}^{(k)}$. 11:
- for each edge server $e \in \mathcal{U}^{(k)}$ do 12:
- $f_e(\mathbf{w}^{(k,c_2,c_1)}) = \text{LossEstimation}(\mathbf{w}^{(k,c_2,c_1)}).$ 13
- Cloud updates weight vector via (7). 14:

 $\mathcal{W} \subseteq \mathbb{R}^d$, and $\mathcal{P} \subseteq \Delta_{N_{\mathcal{E}}-1}$ with $\Delta_{N_{\mathcal{E}}-1}$ being the $(N_{\mathcal{E}}-1)$ -dimensional probability simplex. We allow both W and \mathcal{P} to be any compact convex set.¹

A naive approach to solve (3) is to directly adopt STOCHASTIC-AFL as follows. For each training round, the cloud server broadcasts the current global model to all the clients through the edge servers. Each client updates the model via single-step local SGD and sends it through an edge server to the cloud server for global aggregation. The edge weights are then updated by gradient ascent utilizing the loss estimation of the global model among the edge areas. However, the communication overhead of this approach is high, as communication between the edge areas and the cloud server is required per training round of both the model update and the weight update.

Ideally, a communication-efficient algorithm should allow multistep local model updates at the clients and multi-step client-edge aggregations per training round. However, this will require a way to evenly sample the local loss values between the aggregation instances over the entire training process, which is needed to update the edge weights. We follow this idea and propose HIERMINIMAX, which we will provide provable performance guarantees.

HIERARCHICAL DISTRIBUTED MINIMAX 4 ALGORITHM

The HIERMINIMAX algorithm is carried out over K training rounds. Each training round corresponds to one update of the global model parameters w and the edge weights p by the cloud server. In each training round, τ_2 client-edge model aggregations are performed. Each client-edge aggregation in turn is performed after τ_1 local SGD steps at the clients. Thus, the total number of training time slots is $T = K \tau_1 \tau_2$.

¹Note that if $\mathcal{P} = \Delta_{N_E-1}$, the optimization over \boldsymbol{p} is essentially unconstrained and the optimal minimax solution maximizes performance for the worst-off edge area that has the least favorable loss function. Our formulation here is more general, where \mathcal{P} may represent, e.g., prior knowledge or parameter regularization.



Figure 2: One training round for the updates of model parameters and the weight vector for HIERMINIMAX.

We next discuss the detailed design of HIERMINIMAX, while its formal specification is given in Algorithm 1. As shown in Fig. 2, each training round of HIERMINIMAX contains *two phases*. In Phase 1, the cloud server updates the global model parameters w after τ_2 client-edge model aggregations, which corresponds to $\tau_1 \tau_2$ local SGD steps at the clients. In this phase, the algorithm also obtains a random checkpoint of an intermediate model. In Phase 2, the cloud server updates the edge weights p based on the checkpoint model generated in Phase 1.

4.1 Phase 1: Model Parameter Update

In Phase 1 of each training round k, the cloud server first samples $m_{\mathcal{E}}$ edge servers, denoted by a set $\mathcal{E}^{(k)}$, based on the probability defined by the edge weights $p^{(k)}$ of the current round. The cloud server also samples a *checkpoint* index (c_1, c_2) from $[\tau_1] \times [\tau_2]$ uniformly at random. The cloud server then broadcasts the current global model $w^{(k)}$ and the checkpoint index (c_1, c_2) to each sampled edge server. Each sampled edge server $e \in \mathcal{E}^{(k)}$ runs a MODELUP-DATE procedure based on $w^{(k)}$ and (c_1, c_2) received from the cloud. The MODELUPDATE procedure contains two parts.

Part (a): Edge Model Update. Each edge server $e \in \mathcal{E}^{(k)}$ performs τ_2 client-edge model aggregations. At the $(t_2 + 1)$ -th client-edge aggregation, each edge server e first broadcasts the model $\mathbf{w}_e^{(k,t_2)}$, which is the model at the edge server e after t_2 client-edge aggregations in training round k, to all clients \mathcal{N}_e associated with the edge server. Each client $n \in \mathcal{N}_e$ then performs τ_1 steps of local SGD starting from $\mathbf{w}_e^{(k,t_2)}$, i.e., at each step $t_1 \in \{0, \ldots, \tau_1-1\}$, each client n performs the following SGD update:

$$\mathbf{w}_{n}^{(k,t_{2},t_{1}+1)} = \Pi_{\mathcal{W}}(\mathbf{w}_{n}^{(k,t_{2},t_{1})} - \eta_{w}\nabla_{w}f_{n}(\mathbf{w}_{n}^{(k,t_{2},t_{1})};\xi_{n}^{(t_{1})})), \quad (4)$$

where $\mathbf{w}_n^{(k,t_2,t_1)}$ is the local model after t_1 local model update and t_2 client-edge aggregation in *k*-th training round, η_w is the learning rate on the model update, $\xi_n^{(t_1)}$ is a mini-batch training examples sampled by client *n*, and $\Pi_W(\cdot)$ is the projection onto set W. After finishing τ_1 steps of local SGD in (4), each client $n \in N_e$ sends back $\mathbf{w}_n^{(k,t_2,\tau_1)}$ to the edge server $e \in \mathcal{E}^{(k)}$. Finally, each edge server $e \in \mathcal{E}^{(k)}$ performs client-edge model aggregation via $\mathbf{w}_e^{(k,t_2+1)} = \frac{1}{N_0} \sum_{n \in N_e} \mathbf{w}_n^{(k,t_2,\tau_1)}$.

Part (b): Checkpoint Model Update. Each edge server $e \in \mathcal{E}^{(k)}$ obtains the checkpoint index (c_1, c_2) from the cloud server. It broadcasts c_1 to its clients \mathcal{N}_e at the beginning of the c_2 -th step of local aggregation. Then, when each client $n \in \mathcal{N}_e$ sends back its updated model $w_n^{(k, t_2, c_1)}$ to edge server $e \in \mathcal{E}^{(k)}$, it also sends along the checkpoint model $w_n^{(k, c_2, c_1)}$. Each edge server $e \in \mathcal{E}^{(k)}$ updates its aggregate checkpoint model via $w_e^{(k, c_2, c_1)} = \frac{1}{N_0} \sum_{n \in \mathcal{N}_e} w_n^{(k, c_2, c_1)}$. The checkpoint model at the edge server will be used to approximate the evolution of the local models for updating p in Phase 2 to guarantee convergence.

After parts (a) and (b) of the MODELUPDATE procedure, each sampled edge server $e \in \mathcal{E}^{(k)}$ sends $w_e^{(k,\tau_2)}$ and $w_e^{(k,c_2,c_1)}$ to the cloud server. The cloud server then performs edge-cloud model aggregation to update the global model via

$$\mathbf{w}^{(k+1)} = \frac{1}{m_{\mathcal{E}}} \sum_{e \in \mathcal{E}^{(k)}} \mathbf{w}_{e}^{(k,\tau_{2})}.$$
 (5)

The cloud server also updates a new checkpoint model via

$$\mathbf{w}^{(k,c_2,c_1)} = \frac{1}{m_{\mathcal{E}}} \sum_{e \in \mathcal{E}^{(k)}} \mathbf{w}_e^{(k,c_2,c_1)}.$$
 (6)

4.2 Phase 2: Edge Weight Update

In each training round k, the cloud server uniformly samples a set of edge servers $\mathcal{U}^{(k)}$ of size $m_{\mathcal{E}}$ and broadcasts the checkpoint model $\mathbf{w}^{(k,c_2,c_1)}$ to these edge servers. Note that $\mathcal{U}^{(k)}$ can be different from $\mathcal{E}^{(k)}$ in Phase 1. For each sampled edge server $e \in \mathcal{U}^{(k)}$, it runs the LossEstimation procedure. Specifically, each sampled edge server $e \in \mathcal{U}^{(k)}$ broadcasts the checkpoint model $\mathbf{w}^{(k,c_2,c_1)}$ to its clients N_e . Each client $n \in N_e$ calculates $f_n(\mathbf{w}^{(k,c_2,c_1)};\xi_n^{(k)})$, i.e., the loss on the checkpoint model via a mini-batch of samples $\xi_n^{(k)}$. After receiving the loss calculated from its clients, each edge server $e \in \mathcal{U}^{(k)}$ estimates the loss of the checkpoint model over its data distribution via $f_e(\mathbf{w}^{(k,c_2,c_1)}) = \frac{1}{N_0} \sum_{n \in N_e} f_n(\mathbf{w}^{(k,c_2,c_1)};\xi_n^{(k)})$.

Each sampled edge server $e \in \mathcal{U}^{(k)}$ then sends its estimated loss to the cloud server. After receiving the loss estimations from all the sampled edge servers, the cloud server constructs a vector $v = [v_1, \ldots, v_{N_{\mathcal{E}}}]$ where $v_e = \frac{N_{\mathcal{E}}}{m_{\mathcal{E}}} f_e(\mathbf{w}^{(k,c_2,c_1)})$ if $e \in \mathcal{U}^{(k)}$, else $v_e = 0$. Recall that the minimax optimization objective in (3) is $F(w, p) = \sum_{e \in \mathcal{E}} p_e f_e(w)$. The *e*-th coordinate of its gradient with respect to p is $[\nabla_p F(w, p)]_e = f_e(w), \forall w \in \mathcal{W}$. The expectation of the *e*-th coordinate of v is $\mathbb{E}[v_e] = \frac{m_{\mathcal{E}}}{N_{\mathcal{E}}} \frac{N_{\mathcal{E}}}{m_{\mathcal{E}}} f_e(w^{(k,c_2,c_1)}) + (1 - \frac{m_{\mathcal{E}}}{N_{\mathcal{E}}}) \times 0 = f_e(w^{(k,c_2,c_1)})$. Hence, the constructed v is an unbiased estimator of $\nabla_p F(w^{(k,c_2,c_1)}, p)$. The cloud server then updates the edge weights via projected gradient ascent, given by

$$\boldsymbol{p}^{(k+1)} = \Pi_{\mathcal{P}}(\boldsymbol{p}^{(k)} + \eta_p \tau_1 \tau_2 \boldsymbol{v}), \tag{7}$$

where η_p is the learning rate on p and $\Pi_{\mathcal{P}}(\cdot)$ is the projection onto set \mathcal{P} .

5 CONVERGENCE ANALYSIS

In this section, we show that HIERMINIMAX provides guaranteed convergence performance for both convex and non-convex local loss functions. We note that new analysis techniques are required to capture the three-layer client-edge-cloud architecture, as well as multi-step local SGD updates and multi-step client-edge aggregations. All proof details are given in the appendix with proof prerequisites given in Appendix A.

We make the following assumptions, which are common in existing works on distributed minimax optimization [10, 20, 25].

Assumption 1 (Bounded Domains). The diameters of the compact convex sets W and P are R_W and R_P , respectively.

Assumption 2 (*L*-smoothness). There exists some positive *L* such that $\|\nabla f_n(\mathbf{w}_1) - \nabla f_n(\mathbf{w}_2)\| \le L \|\mathbf{w}_1 - \mathbf{w}_2\|, \forall \mathbf{w}_1, \mathbf{w}_2 \in \mathcal{W}, n \in \mathcal{N} \text{ and } \|\nabla F(\mathbf{w}_1, \mathbf{p}_1) - \nabla F(\mathbf{w}_2, \mathbf{p}_2)\| \le L \|(\mathbf{w}_1, \mathbf{p}_1) - (\mathbf{w}_2, \mathbf{p}_2)\|, \forall (\mathbf{w}_1, \mathbf{p}_1), (\mathbf{w}_2, \mathbf{p}_2) \in (\mathcal{W}, \mathcal{P}).$

Assumption 3 (Bounded Gradients). There exist some positive G_w and G_p such that $\|\nabla_w f_n(w)\|_2 \le G_w$ and $\|\nabla_p F(w, p)\|_2 \le G_p$.

Assumption 4 (Bounded Stochastic Gradient Variance). There exist some positive σ_w^2 and σ_p^2 such that $\mathbb{E}[\|\nabla f_n(w;\xi) - \nabla f_n(w)\|_2^2] \le \sigma_w^2, \forall n \in \mathcal{N}, w \in \mathcal{W}, and \mathbb{E}[\|\nabla_p F(w, p; \xi) - \nabla_p F(w, p)\|_2^2] \le \sigma_p^2, \forall w \in \mathcal{W} and p \in \mathcal{P}.$

Assumption 5 (Bounded Gradient Dissimilarity). Let $\nabla f_e(\mathbf{w})$ be the gradient of the edge loss function $f_e(\mathbf{w})$ for any $e \in \mathcal{E}$. There exists some positive Ψ such that the gradient dissimilarity is bounded by $\sup_{\mathbf{w}\in \mathcal{W}, \mathbf{p}\in \mathcal{P}, e\in \mathcal{E}} \sum_{j\in \mathcal{E}} p_j ||\nabla f_e(\mathbf{w}) - \nabla f_j(\mathbf{w})||^2 \leq \Psi$.

5.1 Convex Loss

We first analyze the performance of HIERMINIMAX when the local loss functions are convex in w, i.e., $f_n(w_1) \ge f_n(w_2) + \nabla f_n(w_2)^T(w_1 - w_2)$, $\forall w_1, w_2 \in W$, $\forall n \in N$. Then the global loss function F(w, p) is convex in w. Since F(w, p) is also linear in p, it is a convex-concave function. The standard way to measure the optimality of a solution (\hat{w}, \hat{p}) for constrained convex-concave optimization is the duality gap, given by

$$\max_{\boldsymbol{p}\in\mathcal{P}}F(\hat{\boldsymbol{w}},\boldsymbol{p}) - \min_{\boldsymbol{w}\in\mathcal{W}}F(\boldsymbol{w},\hat{\boldsymbol{p}}),\tag{8}$$

where $\hat{\boldsymbol{p}} = \frac{1}{K} \sum_{k=0}^{K-1} \boldsymbol{p}^{(k)}$ is the time-averaged edge weights, $\hat{\boldsymbol{w}} = \frac{1}{mT} \sum_{t=0}^{T-1} \sum_{n \in S^{(t)}} \boldsymbol{w}_n^{(t)}$ is the time-averaged model parameters, and $S^{(t)}$ is the set of sampled clients at iteration *t*, whose cardinality

is $m = m_{\mathcal{E}}N_0$. Let $(\boldsymbol{w}^*, \boldsymbol{p}^*)$ be a minimax point, in the sense that a Nash equilibrium is established [26]. The duality gap of $(\boldsymbol{w}^*, \boldsymbol{p}^*)$ is zero as $\max_{\boldsymbol{p}\in\mathcal{P}}F(\boldsymbol{w}^*, \boldsymbol{p}) = F(\boldsymbol{w}^*, \boldsymbol{p}^*) = \min_{\boldsymbol{w}\in\mathcal{W}}F(\boldsymbol{w}, \boldsymbol{p}^*)$ for convex-concave functions by the von Neumann minimax theorem [35]. Hence, the lower the duality gap is, the better the solution $(\hat{\boldsymbol{w}}, \hat{\boldsymbol{p}})$ is, as an approximation of a Nash equilibrium.

We require a key lemma in our convergence analysis which bounds the squared distance between any local model $w_n^{(t)}$ and the virtual global model $w^{(t)} = \frac{1}{m} \sum_{n \in S^{(t)}} w_n^{(t)}$.

Lemma 1 (Bounded Squared Model Divergence). For Algorithm 1 with convex loss, assuming $1-20\eta_w^2 L^2 \tau_1^2 (1+\tau_2^2) \geq \frac{1}{2}$, we have

$$\begin{split} &\frac{1}{mT}\sum_{t=0}^{T-1}\sum_{n\in\mathcal{S}^{(t)}}\mathbb{E}[\|\boldsymbol{w}^{(t)}-\boldsymbol{w}^{(t)}_{n}\|^{2}]\\ &\leq 20\eta_{w}^{2}\tau_{1}^{2}(\frac{m+1}{m}\sigma_{w}^{2}+\Psi)+20\eta_{w}^{2}\tau_{1}^{2}\tau_{2}^{2}(\frac{m\varepsilon+1}{N_{0}}\sigma_{w}^{2}+\Psi). \end{split}$$

PROOF. See Appendix B.

Leveraging Lemma 1, we have the following theorem.

THEOREM 1. The duality gap achieved by HIERMINIMAX for convex loss is upper bounded by

$$\mathbb{E}\left[\max_{\boldsymbol{p}\in\mathcal{P}}F(\hat{\boldsymbol{w}},\boldsymbol{p}) - \min_{\boldsymbol{w}\in\mathcal{W}}F(\boldsymbol{w},\boldsymbol{p})\right]$$

$$\leq \underbrace{\frac{R_{\varphi}^{2}}{2\eta_{p}T} + \frac{\eta_{p}\tau_{1}\tau_{2}}{2}G_{p}^{2} + \frac{\eta_{p}\tau_{1}\tau_{2}}{2m}\sigma_{p}^{2}}_{maximization \ gap} + \underbrace{\frac{N_{\mathcal{E}}R_{\mathcal{W}}^{2}}{2\eta_{w}T} + \frac{\eta_{w}N_{\mathcal{E}}}{2}G_{w}^{2} + \frac{\eta_{w}}{2N_{0}}\sigma_{w}^{2}}_{minimization \ gap}}_{chineter \ edge \ aggregation} + \underbrace{\underbrace{10LN_{\mathcal{E}}\eta_{w}^{2}\tau_{1}^{2}\left(\frac{m+1}{m}\sigma_{w}^{2}+\Psi\right)}_{client-edge \ aggregation} + \underbrace{\underbrace{10LN_{\mathcal{E}}\eta_{w}^{2}\tau_{1}^{2}\tau_{2}^{2}\left(\frac{m_{\mathcal{E}}+1}{N_{0}}\sigma_{w}^{2}+\Psi\right)}_{edge-cloud \ aggregation}.$$

PROOF. We require the following two lemmas: Lemma 3 bounds the update of model w and Lemma 4 bounds the update on weight p. The proof details are given in Appendix C.

Theorem 1 provides a means to tune the tradeoff between the communication complexity and the convergence rate. Suppose we set τ_1 and τ_2 such that $\tau_1\tau_2 \in \Theta(T^{\alpha})$ for any $\alpha \in [0, 1)$. The communication complexity between the edge and cloud servers is $\Theta(T^{1-\alpha})$. Let $\eta_P = \Theta(\frac{1}{T^{(1+\alpha)/2}})$, and let $\eta_w = \Theta(\frac{1}{T^{1-2\alpha}})$ if $\alpha \in (0, \frac{1}{4})$ and $\eta_w = \Theta(\frac{1}{T^{1/2}})$ if $\alpha \in [\frac{1}{4}, 1)$. From Theorem 1, the convergence rate of HIERMINIMAX is $O(\frac{1}{T^{(1-\alpha)/2}})$. Note that α is a tunable value. As α increases, the communication complexity $O(T^{1-\alpha})$ decreases but the convergence rate $O(\frac{1}{T^{(1-\alpha)/2}})$ becomes worse, i.e., we can trade convergence rate to reduce communication complexity.

One extreme case is to let $\tau_1 = \tau_2 = 1$. Then, the communication complexity is O(T) and the convergence rate is $O(\frac{1}{T^{1/2}})$. This recovers the same scaling in *T* for STOCHASTIC-AFL in [25]. Another interesting and more general case is $\tau_2 = 1$. When substituting this into Theorem 1, we recover a duality gap bound that scales in *T* and τ_1 in the same way as DRFA in [10]. We further note that HIER-MINIMAX allows any choice from a wider range of trade-off points beyond these special cases.

ICPP '24, August 12-15, 2024, Gotland, Sweden

5.2 Non-convex Loss

The loss functions of many modern ML applications such as neural networks training are non-convex. A Nash equilibrium may not exist in this setting and the duality gap is no longer a meaningful measure of optimality [20]. We first define $\Phi: \mathcal{W} \to \mathbb{R}$ such that $\Phi(w) = \max_{p \in \mathcal{P}} F(w, p)$. Hence, our optimization formulation (3) is equivalent to $\min_{w \in \mathcal{W}} \Phi(w)$. Since for any given $w \in \mathcal{W}, F(w, \cdot)$ is linear in $p, \Phi(w)$ can be computed efficiently. However, since Φ itself is non-convex in w, the problem of finding the global minimum of Φ is still in general NP-hard. Therefore, in the non-convex optimization literature, the stationary point is commonly used for measuring the optimality of the solutions [4]. However, the function Φ is not necessarily differentiable, making it improper to directly use the gradient $\nabla \Phi$.

We first follow the approach in [8, 20] to define a Moreau envelope to facilitate our convergence analysis. The λ -Moreau envelope of a function Φ with a positive parameter λ is

$$\Phi_{\lambda}(\boldsymbol{w}) = \min_{\boldsymbol{x} \in \mathcal{W}} \left\{ \Phi(\boldsymbol{x}) + (1/2\lambda) \|\boldsymbol{x} - \boldsymbol{w}\|^2 \right\}.$$
(9)

We choose $\lambda = 1/2L$ such that $\Phi_{1/2L}(\cdot)$ is differentiable [20, Lemma 3.6]. Then the stationarity of the function Φ can be approximately measured by $\|\nabla \Phi_{1/2L}(\cdot)\|$.

We observe that the analysis methods in [8, 20, 29] are for algorithms solving minimax optimization of nonconvex-concave objectives with the same updated frequencies on both sets of optimization variables. Therefore, they are not directly applicable to the analysis of HIERMINIMAX, as the edge weights are updated for one time after both multi-step local client model updates and multi-step client-edge aggregations are performed. As part of the new analysis required to bound the performance of HIERMINIMAX, we first present the following lemma to bound the distance between any local model $w_n^{(t)}$ and the virtual global model $w^{(t)} = \frac{1}{m} \sum_{n \in S^{(t)}} w_n^{(t)}$.

LEMMA 2 (BOUNDED DIVERGENCE OF MODEL). For Algorithm 1 with non-convex loss, assuming $1 - 2\eta_w L \tau_1(1 + \tau_2) \ge \frac{1}{2}$, we have

$$\frac{1}{mT} \sum_{t=0}^{T-1} \frac{1}{m} \sum_{n \in S^{(t)}} \mathbb{E} \left[\left\| \boldsymbol{w}^{(t)} - \boldsymbol{w}_n^{(t)} \right\| \right] \\
\leq 2\eta_w \tau_1 \left(\frac{m+1}{m} \sigma_w + \sqrt{\Psi} \right) + 2\eta_w \tau_1 \tau_2 \left(\frac{m_{\mathcal{E}} + 1}{N_0} \sigma_w + \sqrt{\Psi} \right). \quad (10)$$

PROOF. The proof follows similar steps as the proof of Lemma 1, and it is omitted due to space constraint. The difference is that we bound $\mathbb{E}[\|\boldsymbol{w}^{(t)}-\boldsymbol{w}_n^{(t)}\|]$ instead of $\mathbb{E}[\|\boldsymbol{w}^{(t)}-\boldsymbol{w}_n^{(t)}\|^2]$ in Lemma 1. \Box

Leveraging Lemma 2, we have the following theorem.

THEOREM 2. The time-averaged expectation of the squared norm of the (1/2L)-Moreau envelope of Φ , achieved by HIERMINIMAX for non-convex loss, is upper bounded by

$$\begin{split} &\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \Big[\| \nabla \Phi_{1/2L}(\boldsymbol{w}^{(t)}) \|^2 \Big] \\ &\leq \frac{4 \Phi_{1/2L}(\boldsymbol{w}^{(0)})}{\eta_w N_{\mathcal{E}} T} + 16L \sqrt{K} \eta_w \tau_1 \tau_2 G_w \sqrt{G_w^2 + \sigma_w^2} \\ &+ 4L \frac{R_{\mathcal{P}}^2}{\sqrt{K} \eta_p \tau_1 \tau_2} + 8 \eta_p \tau_1 \tau_2 L (G_p^2 + \frac{\sigma_p^2}{m}) \end{split}$$

Xu et al.

$$+\frac{4\eta_{w}}{N_{\mathcal{E}}}(G_{w}^{2}+\frac{\sigma_{w}^{2}}{m})+\frac{8\eta_{w}\tau_{1}R_{W}L^{2}}{N_{\mathcal{E}}}(\frac{m+1}{m}\sigma_{w}+\sqrt{\Psi})$$
$$+\frac{8\eta_{w}\tau_{1}\tau_{2}R_{W}L^{2}}{N_{\mathcal{E}}}(\frac{m_{\mathcal{E}}+1}{N_{0}}\sigma_{w}+\sqrt{\Psi}).$$
(11)

PROOF. The proof consists of the following key steps. We first derive a bound for the update of $w^{(t+1)}$ with respect to $\Phi_{1/2L}$ and the LHS of (10) in Lemma 2. Then we apply Lemma 4 to bound one update of p in the proof of Theorem 1. Finally, we use Lemma 2 and the block technique in general non-convex-concave optimization [20], which controls the gap between p and its optimal value at each t, to complete the proof. Further details are omitted due to space constraint.

Similarly to the convex case, Theorem 2 provides a means to tune the tradeoff between the communication complexity and the convergence rate. Let $\tau_1 \tau_2 \in \Theta(T^{\alpha})$ for any $\alpha \in [0, 1)$. The communication complexity between the edge and cloud servers is $\Theta(T^{1-\alpha})$. Let $\eta_p = \Theta(\frac{1}{T^{(1+\alpha)/4}})$ and $\eta_w = \Theta(\frac{1}{T^{(3+\alpha)/4}})$, the convergence rate is $O(\frac{1}{T^{(1-\alpha)/4}})$. As α increases, the communication complexity $O(T^{1-\alpha})$ decreases but the convergence rate $O(\frac{1}{T^{(1-\alpha)/4}})$ becomes worse.

For the special case $\tau_1 = \tau_2 = 1$, HIERMINIMAX has communication complexity O(T) and convergence rate $O(\frac{1}{T^{1/4}})$. We note that this is the same as the best possible convergence rate for *centralized* minimax optimization with nonconvex loss in [20]. HIERMINIMAX can achieve this convergence rate while incurring maximum communication complexity, but it can also be tuned to reduce its convergence speed in exchange for higher communication efficiency. For the other special case $\tau_2 = 1$, it is easy to check from Theorem 2 that HIERMINIMAX recovers a Moreau envelope bound that scales in *T* and τ_1 in the same way as DRFA in [10].

6 NUMERICAL EXPERIMENTS

In addition to deriving the communication complexity and convergence scaling results shown in the previous section, we further conduct experiments to study the numerical performance of HI-ERMINIMAX. We use distributed machine learning as an example application. Our experiments are conducted via PyTorch version 2.0.1 [28].

We consider the following benchmarks for performance comparison:

- FEDAVG [23]. This is the standard method for FL. It does not use edge servers. It solves the minimization problem (1) via LOCAL SGD with multi-step local model update in each training round.
- STOCHASTIC-AFL [25]. See Section 3 for details. This method does not use edge servers. It solves the minimax problem (3) by single-step local model update in each training round.
- DRFA [10]. This method is similar to STOCHASTIC-AFL, except that it solves the minimax problem (3) by multi-step local model update in each training round.
- HIERFAVG [21]. This method uses the same three-layer clientedge-cloud architecture as HIERMINIMAX but solves the minimization problem (1).

Distirbuted Minimax Fair Optimization over Hierarchical Networks



Figure 3: Average and worst test accuracies with convex loss functions (for EMNIST-Digits).



Figure 4: Average and worst test accuracies with non-convex loss functions (for Fashion-MNIST).

For all methods with multi-step model updates, we set $\tau_1 = 2$. Furthermore, for methods utilizing hierarchical architectures, we set $\tau_2 = 2$, i.e., $\tau_1 = \tau_2 = 2$ for both HIERFAVG and HIERMINIMAX, so that both models have the same amount of model updates.

6.1 Convex Loss Functions

We consider multinomial logistic regression as the model and crossentropy loss as the loss function. They are applied to the image classification dataset EMNIST-Digits [6], which contains 10 classes of hand-written digits. We set $N_{\mathcal{E}} = 10$, $N_0 = 3$, $m_{\mathcal{E}} = 5$, $\mathcal{W} = \mathbb{R}^{7850}$, and $\mathcal{P} = \Delta_9$. To create heterogeneous data distributions, we assign one distinct class of training data to the clients of each edge area. We use SGD of learning rate $\eta_w = 0.001$ with batch size 1 for each local model update and learning rate of $\eta_p = 0.001$ for the weight vector update.

Fig. 3 shows a comparison of the average test accuracy and worst test accuracy among the clients. Clearly, the three methods that solve the minimax problem (3) provide much better worst-case performance, while paying only a small price on the average performance. Furthermore, HIERMINIMAX substantially outperforms the two-layer minimax methods. Specifically, to reach 80% worst accuracy, HIERMINIMAX takes only 8200 communication rounds, compared with 16652 rounds for STOCHASTIC-AFL and 11727 rounds for DRFA, corresponding to overhead reduction of 51% and 30%, respectively. HIERMINIMAX also substantially outperforms the three-layer minimization method HIERFAVG with communication overhead reduction of 55% from 18228 rounds. We further note that FEDAVG does not reach 80% worst accuracy even after 20000 communication rounds.

6.2 Non-convex Loss Functions

For non-convex loss, we consider a two hidden-layer fully-connected neural network, with 300 and 100 neurons in the hidden-layers, ReLU as the activation function, and cross-entropy loss as the loss function. It is applied to perform a more difficult image classification ICPP '24, August 12-15, 2024, Gotland, Sweden

Tabl	e 2: (Comparison	of HIERFAVG	and HIERMINIMAX
------	--------	------------	-------------	-----------------

Methods	Average	Worst	Variance
HierFAVG	0.9070	0.8035	21.0504
HierMinimax	0.8999	0.8348	5.5657
HierFAVG	0.8072	0.4829	206.6945
HierMinimax	0.7631	0.6051	24.7095
HierFAVG	0.8703	0.7572	30.9331
HierMinimax	0.8501	0.7818	20.2926
HIERFAVG	0.8180	0.6453	76.2957
HierMinimax	0.8123	0.7323	16.3589
HIERFAVG	0.7539	0.2102	732.6033
HierMinimax	0.7250	0.2896	478.8593
	Methods HIERFAVG HIERMINIMAX HIERFAVG HIERMINIMAX HIERFAVG HIERMINIMAX HIERFAVG HIERMINIMAX	MethodsAverageHIERFAVG0.9070HIERMINIMAX0.8999HIERFAVG0.8072HIERFAVG0.8073HIERFAVG0.8703HIERFAVG0.8703HIERFAVG0.8180HIERMINIMAX0.8123HIERFAVG0.7539HIERMINIMAX0.7250	Methods Average Worst HIERFAVG 0.9070 0.8035 HIERMINIMAX 0.8909 0.8348 HIERFAVG 0.8072 0.4829 HIERFAVG 0.8072 0.4829 HIERFAVG 0.8703 0.7572 HIERFAVG 0.8703 0.7572 HIERFAVG 0.8180 0.6453 HIERFAVG 0.8180 0.6453 HIERFAVG 0.8180 0.2102 HIERFAVG 0.7539 0.2102 HIERMINIMAX 0.7250 0.2896

task on Fashion-MNIST [37], which contains 10 classes of clothes. We set $N_{\mathcal{E}} = 10$, $N_0 = 3$, $m_{\mathcal{E}} = 2$, $\mathcal{W} = \mathbb{R}^{266610}$, and $\mathcal{P} = \Delta_9$. To create heterogeneous data distributions, we adopt a more controllable way as in [15]: for *s*% similarity we allocate to each edge area *s*% i.i.d. data and the remaining (100 - s)% by sorting according to label. Here we present the case s = 50. We use SGD of learning rate $\eta_w = 0.001$ with batch size 8 for each local model update and learning rate of $\eta_p = 0.0001$ for the weight vector update.

Fig. 4 shows a comparison of the average test accuracy and worst test accuracy among the clients. Again, the three methods that solve the minimax problem (3) provide substantial improvements on the worst-case performance, compared with their counterparts solving the minimization problem (1), while maintaining similar average performance. HIERMINIMAX again substantially outperforms the two-layer minimax methods and three-layer minimization methods. Specifically, to reach 50% worst accuracy, HIERMINIMAX takes 21576 communication rounds, which is a 52% reduction from the 45201 communication rounds needed by STOCHASTIC-AFL, a 23% reduction from the 28087 communication rounds by DRFA, and a 41% reduction from the 36445 communication rounds by HIERFAVG. We further note that FEDAVG does not reach 50% worst accuracy even after 50000 communication rounds.

6.3 Minimax Fairness and Variance

We further investigate the fairness of HIERMINIMAX in terms of test accuracies achieved by different edge areas. In addition to the average and worst test accuracy discussed earlier, we further compare the variance of test accuracies over edge areas to explicitly consider the uniformity of test accuracies. The first part of Table 2 contains the results of logistic regression models solving image classification problems on EMNIST-Digits [6], Fashion-MNIST [37], and MNIST [17]. We observe that HIERMINIMAX typically yields only slightly lower average test accuracy compared with HIERFAVG. However, HIERMINIMAX achieves much higher worst test accuracy, especially for a harder dataset such as Fashion-MNIST. Meanwhile, the variance of test accuracies over different edge areas achieved by HIERMINIMAX is much lower than that of HIERFAVG, as much as one order of magnitude in the case of Fashion-MNIST.

Table 2 additionally shows the results for salary prediction on the Adult dataset [2] and sample classification on the Synthetic dataset [19]. For Adult, we consider 2 edge areas, each with data samples from Doctorate and non-Doctorate groups, respectively. We train a logistic regression model on categorical features with $\eta_w = 0.001$ and $\eta_p = 0.0001$. For Synthetic, we consider 100 edge areas and train a logistic regression model with $\eta_w = \eta_p = 0.0001$. Similar to [19], we report the worst 10% accuracy. We again observe substantial advantage of HIERMINIMAX over HIERFAVG on both Adult and Synthetic.

7 CONCLUSION

We propose HIERMINIMAX for distributed minimax optimization over a three-layer network architecture exemplified by the common client-edge-cloud system in mobile edge computing. Our convergence analysis, for both convex and non-convex loss functions, sheds light on the tradeoff between the communication complexity and the convergence rate. In experiments with standard ML datasets, we further show that HIERMINIMAX has substantial advantages in terms of communication overhead reduction and worst-case performance, compared with the existing two-layer approaches and the three-layer minimization approach.

ACKNOWLEDGMENTS

This work was supported by Ericsson, the Natural Sciences and Engineering Research Council of Canada (NSERC), and Mitacs.

REFERENCES

- Tamer Basar and Geert Jan Olsder. 1998. Dynamic Noncooperative Game Theory. Society for Industrial and Applied Mathematics. https://doi.org/10.1137/1. 9781611971132
- [2] Barry Becker and Ronny Kohavi. 1996. Adult. UCI Machine Learning Repository.[3] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex
- [1] Kim Denne, et al. 2019. Towards federated learning at scale: System design. In Proc. of the SysML Conf.
- [4] Léon Bottou, Frank E. Curtis, and Jorge Nocedal. 2018. Optimization methods for large-scale machine learning. SIAM Rev. 60, 2 (2018), 223–311.
- [5] Timothy Castiglia, Anirban Das, and Stacy Patterson. 2021. Multi-level local SGD: Distributed SGD for heterogeneous hierarchical networks. In Proc. of ICLR.
- [6] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. 2017. EMNIST: Extending MNIST to handwritten letters. In *Proc. of IJCNN*.
 [7] Constantinos Daskalakis and Ioannis Panageas. 2018. The limit points of (opti-
- mistic) gradient descent in min-max optimization. In *Proc. of NeurIPS*.
 [8] Damek Davis and Dmitriy Drusvyatskiy. 2019. Stochastic model-based minimiza-
- tion of weakly convex functions. *SIAM J. on Optim.* 29, 1 (2019), 207–239.
 V.F. Dem'yanov and A.B. Pevnyi. 1972. Numerical methods for finding saddle
- points. Comput. Math. and Math. Phys. 12, 5 (1972), 11-52. [10] Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. 2020. Distri-
- butionally robust federated averaging. In *Proc. of NeurIPS*. [11] John C. Duchi and Hongseok Namkoong. 2021. Learning models with uniform
- [11] John C. Duchi and Hongseok Nankoong. 2021. Learning models with uniform performance via distributionally robust optimization. *The Ann. of Statist.* 49, 3 (2021), 1378 – 1406.
- [12] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, et al. 2014. Generative adversarial nets. In Proc. of NeurIPS.
- [13] Nika Haghtalab, Michael Jordan, and Eric Zhao. 2022. On-demand sampling: Learning optimally from multiple distributions. In Proc. of NeurIPS.
- [14] Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, et al. 2021. Advances and open problems in federated learning. *Found.* and Trends in Mach. Learn. 14, 1-2 (2021), 1–210.
- [15] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, et al. 2020. SCAFFOLD: Stochastic controlled averaging for federated learning. In *Proc. of ICML*.
- [16] Galina M Korpelevich. 1976. The extragradient method for finding saddle points and other problems. *Matecon* 12 (1976), 747–756.
- [17] Yann LeCun, Corinna Cortes, and CJ Burges. 2010. MNIST handwritten digit database. ATT Labs [Online]. Available: http://yann.lecun.com/exdb/mnist (2010).
- [18] Tian Li, Anit Kumar Sahu, Ameet S. Talwalkar, and Virginia Smith. 2020. Federated learning: Challenges, methods, and future directions. *IEEE Signal Process. Mag.* 37 (2020), 50–60.
- [19] Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. 2020. Fair resource allocation in federated learning. In Proc. of ICLR.
- [20] Tianyi Lin, Chi Jin, and Michael Jordan. 2020. On gradient descent ascent for nonconvex-concave minimax problems. In Proc. of ICML.

- [21] Lumin Liu, Jun Zhang, SH Song, and Khaled B Letaief. 2020. Client-edge-cloud hierarchical federated learning. In Proc. of ICC.
- [22] Lumin Liu, Jun Zhang, Shenghui Song, and Khaled B Letaief. 2023. Hierarchical federated learning with quantization: Convergence analysis and system design. *IEEE Trans. on Wireless Commun.* 22, 1 (2023), 2–18.
- [23] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *Proc. of AISTATS*.
- [24] Panayotis Mertikopoulos, Bruno Lecouat, Houssam Zenati, Chuan-Sheng Foo, Vijay Chandrasekhar, et al. 2019. Optimistic mirror descent in saddle-point problems: Going the extra(-gradient) mile. In Proc. of ICLR.
- [25] Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. 2019. Agnostic federated learning. In Proc. of ICML.
- [26] John Nash. 1951. Non-cooperative games. Ann. of Math. (1951), 286-295.
- [27] Noam Nisan, Tim Roughgarden, Éva Tardos, and Vijay V. Vazirani. 2007. Algorithmic game theory. Cambridge University Press.
- [28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, et al. 2019. PyTorch: An imperative style, high-performance deep learning library. In Proc. of NeurIPS.
- [29] Hassan Rafique, Mingrui Liu, Qihang Lin, and Tianbao Yang. 2022. Weaklyconvex-concave min-max optimization: Provable algorithms and applications in machine learning. *Optim. Met. and Softw.* 37, 3 (2022), 1087–1121.
- [30] Hamed Rahimian and Sanjay Mehrotra. 2019. Distributionally robust optimization: A review. arXiv preprint arXiv:1908.05659 (2019).
- [31] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. 2015. Online learning via sequential complexities. J. of Mach. Learn. Res. 16, 6 (2015), 155–186.
- [32] Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. 2020. Distributionally robust neural networks. In Proc. of ICLR.
- [33] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. 2018. Adversarially robust generalization requires more data. In Proc. of NeurIPS.
- [34] Sebastian U Stich. 2019. Local SGD converges fast and communicates little. In Proc. of ICLR.
- [35] John von Neumann. 1928. Zur theorie der gesellschaftsspiele. 119, 1 (1928), 295–320.
- [36] Jiayi Wang, Shiqiang Wang, Rong-Rong Chen, and Mingyue Ji. 2022. Demystifying why local aggregation helps: Convergence analysis of hierarchical SGD. In Proc. of AAAI.
- [37] Han Xiao, Kashif Rasul, and Roland Vollgraf. 2017. Fashion-MNIST: A novel image dataset for benchmarking machine learning algorithms. https://github. com/zalandoresearch/fashion-mnist.
- [38] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. 2019. Federated machine learning: Concept and applications. ACM Trans. on Intell. Syst. and Technol. 10, 2 (jan 2019), 1–19.

A PREREQUISITES TO PROOFS

We first introduce the symbols and notations that will be used throughout our appendices. We use $k \in [K]$ to indicate the *k*-th training round and $t \in [T]$ to indicate the *t*-th training time slots (iterations) with $T = K\tau_1\tau_2$. Let the sampled clients after *k*-th global aggregation be $S^{(k)} = \bigcup_{e \in \mathcal{E}^{(k)}} N_e$. We denote the cardinality of $S^{(k)}$ by $m = m_{\mathcal{E}}N_0$, where $m_{\mathcal{E}}$ is the number of sampled edge servers and N_0 is the number of clients in each edge area. Note that for any $t \in [T]$, we have $S^{(t)} = S^{(k)}$ for any $k\tau_1\tau_2 \le t < (k+1)\tau_1\tau_2$.

To facilitate our proofs, we define the following auxiliary variables. At iteration *t*, let the average model of all selected clients be $\mathbf{w}^{(t)} = \frac{1}{m} \sum_{n \in S^{(t)}} \mathbf{w}_n^{(t)}$. The full gradient of the model $\mathbf{w}_n^{(t)}$ at client *n* is denoted by $g_n^{(t)} = \nabla f_n(\mathbf{w}_n^{(t)})$ and the stochastic gradient of the model $\mathbf{w}_n^{(t)}$ at client *n* is denoted by $h_n^{(t)} = \nabla f_n(\mathbf{w}_n^{(t)}; \xi_n^{(t)})$. The average full and stochastic gradient at edge server *e* are denoted by $g_e^{(t)} = \frac{1}{N_0} \sum_{n \in N_e} g_n^{(t)}, h_e^{(t)} = \frac{1}{N_0} \sum_{n \in N_e} h_n^{(t)}$. The average full and stochastic gradient of all selected clients are denoted by $g^{(t)} = \frac{1}{m} \sum_{n \in S^{(t)}} g_n^{(t)}, h^{(t)} = \frac{1}{m} \sum_{n \in S^{(t)}} h_n^{(t)}$. Recall that at iteration *t*, the weight vector *p* has been updated

Recall that at iteration *t*, the weight vector **p** has been updated $\lfloor \frac{t}{\tau_1 \tau_2} \rfloor$ times, so we denote the global minimax objective at time *t* by $F(\mathbf{w}^{(t)}, \mathbf{p}^{(\lfloor \frac{t}{\tau_1 \tau_2} \rfloor)})$. We note that the stochastic gradient $h^{(t)}$

is an unbiased estimate of the full gradient of $F(\mathbf{w}^{(t)}, \mathbf{p}^{(\lfloor \frac{t}{\tau_1 \tau_2} \rfloor)})$, i.e., $\mathbb{E}_{\mathcal{E}^{(t)}, \{\xi_n^{(t)}\}_{n \in S^{(t)}}}[h^{(t)}] = \mathbb{E}_{\mathcal{E}^{(t)}}[g^{(t)}] = \mathbb{E}[\sum_{e \in \mathcal{E}} p_e^{(\lfloor \frac{t}{\tau_1 \tau_2} \rfloor)} \nabla f_e(\mathbf{w}_e^{(t)})].$ Furthermore, from Assumption 4, we can show that the variance of $h^{(t)}$ is bounded by $\mathbb{E}\left[\|h^{(t)} - g^{(t)}\|^2\right] \leq \frac{\sigma_w^2}{m}$.

Now, we derive the bias and variance of the estimation of p. Recall that $v^{(k)}$ is the vector constructed in Phase 2 of training round k. Denote the stochastic gradient of p in this training round by $\boldsymbol{u}^{(k)} = \tau_1 \tau_2 \boldsymbol{v}^{(k)}$. The full gradient of $F(\boldsymbol{w}^{(t)}, \boldsymbol{p})$ with respect to \boldsymbol{p} is $\bar{v}^{(t)}$. Its *e*-th coordinate is $[\bar{v}^{(t)}]_e = f_e(\mathbf{w}^{(t)}) = [\nabla_p F(\mathbf{w}^{(t)}, \mathbf{p})]_e$ for all $e \in \mathcal{E}$. We further define $\bar{u}^{(k)} = \sum_{t=k\tau_1\tau_2+1}^{(k+1)\tau_1\tau_2} \bar{v}^{(t)}$, which is the summation of the virtual full gradient with respect to p over $\tau_1 \tau_2$ iterations within the k-th round. Then the stochastic gradient of \pmb{p} in our algorithm is unbiased, i.e., $\mathbb{E}\left[\boldsymbol{u}^{(k)}\right] = \mathbb{E}\left[\eta_p \tau_1 \tau_2 \boldsymbol{v}^{(k)}\right] = \bar{\boldsymbol{u}}^{(k)}$, as we sample a timestamp from an interval of length $\tau_1 \tau_2$ and sample each edge server uniformly and independently. Note that $u^{(k)}$ also has bounded variance $\mathbb{E}\left[\left\|\boldsymbol{u}^{(k)} - \bar{\boldsymbol{u}}^{(k)}\right\|^2\right] \leq \frac{\tau_1^2 \tau_2^2}{m} \sigma_p^2$.

B PROOF OF LEMMA 1

PROOF. We first bound each $A^{(t)}$ for $t \in [k\tau_1\tau_2, (k+1)\tau_1\tau_2]$, i.e., $t = k\tau_1\tau_2 + t_2\tau_1 + t_2$ where $0 \le t_2 < \tau_2$ and $0 \le t_1 \le \tau_1$. Let *e* be the edge server with which client n is associated. We define

$$C1 = \sum_{i=0}^{t_2-1} \frac{1}{N_0} \sum_{n' \in \mathcal{N}_e} \sum_{j=0}^{\tau_1-1} \eta_w h_{n'}^{(k,i,j)} \text{ and } C2 = \sum_{j=0}^{t_1-1} \eta_w h_n^{(k,t_2,j)}.$$

We have $\boldsymbol{w}_n^{(t)} = \boldsymbol{w}^{(k,0,0)} - C1 - C2$. We further define

$$C3 = \frac{1}{m_{\mathcal{E}}} \sum_{e' \in \mathcal{E}^{(t)}} \sum_{i=0}^{t_2-1} \frac{1}{N_0} \sum_{n'' \in \mathcal{N}_{e'}} \sum_{j=0}^{\tau_1-1} \eta_w h_{n''}^{(k,i,j)}$$
$$C4 = \frac{1}{m} \sum_{n' \in S^{(t)}} \sum_{j=0}^{t_1-1} \eta_w h_n^{(k,t_2,j)}.$$

We have $\mathbf{w}^{(t)} = \mathbf{w}^{(k,0,0)} - C3 - C4$. We define $A^{(t)} = \frac{1}{m} \sum_{n \in S^{(t)}} \mathbb{E}[\|\mathbf{w}_n^{(t)} - \mathbf{w}^{(t)}\|^2]$. Hence,

$$A^{(t)} \stackrel{(a)}{\leq} \frac{2}{m} \sum_{n \in S^{(t)}} \mathbb{E}[\|C1 - C3\|^2] + \frac{2}{m} \sum_{n \in S^{(t)}} \mathbb{E}[\|C2 - C4\|^2], \quad (12)$$

where (a) follows from $\|x + y\|^2 \le 2(\|x\|^2 + \|y\|^2)$.

We now bound the right-hand side (RHS) of (12). For the first term, we have

$$\frac{2}{m} \sum_{n \in S^{(t)}} \mathbb{E} \left[\|C1 - C3\|^{2} \right]
\stackrel{(a)}{\leq} \frac{2\eta_{w}^{2} t_{2} \tau_{1}}{m} \sum_{n \in S^{(t)}} \sum_{i=0}^{t_{2}-1} \sum_{j=0}^{\tau_{1}-1} \mathbb{E} \left[\|\sum_{n' \in \mathcal{N}_{e}} \frac{h_{n'}^{(k,i,j)}}{N_{0}} - \sum_{n'' \in S^{(t)}} \frac{h_{n''}^{(k,i,j)}}{m} \|^{2} \right]
\stackrel{(b)}{\leq} 10\eta_{w}^{2} t_{2} \tau_{1} \sum_{i=0}^{t_{2}-1} \sum_{j=0}^{\tau_{1}-1} \left(\frac{m_{\mathcal{E}}+1}{N_{0}} \sigma_{w}^{2} + 2L^{2} A^{(k,i,j)} + \Psi \right)
\stackrel{(c)}{\leq} 10\eta_{w}^{2} \tau_{2} \tau_{1} \sum_{i=0}^{\tau_{2}-1} \sum_{j=0}^{\tau_{1}-1} \left(\frac{m_{\mathcal{E}}+1}{N_{0}} \sigma_{w}^{2} + 2L^{2} A^{(k,i,j)} + \Psi \right), \quad (13)$$

ICPP '24, August 12-15, 2024, Gotland, Sweden

where (a) is by Jensen's inequality; (b) is by first both adding and subtracting four terms $\frac{1}{N_0} \sum_{n' \in N_e} g_{n'}^{(k,i,j)}, \frac{1}{N_0} \sum_{n' \in N_e} \nabla f_{n'}(\boldsymbol{w}^{(k,i,j)}),$ $\frac{1}{m}\sum_{n''\in S^{(t)}}\nabla f_{n''}(\boldsymbol{w}^{(k,i,j)}), \text{ and } \frac{1}{m}\sum_{n''\in S^{(t)}}g_{n''}^{(k,i,j)}, \text{ and then applying Jensen's inequality, along with Assumptions 2, 4, and 5; and$ (c) is by upper bounding t_2 by τ_2 . Applying the same techniques to the second term on the RHS of (12), we can show that

$$\frac{2}{m} \sum_{n \in S^{(t)}} \mathbb{E} \left[\left\| C2 - C4 \right\|^2 \right] \\
\leq 10 \eta_w^2 \tau_1 \sum_{j=0}^{\tau_1 - 1} \left(\frac{m+1}{m} \sigma_w^2 + 2L^2 A^{(k, t_2, j)} + \Psi \right).$$
(14)

Substituting (13) and (14) into (12) and summing the resulting inequality over t from $k\tau_1\tau_2$ to $(k + 1)\tau_1\tau_2 - 1$, we have

$$\begin{split} \sum_{t_2=0}^{\tau_2-1} \sum_{t_1=0}^{\tau_1-1} A^{(k,t_2,t_1)} \\ &= 10\eta_w^2 \tau_1^2 \sum_{t_2=0}^{\tau_2-1} \sum_{j=0}^{\tau_1-1} \left(\frac{m+1}{m} \sigma_w^2 + 2L^2 A^{(k,t_2,j)} + \Psi \right) \\ &+ 10\eta_w^2 \tau_1^2 \tau_2^2 \sum_{i=0}^{\tau_2-1} \sum_{j=0}^{\tau_1-1} \left(\frac{m\varepsilon+1}{N_0} \sigma_w^2 + 2L^2 A^{(k,i,j)} + \Psi \right) \\ &\stackrel{(a)}{\leq} 20\eta_w^2 \tau_1^3 \tau_2 \left(\frac{m+1}{m} \sigma_w^2 + \Psi \right) + 20\eta_w^2 \tau_1^3 \tau_2^3 \left(\frac{m\varepsilon+1}{N_0} \sigma_w^2 + \Psi \right), \quad (15) \end{split}$$

where (a) is by assuming $1 - 20\eta_w^2 L^2 \tau_1^2 (1 + \tau_2^2) \ge \frac{1}{2}$. Summing (15) over k from 0 to K - 1 and diving both sides by $T = K\tau_1\tau_2$, we complete the proof.

PROOF OF THEOREM 1 С

We first prove Lemma 3 and Lemma 4 to facilitate the proof of Theorem 1.

LEMMA 3 (Upper Bound for One Iteration of w). For all $w \in$ W and $t \in [T]$, we have

$$\mathbb{E}\left[\left\|\boldsymbol{w}^{(t+1)} - \boldsymbol{w}\right\|^{2}\right] - \mathbb{E}\left[\left\|\boldsymbol{w}^{(t)} - \boldsymbol{w}\right\|^{2}\right] \leq L\eta_{w}A^{(t)} + \eta_{w}^{2}G_{w}^{2} + \frac{\eta_{w}^{2}}{m}\sigma_{w}^{2} - \frac{2\eta_{w}}{N_{\mathcal{E}}}\mathbb{E}\left[F(\boldsymbol{w}^{(t)}, \boldsymbol{p}^{\left(\lfloor\frac{t}{\tau_{1}\tau_{2}}\rfloor\right)}) - F(\boldsymbol{w}, \boldsymbol{p}^{\left(\lfloor\frac{t}{\tau_{1}\tau_{2}}\rfloor\right)})\right].$$
(16)

PROOF. We have

$$\mathbb{E}[\|\mathbf{w}^{(t+1)} - \mathbf{w}\|^{2}] \stackrel{(a)}{\leq} \mathbb{E}[\|\mathbf{w}^{(t)} - \eta_{w}h^{(t)} - \mathbf{w}\|^{2}] \\
\stackrel{(b)}{=} \mathbb{E}[\|\mathbf{w}^{(t)} - \eta_{w}g^{(t)} - \mathbf{w}\|^{2}] + \eta_{w}^{2}\mathbb{E}[\|g^{(t)} - h^{(t)}\|^{2}] \\
\stackrel{(c)}{\leq} \mathbb{E}[\|\mathbf{w}^{(t)} - \mathbf{w}\|^{2}] + A_{1}^{(t)} + A_{2}^{(t)} + \frac{\eta_{w}^{2}}{m}\sigma_{w}^{2}, \quad (17)$$

where (a) follows from the update rule in (4) and projection onto a closed convex set, (b) follows from both adding and subtracting $\eta_w q^{(t)}$, expanding the squares and eliminating the zero-valued cross-term, and (c) is from bounded variance of $h^{(t)}$ and definitions of $A^{(t)} = n^2 \mathbb{E}[||a^{(t)}||^2]$ and $A^{(t)} = -2n \mathbb{E}[(a^{(t)})^T (w^{(t)} - w)]$

$$\begin{aligned} & X_2 = \eta_{\mathcal{W}} \mathbb{E}[\|g(\nabla)\|^2] \text{ and } X_1 = -2\eta_{\mathcal{W}} \mathbb{E}[(g(\nabla))^2 (\mathbf{w}(\nabla - \mathbf{w}))]. \\ & \text{We now bound } A_1^{(t)} \text{ on the RHS of (17). We have } A_1^{(t)} \\ & A_1^{(t)} \stackrel{(a)}{\leq} -\frac{2\eta_{\mathcal{W}}}{m} \mathbb{E}\Big[\sum_{e \in \mathcal{E}^{(t)}} \sum_{n \in \mathcal{N}_e} \left(\nabla f_n(\mathbf{w}_n^{(t)})\right)^T (\mathbf{w}^{(t)} - \mathbf{w}_n^{(t)})\Big] \end{aligned}$$

ICPP '24, August 12-15, 2024, Gotland, Sweden

$$-\frac{2\eta_{w}}{m}\mathbb{E}\Big[\sum_{e\in\mathcal{E}^{(t)}}\sum_{n\in\mathcal{N}_{e}}\left(\nabla f_{n}(\boldsymbol{w}_{n}^{(t)})\right)^{T}(\boldsymbol{w}_{n}^{(t)}-\boldsymbol{w})\Big]$$

$$\stackrel{(b)}{\leq}\frac{2\eta_{w}}{m}\mathbb{E}\Big[\sum_{e\in\mathcal{E}^{(t)}}\sum_{n\in\mathcal{N}_{e}}f_{n}(\boldsymbol{w}_{n}^{(t)})-f_{n}(\boldsymbol{w}^{(t)})+\frac{L}{2}\|\boldsymbol{w}^{(t)}-\boldsymbol{w}_{n}^{(t)}\|^{2}\Big]$$

$$+\frac{2\eta_{w}}{m}\mathbb{E}\Big[\sum_{e\in\mathcal{E}^{(t)}}\sum_{n\in\mathcal{N}_{e}}f_{n}(\boldsymbol{w})-f_{n}(\boldsymbol{w}_{n}^{(t)})\Big]$$

$$\stackrel{(c)}{=}-\frac{2\eta_{w}}{N_{\mathcal{E}}}\mathbb{E}\Big[F(\boldsymbol{w}^{(t)},\boldsymbol{p}^{(\lfloor\frac{t}{\tau_{1}\tau_{2}}\rfloor)})-F(\boldsymbol{w},\boldsymbol{p}^{(\lfloor\frac{t}{\tau_{1}\tau_{2}}\rfloor)})\Big]+\eta_{w}LA^{(t)}, \quad (18)$$

where (a) follows from the definition of $g^{(t)}$, (b) is from the Lsmoothness and convexity of f_n for any $n \in N$, and (c) holds because $f_e = \frac{1}{N_0} \sum_{n \in N_e} f_n$ and we sample the edge servers based on $\boldsymbol{p}^{(\lfloor \frac{t}{\tau_1 \tau_2} \rfloor)}$.

We then bound $A_2^{(t)}$ on the RHS of (17). We have

$$A_{2}^{(t)} \stackrel{(a)}{\leq} \frac{\eta_{w}^{2}}{m} \sum_{e \in \mathcal{E}^{(t)}} \sum_{n \in \mathcal{N}_{e}} \mathbb{E}\left[\left\| \nabla f_{n}(\boldsymbol{w}_{n}^{(t)}) \right\|^{2} \right] \stackrel{(b)}{\leq} \eta_{w}^{2} G_{w}^{2}, \quad (19)$$

where (a) is from Jensen's inequality and (b) is from Assumption 3. Substituting (18) and (19) into (17), we obtain (16). П

Lemma 4 (Upper Bound for One Update of p). For all $p \in \mathcal{P}$ and $k \in [K]$, we have

$$\mathbb{E}\left[\left\|\boldsymbol{p}^{(k+1)} - \boldsymbol{p}\right\|^{2}\right] - \mathbb{E}\left[\left\|\boldsymbol{p}^{(k)} - \boldsymbol{p}\right\|^{2}\right] \le \eta_{p}^{2}\tau_{1}^{2}\tau_{2}^{2}G_{p}^{2} + \frac{\eta_{p}^{2}\tau_{1}^{2}\tau_{2}^{2}}{m}\sigma_{p}^{2} - 2\eta_{p}\sum_{t=k\tau_{1}\tau_{2}+1}^{(k+1)\tau_{1}\tau_{2}}\mathbb{E}\left[F(\boldsymbol{w}^{(t)}, \boldsymbol{p}) - F(\boldsymbol{w}^{(t)}, \boldsymbol{p}^{(k)})\right].$$
(20)

PROOF. Following similar arguments of Lemma 3, we have

$$\mathbb{E}[\|\boldsymbol{p}^{(k+1)} - \boldsymbol{p}\|^2] \leq \mathbb{E}[\|\boldsymbol{p}^{(k)} - \boldsymbol{p}\|^2] + B_1^{(k)} + B_2^{(k)} + \frac{\eta_p^2 \tau_1^2 \tau_2^2}{m} \sigma_p^2, \quad (21)$$

where we define $B_1^{(k)} = \mathbb{E}[2(\eta_p \bar{u}^{(k)})^T (\boldsymbol{p}^{(k)} - \boldsymbol{p})] \text{ and } B_2^{(k)} = \mathbb{E}[\|\eta_p \bar{u}^{(k)}\|^2].$

 $\mathbb{E}[\|\eta_p \bar{u}^{(k)}\|^2].$ Next, we bound C6. Rearranging the term We now bound the RHS of (21). By substituting $\bar{u}^{(k)} = \sum_{t=k\tau_1\tau_2+1}^{(k+1)\tau_1\tau_2} \bar{v}^{(t)}$ and multiplying both sides by $\frac{N_{\mathcal{E}}}{2\eta_w}$, we have into $B_1^{(k)}$, we have

$$B_{1}^{(k)} \stackrel{(a)}{=} 2\eta_{p} \sum_{t=k\tau_{1}\tau_{2}+1}^{(k+1)\tau_{1}\tau_{2}} \mathbb{E}\left[\left(\nabla_{p}F(\boldsymbol{w}^{(t)},\boldsymbol{p})\right)^{T}(\boldsymbol{p}^{(k)}-\boldsymbol{p})\right]$$
$$\stackrel{(b)}{=} -2\eta_{p} \sum_{t=k\tau_{1}\tau_{2}+1}^{(k+1)\tau_{1}\tau_{2}} \mathbb{E}\left[F(\boldsymbol{w}^{(t)},\boldsymbol{p})-F(\boldsymbol{w}^{(t)},\boldsymbol{p}^{(k)})\right], \quad (22)$$

where (*a*) follows from the definition of $\nabla_{\boldsymbol{p}} F(\boldsymbol{w}, \boldsymbol{p})$ and the linearity of expectation, and (*b*) is from the linearity of F(w, p) in p given any w.

Substituting $\bar{u}^{(k)}$ into $B_2^{(k)},$ we have

$$B_{2}^{(k)} \stackrel{(a)}{\leq} \eta_{p}^{2} \tau_{1} \tau_{2} \sum_{t=k\tau_{1}\tau_{2}+1}^{(k+1)\tau_{1}\tau_{2}} \mathbb{E}\left[\left\|\bar{v}^{(t)}\right\|^{2}\right] \stackrel{(b)}{\leq} \eta_{p}^{2} \tau_{1}^{2} \tau_{2}^{2} G_{p}^{2}, \quad (23)$$

where (a) is from Jensen's inequality and (b) is from Assumption 3. Substituting (22) and (23) into (21), we prove (20).

We now prove Theorem 1.

PROOF. We start from bounding the expected duality gap

$$\max_{\boldsymbol{p}\in\mathcal{P}} \mathbb{E}\left[F(\hat{\boldsymbol{w}}, \boldsymbol{p})\right] - \min_{\boldsymbol{w}\in\mathcal{W}} \mathbb{E}\left[F(\boldsymbol{w}, \hat{\boldsymbol{p}})\right]$$

$$= \max_{\boldsymbol{w}\in\mathcal{W}, \boldsymbol{p}\in\mathcal{P}} \left(\mathbb{E}\left[F(\hat{\boldsymbol{w}}, \boldsymbol{p})\right] - \mathbb{E}\left[F(\boldsymbol{w}, \hat{\boldsymbol{p}})\right]\right)$$

$$\stackrel{(a)}{\leq} \max_{\boldsymbol{w}\in\mathcal{W}, \boldsymbol{p}\in\mathcal{P}} \mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T} \left(F(\boldsymbol{w}^{(t)}, \boldsymbol{p}) - F(\boldsymbol{w}, \boldsymbol{p}^{\left(\lfloor\frac{t}{\tau_{1}\tau_{2}}\rfloor\right)})\right)\right]$$

$$\stackrel{(b)}{=} \max_{\boldsymbol{w}\in\mathcal{W}, \boldsymbol{p}\in\mathcal{P}} (C5 + C6), \qquad (24)$$

where (a) is from the convexity in w and linearity in p and (b) is from the definitions of

$$C5 = \mathbb{E}\Big[\frac{1}{T}\sum_{k=0}^{K-1}\sum_{t=k\tau_{1}\tau_{2}+1}^{(k+1)\tau_{1}\tau_{2}}F(\boldsymbol{w}^{(t)},\boldsymbol{p}) - F(\boldsymbol{w}^{(t)},\boldsymbol{p}^{(k)})\Big]$$
$$C6 = \mathbb{E}\Big[\frac{1}{T}\sum_{t=1}^{T}F(\boldsymbol{w}^{(t)},\boldsymbol{p}^{(\lfloor\frac{t}{\tau_{1}\tau_{2}}\rfloor)}) - F(\boldsymbol{w},\boldsymbol{p}^{(\lfloor\frac{t}{\tau_{1}\tau_{2}}\rfloor)})\Big].$$

We now bound C5. Rearranging the terms of (20) in Lemma 4 and dividing both sides by $\frac{1}{2\eta_p}$, we have

$$\sum_{t=k\tau_{1}\tau_{2}+1}^{(k+1)\tau_{1}\tau_{2}} \mathbb{E}[F(\boldsymbol{w}^{(t)},\boldsymbol{p}) - F(\boldsymbol{w}^{(t)},\boldsymbol{p}^{(k)})]$$

$$\leq \frac{1}{2\eta_{p}} (\mathbb{E}[\|\boldsymbol{p}^{(k)} - \boldsymbol{p}\|^{2}] - \mathbb{E}[\|\boldsymbol{p}^{(k+1)} - \boldsymbol{p}\|^{2}])$$

$$+ \frac{\eta_{p}\tau_{1}^{2}\tau_{2}^{2}}{2}G_{p}^{2} + \frac{\eta_{p}\tau_{1}^{2}\tau_{2}^{2}}{2m}\sigma_{p}^{2}.$$
(25)

Summing (25) over k and dividing both sides by T, we have

$$C5 \leq \frac{1}{2\eta_{p}T} \mathbb{E}[\|\boldsymbol{p}^{(0)} - \boldsymbol{p}\|^{2}] + \frac{\eta_{p}\tau_{1}\tau_{2}}{2}G_{p}^{2} + \frac{\eta_{p}\tau_{1}\tau_{2}}{2m}\sigma_{p}^{2}$$
$$\leq \frac{R_{\mathcal{P}}^{2}}{2\eta_{p}T} + \frac{\eta_{p}\tau_{1}\tau_{2}}{2}G_{p}^{2} + \frac{\eta_{p}\tau_{1}\tau_{2}}{2m}\sigma_{p}^{2}.$$
(26)

Next, we bound C6. Rearranging the terms of (16) in Lemma 3

$$\mathbb{E}\left[F(\boldsymbol{w}^{(t)}, \boldsymbol{p}^{\left(\lfloor\frac{t}{\tau_{1}\tau_{2}}\rfloor\right)}) - F(\boldsymbol{w}, \boldsymbol{p}^{\left(\lfloor\frac{t}{\tau_{1}\tau_{2}}\rfloor\right)})\right] \leq \frac{\eta_{w}N_{\mathcal{E}}}{2}G_{w}^{2} + \frac{\eta_{w}^{2}}{N_{0}}\sigma_{w}^{2} + \frac{LN_{\mathcal{E}}}{2\eta_{w}}A^{(t)} + \frac{N_{\mathcal{E}}}{2\eta_{w}}\mathbb{E}\left[\left\|\boldsymbol{w}^{(t)} - \boldsymbol{w}\right\|^{2}\right] - \frac{N_{\mathcal{E}}}{2\eta_{w}}\mathbb{E}\left[\left\|\boldsymbol{w}^{(t+1)} - \boldsymbol{w}\right\|^{2}\right]. \quad (27)$$

Summing (27) over t and dividing both sides by T, we have

$$C6 \leq \frac{N_{\mathcal{E}}\mathbb{E}[\|w^{(0)} - w\|^{2}]}{2\eta_{w}T} + \frac{LN_{\mathcal{E}}}{2T} \sum_{t=1}^{T} A^{(t)} + \frac{\eta_{w}N_{\mathcal{E}}G_{w}^{2}}{2} + \frac{\eta_{w}\sigma_{w}^{2}}{2N_{0}}$$

$$\stackrel{(a)}{\leq} \frac{N_{\mathcal{E}}R_{W}^{2}}{2\eta_{w}T} + \frac{\eta_{w}N_{\mathcal{E}}}{2}G_{w}^{2} + 10LN_{\mathcal{E}}\eta_{w}^{2}\tau_{1}^{2}(\frac{m+1}{m}\sigma_{w}^{2} + \Psi)$$

$$+ \frac{\eta_{w}}{2N_{0}}\sigma_{w}^{2} + 10LN_{\mathcal{E}}\eta_{w}^{2}\tau_{1}^{2}\tau_{2}^{2}(\frac{m_{\mathcal{E}}+1}{N_{0}}\sigma_{w}^{2} + \Psi), \qquad (28)$$

where (a) is from the bound on $A^{(t)}$ in Lemma 1. Substituting the bound on C5 in (26) and the bound on C6 in (28) into (24), we complete the proof.

Xu et al.