

# FLARE: A New Federated Learning Framework with Adjustable Learning Rates over Resource-Constrained Wireless Networks

Bingnan Xiao, Jingjing Zhang, *Member, IEEE*, Wei Ni, *Fellow, IEEE* and Xin Wang, *Fellow, IEEE*

**Abstract**—Wireless federated learning (WFL) suffers from heterogeneity prevailing in the data distributions, computing powers, and channel conditions of participating devices. This paper presents a new Federated Learning with Adjusted leaRning rate (FLARE) framework to mitigate the impact of the heterogeneity. The key idea is to allow the participating devices to adjust their individual learning rates and local training iterations, adapting to their instantaneous computing powers. The convergence upper bound of FLARE is established rigorously under a general setting with non-convex models in the presence of non-i.i.d. datasets and imbalanced computing powers. By minimizing the upper bound, we further optimize the scheduling of FLARE to exploit the channel heterogeneity. A nested problem structure is revealed to facilitate iteratively allocating the bandwidth with binary search and selecting devices with a new greedy method. A linear problem structure is also identified and a low-complexity linear programming scheduling policy is designed when training models have large Lipschitz constants. Experiments demonstrate that FLARE consistently outperforms the baselines in test accuracy, and converges much faster with the proposed scheduling policy.

**Index Terms**—Federated learning, wireless networks, device scheduling, resource allocation, and convergence analysis.

## I. INTRODUCTION

The upcoming sixth-generation (6G) wireless communication systems will be characterized by pervasive connectivity. The volume of data generated by edge devices—spanning smartphones, wireless sensors, and wearable devices is expected to experience exponential growth [1], [2]. This surge in data production will catalyst widespread adoption and proliferation of artificial intelligence (AI) [3], and give rise to new federated learning (FL) techniques [4]. In a standard FL algorithm, namely FedAvg [5], a fixed number  $\tau$  of local stochastic gradient descent (SGD) updates are performed at a selected set of devices in each training round before the updated parameters are sent to a server, e.g., a base station (BS), for aggregation. This process repeats until an adequate global model is obtained.

As an extension of FL, wireless FL (WFL) supports data processing and model training in wireless networks [6], and has been applied to, e.g., online path control for massive

unmanned aerial vehicle (UAV) networks [7], distributed localization [8], and content recommendations for mobile devices [9]. While FedAvg has demonstrated good empirical performance [10], it requires some desirable conditions as a prerequisite, which is hardly possible in real-world wireless networks. One challenge for directly applying FedAvg in WFL is non-independent and identically distributed (non-i.i.d.) datasets among devices, also known as data heterogeneity. The data heterogeneity tends to induce significant performance degradation [11]. Another challenge is that most FL frameworks only allow the participants to perform the same number of local SGD iterations per communication round at a consistent training speed. However, various edge devices may differ substantially in hardware specifications (e.g., CPU, memory, battery, etc.), resulting in different training speeds and local training iteration numbers [12]. Moreover, the transmission capability, such as transmit power, varies among devices in many wireless systems. FL algorithms, typically developed under the assumption of random uniform device selection, may not be suitable, especially when the system has a limited bandwidth. Devices with poor channel conditions can be prevented from contributing to FL training [13].

### A. Related Work

Recent studies have devoted significant attention to the heterogeneity of FL. As for data heterogeneity, the authors of [14] enhanced the training performance by assigning a globally shared dataset across edge devices. In [15], a variance reduction algorithm called Scaffold was developed with local correction terms appended to mitigate the gradient drift. The authors of [16] utilized contrastive learning at a model level to rectify local updates. Yet, these works overlooked the device heterogeneity by assuming a consistent training speed across devices. Moreover, data sharing [14] and correction terms [15] would increase CPU and memory burden, making it hard to deploy FL at low-power edge devices.

To mitigate device heterogeneity, FedProx was proposed in [17] by adding a proximal term to local training, where the weight of the proximal term needs to be carefully tuned for different training tasks. It was shown in [18] that device heterogeneity could cause objective inconsistency. Then, FedNova was designed to alleviate such adverse impacts to some extent. In [19], FedLin was proposed with correction terms similar to stochastic variance reduction gradient (SVRG); yet, the additional gradient computation and transmission was non-negligible. In general, the above works [18], [19] relied on full

Bingnan Xiao, Jingjing Zhang, and Xin Wang are with the Department of Communication Science and Engineering, School of Information Science and Technology, Fudan University, Shanghai 200433, China (e-mail: 22110720061@m.fudan.edu.cn, {jingjingzhang, xwang11}@fudan.edu.cn).

Wei Ni is with Data61, CSIRO, Marsfield, NSW 2122, Australia, and the School of Computing Science and Engineering, the University of New South Wales (UNSW), Kensington, NSW 2052, Australia (e-mail: Wei.Ni@data61.csiro.au).

or uniformly random device selection, and their convergence under arbitrary or unbalanced device selection cannot be guaranteed. This would hinder their application in wireless edge networks, where devices often need to be scheduled based on their instantaneous or statistical channel conditions, as well as computational and communication abilities.

To tackle channel heterogeneity, efforts have been directed toward achieving efficient aggregation and resource scheduling. Recently, a new analog aggregation technique named over-the-air computation (OTA) has been utilized in FL to alleviate the bandwidth crunch from orthogonal transmission by exploiting the feature of multiple-access channel superposition [20]. In [21], an optimal channel-adaptive power control strategy for OTA-FL was investigated. For digital schemes, the authors of [22] minimized the weighted sum of computation delay and energy consumption in each communication round under full device participation. In [23], a greedy strategy was studied to schedule as many devices as possible per round. Neither of these policies provided provable convergence. In the absence of communication resource constraints, the authors of [24] compared three scheduling policies: random scheduling, round-robin, and proportional fairness. A heuristic scheduling policy was designed in [25] with channel and update importance jointly considered. In [26], a channel- and data-aware scheduling strategy was developed under asynchronous aggregation, where only coarse-grained (e.g., evenly distributed) bandwidth allocation was considered. Joint device scheduling and resource allocation were studied in [27] to improve the trade-off between training rounds and participating devices. By assuming full gradient descent (GD) at the devices, a joint optimization problem for device selection and bandwidth allocation was considered in [28]. By contrast, gradient divergence was used in [29] to define the importance of updates and develop a scheduling strategy based on such importance and channel equality. Considering data and device heterogeneity, the authors of [30] optimized the sampling strategy in the ideal lossless transmission environments.

## B. Contribution and Organization

In this paper, we design a new scheduling framework to address the device, data and channel heterogeneity in WFL, where the learning rates of participating devices can be adjusted online, adapting to the instantaneous availability of their computing power. A new device selection and bandwidth allocation strategy is developed to facilitate the convergence of WFL under the new framework. The main contributions of the paper are summarized as follows.

- We propose a new framework, **Federated Learning with Adjusted learning rate (FLARE)**, to address device heterogeneity in WFL. By adaptively adjusting the learning rates of the participating devices, FLARE allows for consistent training progress among the participating devices with substantially different computing powers, hence accelerating model convergence.
- Under FLARE, we develop a general convergence analysis of WFL with non-convex models in the presence of data and device heterogeneity, and arbitrary device scheduling strategy.

- Given the convergence upper bound, we design a new WFL scheduling policy under the FLARE framework, when the communication bandwidth is limited. To tackle the non-convexity of the convergence bound and the coupling of device selection and bandwidth allocation, we reveal a nested problem structure and decouple the problem to iteratively allocate the bandwidth with binary search and select devices with a new greedy strategy.
- We further reveal a linear problem structure of device selection and bandwidth allocation for models with large Lipschitz constants. A low-complexity linear programming based policy is designed accordingly.

Experiments demonstrate the effectiveness of FLARE in accelerating model convergence under different data distributions and scheduling strategies. Under the FLARE framework, the proposed scheduling policy outperforms the state-of-the-art strategies even without the learning rate adjustments and exhibits robustness under various system parameter settings.

The rest of this paper is organized as follows. Section II presents the system model. Section III elaborates on the FLARE framework tailored to address device heterogeneity, analyzes the communication of FLARE under data heterogeneity, and formulates the problem to facilitate convergence. Section IV reformulates the problem and presents efficient algorithms. Numerical results are provided in Section V, followed by conclusions in Section VI.

*Notation:* Calligraphic letters represent sets. Boldface lowercase letters indicate vectors.  $|\mathcal{A}|$  denotes the cardinality of set  $\mathcal{A}$ , and  $\setminus$  denotes the set difference operation.  $\mathbb{R}$  denotes the real number field.  $\mathbb{E}[\cdot]$  denotes statistical expectation.  $\nabla$  stands for gradient.  $\|\mathbf{w}\|_2$  is the  $\ell_2$ -norm of vector  $\mathbf{w}$ ,  $\langle \mathbf{a}, \mathbf{b} \rangle$  represents the inner product operation, and  $\text{mean}\{\cdot\}$  denotes a mean-value function.

## II. SYSTEM MODEL

In this section, we elucidate the system model, where the participating devices can have unbalanced data, computing, and communication abilities.

### A. Federated Learning Model

The considered WFL system comprises a BS and  $K$  edge devices collected by  $\mathcal{K} \triangleq \{1, 2, \dots, K\}$ . The objective of WFL is to minimize the empirical loss function on a given dataset, as given by

$$F(\mathbf{w}) := \sum_{i=1}^K c_i f_i(\mathbf{w}), \quad (1)$$

where  $\mathbf{w} \in \mathbb{R}^d$  is the  $d$ -dimensional model parameter;  $f_i(\mathbf{w}) \triangleq \mathbb{E}_{\xi_i \sim \mathcal{D}_i} [f_i(\mathbf{w}, \xi_i)]$  represents the local loss function concerning the local dataset  $\mathcal{D}_i$  of device  $i$ ; and  $c_i$  denotes the aggregation weight of device  $i$ , satisfying  $\sum_{i=1}^K c_i = 1$ .

In this paper, we study a generic scenario where heterogeneity prevails in the data, computing and communication conditions of devices. Specifically, the data may not be i.i.d. among the devices; the devices may have substantially different computing powers and hence train their

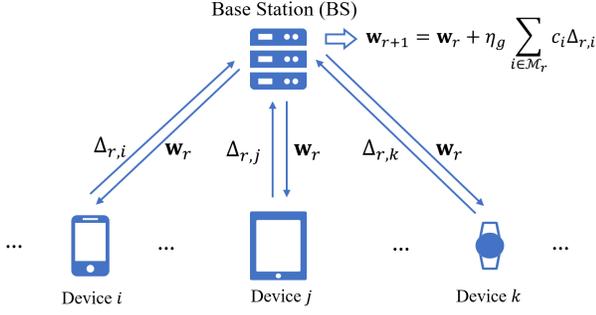


Fig. 1. The architecture of a WFL system with  $M_r$  selected devices at round  $r$ .

local models at different paces; and the communication from the devices to the BS is constrained by limited bandwidth and undergoes non-negligible latency. Define a binary vector  $\mathbf{a} := [a_{r,1}, \dots, a_{r,K}]^\top$ , where  $a_{r,i} \in \{0, 1\}$  indicates whether device  $i$  is scheduled ( $a_{r,i} = 1$ ) or not ( $a_{r,i} = 0$ ) in round  $r$ .

At each training round  $r$ , the following steps are executed, as illustrated in Fig. 1:

- Device scheduling and downlink transmission: The BS selects a subset of devices, denoted by  $\mathcal{M}_r$ , with  $M_r$  devices, and multicasts the global parameter  $\mathbf{w}_r$  to every selected device  $i \in \mathcal{M}_r$ .
- Local update: By setting  $\mathbf{w}_{r,i}^0 = \mathbf{w}_r$ , each selected device  $i$  begins  $\tau_{r,i}$  local updates according to its local computational capacity. The updating rule for SGD is

$$\mathbf{w}_{r,i}^{j+1} = \mathbf{w}_{r,i}^j - \eta_l g_{r,i}^j, \quad j = 0, 1, \dots, \tau_{r,i} - 1, \quad (2)$$

where  $\eta_l$  denotes the local learning rate of the devices, and  $g_{r,i}^j = \nabla f_i(\mathbf{w}_{r,i}^j, \xi_{r,i})$  is the stochastic gradient of device  $i$  with respect to the mini-batch  $\xi_{r,i}$  with size  $D$ , sampled from the local dataset  $\mathcal{D}_i$ .

- Uploading and aggregation: After local computation, each selected device  $i$  uploads its cumulative local gradient  $\Delta_{r,i} = \mathbf{w}_{r,i}^{\tau_{r,i}} - \mathbf{w}_{r,i}^0$ . With the global learning rate  $\eta_g$ , the BS updates the global parameter  $\mathbf{w}_{r+1}$ , as follows.

$$\begin{aligned} \mathbf{w}_{r+1} &= \mathbf{w}_r + \eta_g \sum_{i \in \mathcal{M}_r} c_i \Delta_{r,i} \\ &= \mathbf{w}_r - \eta_g \sum_{i \in \mathcal{M}_r} \frac{1}{M_r} \sum_{j=0}^{\tau_{r,i}-1} \eta_l g_{r,i}^j. \end{aligned} \quad (3)$$

Then, the next  $(r+1)$ -th training round starts. This repeats until the prespecified termination criteria are met.

### B. Local Computation-Communication Latency Model

For each round  $r$ , the total latency  $t_{r,i} = t_{r,i}^{\text{comp}} + t_{r,i}^{\text{comm}}$  undergone by each selected device  $i$  consists of the local computation delay  $t_{r,i}^{\text{comp}}$  and transmission delay  $t_{r,i}^{\text{comm}}$ .

1) *Local Computation*: Let  $f_{r,i}$  represent the computation capacity of device  $i$  in round  $r$ , measured in the number of CPU cycles per second. The computation time needed for the  $\tau_{r,i}$  local training iterations at device  $i$  is given by

$$t_{r,i}^{\text{comp}} = \frac{\tau_{r,i} C_i D}{f_{r,i}}, \quad (4)$$

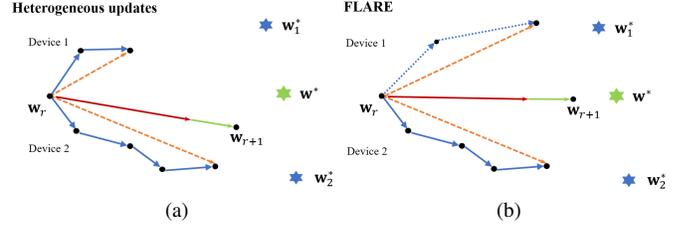


Fig. 2. An illustration of model updates of a heterogeneous setting with (a) equal learning rates and (b) FLARE. The green and blue marks represent the minima of global and local objectives, respectively.

where  $C_i$  (cycles/sample) is the number of CPU cycles required to compute a data sample. Assume that the computation latency needed for model aggregation at the BS is relatively negligible due to its abundant computational capacity and the considerably low complexity of the aggregation process.

2) *Wireless Transmission*: After local computation, the selected devices upload their local updates via frequency-division multiple access (FDMA) with a total bandwidth  $B$ . We consider a line of sight (LoS)-dominating environment, since future wireless systems will predominantly operate in mmWave/THz with quasi-light propagations. For device  $i$ , the achievable uplink rate is given by

$$u_{r,i} = b_{r,i} \log_2 \left( 1 + \frac{p_{r,i} h_{r,i}^2}{b_{r,i} N_0} \right), \quad (5)$$

where  $b_{r,i}$  is the bandwidth allocated to device  $i$ , satisfying  $\sum_{i \in \mathcal{K}} a_{r,i} b_{r,i} \leq B$ ;  $p_{r,i}$  is the transmit power of device  $i$  in round  $r$ ;  $h_{r,i}$  denotes the corresponding channel amplitude gain;  $N_0$  is the power spectral density (PSD) of the additive white Gaussian noise (AWGN) at the BS.

Thus, the transmission latency of device  $i$  is given by

$$t_{r,i}^{\text{comm}} = \frac{S}{u_{r,i}}, \quad (6)$$

where  $S$  is the size of  $\Delta_{r,i}$  (in bits). We assume that the transmissions between the BS and devices are reliable and free of packet errors. Consider synchronous aggregation, then the total latency per round is determined by the slowest of the selected devices. Thus, we require

$$\max_{i \in \mathcal{K}} \{a_{r,i} t_{r,i}\} = \max_{i \in \mathcal{K}} \{a_{r,i} (t_{r,i}^{\text{comp}} + t_{r,i}^{\text{comm}})\} \leq t_{\text{thr}}, \quad (7)$$

where  $t_{\text{thr}}$  denotes the (maximum allowable) overall latency threshold per round.

## III. PROBLEM FORMULATION OF FEDERATED LEARNING WITH DYNAMICALLY ADJUSTED LEARNING RATES

In this section, we propose the FLARE framework, which mitigates the impact of inconsistent local training iterations between devices and expedites the convergence by dynamically adjusting the learning rates of devices involved.

### A. Overview of FLARE

Within the FL architecture, the significant difference in the training loss landscape of different local models, resulting

from the distinct learning capabilities of different devices, can lead to a global objective deviation and hinder the overall convergence [31], [32]. As illustrated in Fig. 2(a), substantially imbalanced updating progresses between different devices (e.g., Device 2 updates its local model much faster than Device 1) are likely to lead to non-negligible deviation of the updated global model  $\mathbf{w}_{r+1}$  from its optimal  $\mathbf{w}^*$ . It is imperative to enhance the consistency of local training progress towards the optimal global parameter  $\mathbf{w}^*$  in a resource-constrained WFL system.

We propose to *dynamically adjust the learning rates of the selected devices based on their computation capabilities* so that the local SGD updates of the devices can progress at consistent rates per round. The local learning rates are linearly scaled to prevent the global objective deviation caused by unbalanced local updates [18]. In other words, resource-constrained devices can conduct a relatively smaller number of local training iterations at faster learning rates. As a result, they can make consistent progress in local training with their more computationally powerful peers, and avoid becoming stragglers, as illustrated in Fig. 2(b).

At any round  $r$ , each device reports its available instantaneous computational resources to the BS. This report incurs negligible communication overhead compared to the local model updates. After a device subset  $\mathcal{M}_r$  is selected, the BS determines  $\bar{\tau}_r$  for the selected devices to adjust their local training iteration numbers. Specifically, the BS sends  $\bar{\tau}_r$  and the global model  $\mathbf{w}_r$  to the selected devices. By setting its learning rate to  $\tilde{\eta}_i = \eta_i \bar{\tau}_r / \tau_{r,i}$ , device  $i \in \mathcal{M}_r$  executes  $\tau_{r,i}$  local SGD iterations per round  $r$ , as given by

$$\begin{aligned} \mathbf{w}_{r,i}^j &= \mathbf{w}_{r,i}^{j-1} - \tilde{\eta}_i g_{r,i}^{j-1} \\ &= \mathbf{w}_{r,i}^{j-1} - \eta_i \frac{\bar{\tau}_r}{\tau_{r,i}} g_{r,i}^{j-1}, \quad j = 0, 1, \dots, \tau_{r,i} - 1. \end{aligned} \quad (8)$$

After local computations, each selected device  $i$  uploads the gradient update  $\Delta_{r,i} = \mathbf{w}_{r,i}^{\tau_{r,i}} - \mathbf{w}_r$  to the BS. The BS then performs the following aggregation:

$$\mathbf{w}_{r+1} = \mathbf{w}_r - \eta_g \sum_{i \in \mathcal{M}_r} \frac{1}{M_r} \sum_{j=0}^{\tau_{r,i}-1} \eta_i \frac{\bar{\tau}_r}{\tau_{r,i}} g_{r,i}^{j-1}. \quad (9)$$

**Remark 1:** FLARE can be viewed as a generalization version of FedAvg by mitigating the imbalance in local SGD iterations. It reduces to FedAvg if  $\mathcal{M}_r$  is randomly uniformly selected with  $\tau_{r,1} = \dots = \tau_{r,K}$  and  $\eta_g = 1$  at each round  $r$ .

### B. Approximation Error Analysis

We analyze the approximation errors caused by the proposed learning rate adjustment in FLARE, with respect to performing a consistent number of local training iterations (e.g.,  $\bar{\tau}_r$  iterations) across all devices. We establish the following Theorem.

**Theorem 1.** *Suppose that  $\mathbf{w}$  is obtained by taking  $\bar{\tau}_r$  SGD steps with a learning rate  $\eta_l$ , and  $\tilde{\mathbf{w}}$  is obtained by taking  $\tau_r$  SGD steps with a learning rate  $\tilde{\eta}_l = \eta_l \bar{\tau}_r / \tau_r$ . Starting from the same initial point with a non-convex function  $f$ , the mean squared error (MSE) of  $\tilde{\mathbf{w}}$  with respect to  $\mathbf{w}$  is bounded by*

$$\mathbb{E} \left[ \|\tilde{\mathbf{w}} - \mathbf{w}\|_2^2 \right] \leq 8\eta_l^2 \bar{\tau}_r^2 (g + \sigma^2), \quad (10)$$

where  $\sigma^2$  and  $g$  denote the SGD variance and gradient upper bound, respectively.

*Proof.* See Appendix A.  $\square$

From **Theorem 1**, the RHS of (10) depends only on the constants,  $g$  and  $\sigma^2$ , when the local learning rate  $\eta_l$  is inversely proportional to  $\bar{\tau}_r$ , i.e.,  $\eta_l = \mathcal{O}\left(\frac{1}{\bar{\tau}_r}\right)$ . As a consequence, the MSE,  $\mathbb{E} \left[ \|\tilde{\mathbf{w}} - \mathbf{w}\|_2^2 \right]$ , is bounded under FLARE. Also, such approximation errors decrease with the progress of model training since the corresponding gradient paradigm becomes smaller, leading to a tighter gradient upper bound  $g$ . In Section III-C, we demonstrate that the relationship between  $\eta_l$  and  $\bar{\tau}_r$  plays a critical role in reducing the convergence upper bound.

### C. WFL Problem Formulation

For WFL with constrained computing and communication resources, the objective is to minimize the training loss. A general formulation of such a problem is given by

$$\mathbf{P1}: \quad \min_{\mathbf{a}, \mathbf{b}} F(\mathbf{w}) \quad (11a)$$

$$\text{s.t.} \quad a_{r,i} \in \{0, 1\}, \quad (11b)$$

$$\sum_{i \in \mathcal{K}} a_{r,i} b_{r,i} \leq B, \quad (11c)$$

$$a_{r,i} t_{r,i} \leq t_{\text{thr}}, \quad \forall i \in \mathcal{K}, \forall r \in [R]. \quad (11d)$$

where  $\mathbf{b} = [b_{r,1}, \dots, b_{r,K}]^\top$  collects the bandwidth allocations of the  $K$  devices in round  $r$ . Constraint (11c) specifies that the total bandwidth  $B$  must not be exceeded; and (11d) dictates that the computation and transmission must be completed before a proposed latency threshold  $t_{\text{thr}}$  in each round.

During the training process, the BS dynamically determines a subset of devices and allocates bandwidths to achieve the minimal global loss while satisfying constraints (11c) and (11d). In order to analyze how the convergence of FLARE is affected by device selection and resource allocation, we next establish an asymptotic convergence upper bound for  $F(\mathbf{w})$  and minimize the derived upper bound to seek an efficient solution for (11).

### D. Convergence Analysis

We analyze the convergence of FLARE, which contributes to the improved tractability of Problem **P1**, as will be described in Section IV. This starts with the following assumptions.

**Assumption 1.** (Unbiasedness and Variance Boundedness). For a model parameter  $\mathbf{w}$ , the stochastic gradient of each device is an unbiased estimator of the true local gradient, i.e.,

$$\mathbb{E}[\nabla f_i(\mathbf{w}; \xi)] = \nabla f_i(\mathbf{w}), \quad \forall i \in \mathcal{K}. \quad (12)$$

Moreover, there exists a bounded SGD variance  $\sigma^2$  such that

$$\mathbb{E} \left[ \|\nabla f_i(\mathbf{w}, \xi) - \nabla f_i(\mathbf{w})\|^2 \right] \leq \sigma^2, \quad \forall i \in \mathcal{K}. \quad (13)$$

**Assumption 2.** (Smoothness) The local objective function  $f_i$  is Lipschitz smooth; i.e., there exists a constant  $L > 0$  satisfying

$$\|\nabla f_i(\mathbf{w}_2) - \nabla f_i(\mathbf{w}_1)\| \leq L \|\mathbf{w}_2 - \mathbf{w}_1\|, \quad \forall i \in \mathcal{K}. \quad (14)$$

**Assumption 3.** (Gradient Dissimilarity Boundedness) There exist two constants  $G \geq 0$  and  $H \geq 1$  such that

$$\|\nabla f_i(\mathbf{w})\|^2 \leq G^2 + H^2 \|\nabla F(\mathbf{w})\|^2, \quad \forall i \in \mathcal{K}. \quad (15)$$

**Assumptions 1 – 3** have been extensively considered in the literature, e.g., [15], [18], [19]. Under these assumptions, we analyze the one-round convergence of FLARE as follows.

**Lemma 1.** For each round  $r$  and the subset of selected devices  $\mathcal{M}_r$ , the upper bound of the one-round convergence of FLARE with a non-convex loss function is given by

$$\mathbb{E}[F(\mathbf{w}_{r+1}) - F(\mathbf{w}_r)] \leq -\frac{\bar{\tau}_r \eta_g \eta_l q_r}{2} \|\nabla F(\mathbf{w}_r)\|^2 + \phi_1 + \phi_2, \quad (16)$$

where  $\phi_1 = 1.2L^2\bar{\tau}_r^3\eta_g\eta_l^3G^2$ ,  $\phi_2 = \frac{L\bar{\tau}_r^2\eta_g\eta_l^2}{2} \left( \frac{L\bar{\tau}_r\eta_l}{M_r} + \frac{\eta_g}{M_r^2} \right) \sigma^2 \sum_{i \in \mathcal{M}_r} \frac{1}{\tau_{r,i}}$ , and  $q_r$  is a round-specific positive constant.

*Proof.* See Appendix B.  $\square$

**Lemma 1** reveals that the convergence bound of any round is affected by  $\phi_1$  and  $\phi_2$ . Here,  $\phi_1$  is determined by  $G^2$ , which reflects the degree of inconsistency between the local and global gradients, as stated in **Assumption 3**. By contrast,  $\phi_2$  can be further controlled as it depends heavily on the numbers of selected devices  $M_r$  and local training iterations  $\tau_{r,i}$ . We note that  $\phi_2$  can be reduced by selecting an adequate  $M_r$ , and scheduling each device to perform a sufficient number  $\tau_{r,i}$  of local updates in round  $r$ . This approach is more effective in decreasing the upper bound  $\mathbb{E}[F(\mathbf{w}_{r+1}) - F(\mathbf{w}_r)]$  than only maximizing the number of selected devices (as done in, e.g., [23], [33]).

Given the one-round convergence in **Lemma 1**, we elucidate the convergence performance of FLARE in **Theorem 2**. For illustration convenience, we assume  $\bar{\tau}_r = \bar{\tau}$ ,  $\forall r = 1, \dots, R$  in the theorem. (Nevertheless, it is straightforward to extend **Theorem 2** to the situation where  $\bar{\tau}_r$ ,  $\forall r$  changes over rounds, i.e., by applying additional round-specific scaling to (35) before summing (34) up from  $r = 1$  to  $R$  in Appendix C.)

**Theorem 2.** Given an initial global model  $\mathbf{w}_1$  and the optimal global model  $\mathbf{w}^*$ , the convergence upper bound of FLARE with a total of  $R$  training rounds is given by

$$\frac{1}{R} \sum_{r=1}^R \mathbb{E} [\|\nabla F(\mathbf{w}_r)\|^2] \leq \frac{2[F(\mathbf{w}_1) - F(\mathbf{w}^*)]}{\bar{\tau}\eta_g\eta_l q R} + \Phi_1 + \Phi_2, \quad (17)$$

where  $\Phi_1 = \frac{2.4}{q} L^2 \bar{\tau}^2 \eta_l^2 G^2$ ,  $\Phi_2 = \frac{L\eta_l}{q} (L\bar{\tau}\eta_l + \frac{\kappa\eta_g}{M}) \sigma^2$ ,  $\kappa = \max_{r,i} \{\frac{\bar{\tau}}{\tau_{r,i}}\}$ ,  $M \triangleq \min_r \{M_r\}$ , and  $q$  is a positive constant.

*Proof.* See Appendix C.  $\square$

Compared with the existing studies [14]–[16], the convergence analysis of FLARE in **Theorem 2** is generic in the sense that it supports both non-i.i.d. data and imbalanced local training, i.e., data and device heterogeneity with non-convex training models.

- By observing the RHS of (17) in comparison with (36), we can reveal the impact of device selection on the convergence of FLARE: In any round  $r$ , a larger  $M_r$  leads

to a tighter upper bound in (36) since  $M \triangleq \min_r \{M_r\}$ , where  $M$  denotes the minimum number of selected devices in  $R$  rounds. Moreover,  $\frac{\kappa L \eta_l \eta_g}{q M}$  in  $\Phi_2$  decreases linearly with the number  $M$  of selected devices. Under full participation at each round, i.e.,  $M = K$ , the convergence gap is the minimum.

- The parameter  $\kappa$  in  $\Phi_2$  represents the maximal difference of local training iterations relative to  $\bar{\tau}$ , which illustrates the convergence error introduced by the imbalanced local SGD iterations. It would decrease if the selected devices perform a consistent number of local iterations, i.e.,  $\tau_{r,i}$  is equal,  $\forall i$ . The corresponding error  $\frac{L\eta_l\kappa\eta_g}{qM}\sigma^2$  in  $\Phi_2$  can be reduced if  $M$  is larger.
- As shown in (17),  $\Phi_1$  and  $\Phi_2$  do not vanish with the increase of training rounds  $R$ . They have similar structures and grow with  $G^2$  and  $\sigma^2$ . To effectively bound  $\bar{\tau}^2\eta_l^2$  in  $\Phi_1$  and  $\bar{\tau}\eta_l^2$  in  $\Phi_2$ , we design  $\eta_l = \mathcal{O}(\frac{1}{\bar{\tau}})$ . As a result, both  $\Phi_1$  and  $\Phi_2$  are upper bounded in the sense that they do not grow with  $\bar{\tau}$ . In turn, the approximation error in (10) is upper bounded, as described in **Theorem 1**.

We also establish the convergence rate of FLARE through an appropriate configuration of global and local learning rates, as delineated in the following corollary.

**Corollary 1.** Let  $\eta_l = \frac{1}{\bar{\tau}\sqrt{R}}$  and  $\eta_g = \sqrt{\bar{\tau}M}$ . The convergence rate of FLARE is given by

$$\frac{1}{R} \sum_{r=1}^R \mathbb{E} [\|\nabla F(\mathbf{w}_r)\|^2] \leq \mathcal{O} \left( \frac{1}{\sqrt{\bar{\tau}MR}} + \frac{1}{R} \right). \quad (18)$$

*Proof.* By plugging  $\eta_l = \frac{1}{\bar{\tau}\sqrt{R}}$  and  $\eta_g = \sqrt{\bar{\tau}M}$  into (17) and suppressing the low-order terms, (18) is obtained.  $\square$

**Remark 2:** It is indicated in **Corollary 1** that FLARE with the local and global learning rates of  $\eta_l = \frac{1}{\bar{\tau}\sqrt{R}}$  and  $\eta_g = \sqrt{\bar{\tau}M}$ , can achieve a linear convergence speedup with the number  $M$  of devices<sup>1</sup>. The convergence rate yields  $\mathcal{O}(\frac{1}{\sqrt{\bar{\tau}MR}})$  when  $R > \sqrt{\bar{\tau}M}$ . To the best of our knowledge, this is the first work to achieve such a convergence rate under non-i.i.d. data, imbalanced local updates (e.g., in terms of local iteration number per round), and arbitrary device selections in a general non-convex case, thereby validating the effectiveness of FLARE.

**Remark 3:** Compared to the existing methods, e.g., [11], [15], [34], FLARE is advantageous in both convergence rate and communication overhead. In non-convex cases, the convergence rate of FLARE, i.e.,  $\mathcal{O}(\frac{1}{\sqrt{\bar{\tau}MR}})$ , surpasses those of prior works, e.g.,  $\mathcal{O}(\frac{\bar{\tau}}{\sqrt{MR}})$  in [11], [34] by properly designing the learning rates. The convergence rate of FLARE increases with the number  $\bar{\tau}$  of SGD iterations per round. In variance reduction methods, e.g., Scaffold [15], devices perform local updates with correction terms of the same

<sup>1</sup>To attain the accuracy of  $\epsilon$ ,  $\mathcal{O}(\frac{1}{\epsilon^2})$  steps are needed with a convergence rate  $\mathcal{O}(\frac{1}{\sqrt{R}})$ , while  $\mathcal{O}(\frac{1}{M\epsilon^2})$  steps are needed if the convergence rate is  $\mathcal{O}(\frac{1}{\sqrt{MR}})$ . In this sense, we achieve a linear speedup with increasingly selected devices.

dimension as the model parameter  $\mathbf{w}$ . The BS aggregates these terms to update the global model. FLARE achieves a comparable convergence rate, while device  $i$  only needs to report  $\tau_{r,i}$  to the BS, which is relatively negligible compared to the correction terms used in variance reduction methods.

#### IV. PROPOSED SCHEDULING UNDER FLARE

In light of (16) and (17), we optimize bandwidth allocation and device selection by minimizing the convergence upper bound. We first develop an iterative algorithm that decouples bandwidth allocation and device selection to alleviate communication bottlenecks and accelerate convergence. We further reveal a convex structure of Problem **P2** under certain scenarios and propose an efficient solution with linear programming.

##### A. Problem Reformulation

Minimizing the one-round convergence upper bound in (16) contributes to decreased convergence upper bound in **Theorem 2**, i.e., (17). Consequently, the minimization of  $F(\mathbf{w})$  can be approximated by minimizing the RHS of (16), which can be efficiently achieved by minimizing  $\phi_2$ . Problem **P1** can be rewritten as

$$\mathbf{P2}: \min_{\mathbf{a}, \mathbf{b}} \left( \frac{1}{M_r} + \frac{\gamma}{M_r^2} \right) \sum_{i \in \mathcal{M}_r} \frac{1}{\tau_{r,i}}, \quad (19a)$$

$$\text{s.t. } a_{r,i} \in \{0, 1\}, \quad (19b)$$

$$\sum_{i \in \mathcal{K}} a_{r,i} b_{r,i} \leq B, \quad (19c)$$

$$a_{r,i} t_{r,i} \leq t_{\text{thr}}, \forall i \in \mathcal{K}, \quad \forall r = 1, \dots, R. \quad (19d)$$

where we define  $\gamma = \frac{\eta_{\text{E}}}{L}$  since  $\eta_{\text{I}} \sim \mathcal{O}(\frac{1}{r})$ , according to **Theorem 1**. Problem **P2** implies that the BS should optimally specify the number of scheduled devices and their respective local computation and transmission capabilities to minimize (19a), while adhering to the bandwidth and delay constraints (19c) and (19d). Problem **P2** is a mixed integer nonlinear programming (MINLP) problem.

##### B. Joint Bandwidth Allocation and Device Participation

According to (19d), selecting devices with smaller delays helps admit more devices into an FL training round (i.e., increasing  $\mathcal{M}_r$ ) and, in turn, decrease the objective function in (19a). We can view the MINLP Problem **P2** as a nested problem by treating  $t_{r,i}, \forall i, r$  as a function of  $\mathbf{a}$  and  $\mathbf{b}$ , albeit  $t_{r,i}$  does not appear explicitly in (19a). This is because minimizing  $t_{r,i}$  can decrease (19a) by expanding the feasible area of Problem **P2**, i.e., with (19d) satisfied. For this reason, we decompose **P2** into two manageable sub-problems. We solve the optimal bandwidth allocation strategy  $\mathbf{b}^*$  for any given (feasible) select of selected devices  $\mathcal{M}_r$  to minimize the total latency in round  $r$ . Then, we develop a greedy method to find  $\mathcal{M}_r$  that minimizes (19a). These two steps iterate till convergence, achieving the overall policy.

1) *Bandwidth Allocation*: Given a device subset  $\mathcal{M}_r$ , we ensure the total latency of the selected devices satisfies (19d) by carefully allocating the available bandwidth  $B$  to the devices, reducing the overall latency per round. The bandwidth allocation is cast as the following min-max problem:

$$\mathbf{P3}: \min_{\mathbf{b}} \max_i t_{r,i} \quad (20a)$$

$$\text{s.t. } \sum_{i \in \mathcal{M}_r} b_{r,i} \leq B, \quad (20b)$$

$$b_{r,i} \geq 0, \forall i \in \mathcal{M}_r, \quad \forall r = 1, \dots, R, \quad (20c)$$

where (19c) is simplified to (20b) as  $a_{r,i} = 1, \forall i \in \mathcal{M}_r$ .

Let  $t_r^* = \min_{\mathbf{b}} \max_i t_{r,i}$  denote the optimal value of (20a). Given  $\mathcal{M}_r$ , the optimal bandwidth allocation  $\mathbf{b}^*$  and  $t_r^*$  have the following relationship.

**Theorem 3.** *The optimal bandwidth allocation of Problem **P3** satisfies*

$$b_{r,i}^* = \frac{S \ln 2}{(t_r^* - t_{r,i}^{\text{comp}})(v_{r,i} + W(-v_{r,i} e^{-v_{r,i}}))}, \quad (21)$$

where  $v_{r,i} = \frac{SN_0 \ln 2}{p_{r,i} h_{r,i}^2 (t_r^* - t_{r,i}^{\text{comp}})}$ ,  $W(\cdot)$  denotes the Lambert-W function, and  $t_r^*$  satisfies

$$\sum_{i \in \mathcal{M}_r} b_{r,i}^* = B. \quad (22)$$

*Proof.* See Appendix D.  $\square$

Based on (21) and (22), a binary search method can be employed to obtain  $t_r^*$  numerically and  $b_{r,i}^*, \forall i$  subsequently.

2) *Device Selection*: Given the fixed bandwidth allocation strategy  $\mathbf{b}^*$ , device selection can be written as

$$\mathbf{P4}: \min_{\mathbf{a}} \left( \frac{1}{M_r} + \frac{\gamma}{M_r^2} \right) \sum_{i \in \mathcal{M}_r} \frac{1}{\tau_{r,i}} \quad (23a)$$

$$\text{s.t. } t_r^* \leq t_{\text{thr}}, \forall \mathcal{M}_r \subset \mathcal{K}, \quad \forall r = 1, \dots, R. \quad (23b)$$

Even with the fixed  $\mathbf{b}^*$ , Problem **P4** is an MINLP. We develop an iterative greedy algorithm to solve Problem **P4** efficiently. In the  $k$ -th iteration of the greedy algorithm, we first specify an ‘‘available’’ subset  $\mathcal{M}_r$  of candidate devices that decreases the objective of (23a). Then, we find the specific device that minimizes the increase of the total latency  $t_r$  in the available subset solving Problem **P3**.

To determine the available subset of devices, denoted by  $\Pi_k$ , for the BS in the  $k$ -th iteration of device selection, we come up with the following Theorem.

**Theorem 4.** *Assuming that a device subset  $\mathcal{Q}_k$  with  $Q$  devices has been selected for up to the  $k$ -th iteration, the BS shall continue to select devices from the set that satisfies the following condition:*

$$\frac{1}{\tau_{r,i}} < \frac{Q^2 + (2\gamma + 1)Q + \gamma}{Q^2(Q + \gamma + 1)} \sum_{q \in \mathcal{Q}_r} \frac{1}{\tau_{r,q}}, \quad \exists i \in \mathcal{N}, \quad (24)$$

where  $\mathcal{N} = \mathcal{K} \setminus \mathcal{Q}_k$ .

*Proof.* See Appendix E.  $\square$

**Algorithm 1** Iterative Scheduling Algorithm

---

```

1: Initialize  $\mathcal{M}_r \leftarrow \emptyset$ ,  $\mathcal{Q}_1 \leftarrow \emptyset$ 
2: Initial determination:  $x \leftarrow \arg \min_{i \in \mathcal{K}} \frac{1+\gamma}{\tau_{r,i}}$ ,  $\mathcal{Q}_1 \leftarrow \mathcal{Q}_1 \cup \{x\}$ 

3: for  $k = 1, 2, \dots$  do
4:    $\mathcal{N}_k \leftarrow \mathcal{K} \setminus \mathcal{Q}_k$ 
5:   For  $\mathcal{N}_k$ , determine the available subset  $\Pi_k$  with (24)
6:   if  $\Pi_k = \emptyset$  then
7:     break
8:   end if
9:   for each device  $y \in \Pi_k$  do
10:    With  $\mathcal{Q}_k \cup \{y\}$ , calculate  $t_{r,y}^*$  with a binary search
11:   end for
12:   if  $\min_y t_{r,y}^* \leq t_{\text{thr}}$  then
13:      $z \leftarrow \arg \min_{y \in \Pi_k} t_{r,y}^*$ ,  $\mathcal{Q}_k \leftarrow \mathcal{Q}_k \cup \{z\}$ 
14:   else if  $\min_y t_{r,y}^* > t_{\text{thr}}$  then
15:     break
16:   end if
17: end for
18: return  $\mathcal{M}_r \leftarrow \mathcal{Q}_k$ 

```

---

**Theorem 4** ensures that the selected device at each device selection iteration decreases the value of the objective function in (23a). By employing (24), the BS does not need to evaluate all unselected devices at every step, thereby improving search efficiency. The size of the subset  $\Pi_k$  also decreases gradually. This iterative device selection terminates when a stopping condition is met, i.e.,  $|\Pi_k| = 0$ .

For each device in the subset  $\Pi_k$ , a binary search method is employed for solving Problem **P3**, and the device  $z$  that incurs the smallest increase in the total latency  $t_{r,z}^*$  is selected to be included in the set  $\mathcal{Q}_k$ , as described in Section IV-B(1). This iterative operation is terminated when no device in  $\Pi_k$  satisfies (24), or the inclusion of any other devices in  $\mathcal{Q}_k$  would violate the latency threshold  $t_{\text{thr}}$ . This proposed scheduling policy is summarized in Algorithm 1, where **Theorem 4** is leveraged to ensure the monotonic decrease of (23a), until  $|\Pi_k| = 0$  or constraint (23b) is satisfied. A suboptimal solution is then attained.

The binary search method adopted in Alg. 1 is known to incur a computational complexity of  $\mathcal{O}(|\mathcal{Q}_k| \log_2 \frac{t_{\text{thr}}}{\epsilon})$  to achieve  $\epsilon$ -accuracy. Since the number of device selection iterations must not exceed  $K$  and no more than  $K$  binary searches are performed at each device selection iteration, the worst-case complexity of Alg. 1 is  $\mathcal{O}(K^2 |\mathcal{Q}_k| \log_2 \frac{t_{\text{thr}}}{\epsilon})$ , which is no more than  $\mathcal{O}(K^3 \log_2 \frac{t_{\text{thr}}}{\epsilon})$ . Compared to the exhaustive search with a complexity of  $\mathcal{O}(2^K |\mathcal{Q}_k| \log_2 \frac{t_{\text{thr}}}{\epsilon})$ , Alg. 1 effectively reduces the computational complexity, especially in the presence of a large number of devices.

### C. Linear Programming Solution with Special Structures

According to (19), as  $\gamma = \frac{\eta_g}{L} \rightarrow 0$ , especially when the Lipschitz gradient constant  $L$  is large due to a deeper neural network structure and  $\eta_g$  is relatively small, (19a) is readily approximated by  $\frac{1}{M_r} \sum_{i \in \mathcal{M}_r} \frac{1}{\tau_{r,i}}$  since  $\gamma \ll M_r$ . In this case,

**Algorithm 2** The Proposed Scheduling Algorithm for WFL

---

```

1: Initialize  $\mathbf{w}_1$ ,  $t_{\text{thr}}$ ,  $\eta_g$ ,  $\eta_l$  and  $\gamma$ 
2: for  $r = 1, 2, \dots, R$  do
3:   Each device  $i \in \mathcal{K}$  sends  $f_{r,i}$ ,  $\tau_{r,i}$ ,  $p_{r,i}$ ,  $h_{r,i}$  to the BS
4:   The BS determines the value of  $\bar{\tau}_r$ .
5:   The BS derives the scheduled subset  $\mathcal{M}_r$  with Alg. 1 or linear programming ( $\gamma \rightarrow 0$ ) in Section IV-C.
6:   The BS broadcasts the global model  $\mathbf{w}_r$ ,  $\bar{\tau}_r$  and the allocated bandwidth  $b_{r,i}^*$  to the selected devices.
7:   for each scheduled device  $i \in \mathcal{M}_r$  in parallel do
8:     Make the learning rate adjustment  $\tilde{\eta}_l = \frac{\bar{\tau}_r}{\tau_{r,i}} \eta_l$  with the received  $\bar{\tau}_r$ 
9:     for  $j = 0, \dots, \tau_{r,i} - 1$  do
10:      Perform local SGD update according to (8)
11:     end for
12:     Send  $\tilde{\Delta}_{r,i} = -\sum_{j=0}^{\tau_{r,i}-1} \tilde{\eta}_l g_{r,i}^j$  back to the BS
13:   end for
14:   BS receives  $\tilde{\Delta}_r = \frac{1}{M} \sum_{i \in \mathcal{M}_r} \tilde{\Delta}_{r,i}$  and update the global model with  $\mathbf{w}_{r+1} = \mathbf{w}_r + \eta_g \tilde{\Delta}_r$ 
15: end for
16: return  $\mathbf{w}_{R+1}$ 

```

---

we define  $\boldsymbol{\nu} = \left[ \frac{1}{\tau_{r,1}}, \dots, \frac{1}{\tau_{r,K}} \right]^\top$ ,  $\theta = \frac{1}{\mathbf{1}^\top \mathbf{a}}$ ,  $\boldsymbol{\alpha} = \theta \mathbf{a}$ , and  $\boldsymbol{\beta} = \theta \mathbf{b}$ . Problem **P2** can be rewritten in a linear form:

$$\mathbf{P5} : \min_{\boldsymbol{\alpha}, \boldsymbol{\beta}, \theta} \boldsymbol{\nu}^\top \boldsymbol{\alpha} \quad (25a)$$

$$\text{s.t. } \alpha_{r,i} c_{r,i} \leq \beta_{r,i} \leq \alpha_{r,i} B, \quad (25b)$$

$$\sum_{i \in \mathcal{K}} \beta_{r,i} \leq \theta B, \quad \sum_{i \in \mathcal{K}} \alpha_{r,i} = 1, \quad (25c)$$

$$0 \leq \alpha_{r,i} \leq 1, 0 < \theta \leq 1, \forall i \in \mathcal{K}, \forall r = 1, \dots, R, \quad (25d)$$

Here the term  $c_{r,i}$  shares the same form as  $b_{r,i}^*$  in (21), with  $t_r^* = t_{\text{thr}}$  substituted. To this end, the minimum bandwidth is allocated when any device  $i$  is selected.

It is observed that Problem **P5** is a convex linear programming problem. This convex linear program can be effectively solved using, e.g., the Matlab linprog function. Once the optimal  $\boldsymbol{\alpha}^*$ ,  $\boldsymbol{\beta}^*$  and  $\theta^*$  are obtained, the optima of  $\mathbf{a}^*$  and  $\mathbf{b}^*$  can be approximately identified as  $\mathbf{a}^* = \text{round} \left( \frac{\boldsymbol{\alpha}^*}{\theta^*} \right)$  and  $\mathbf{b}^* = \text{round} \left( \frac{\boldsymbol{\beta}^*}{\theta^*} \right)$ , where  $\text{round}(\cdot)$  stands for rounding. This linear programming solution complements Alg. 1 to achieve a fast approximate solution with polynomial complexity when  $\gamma$  is small.

### D. Overall Policy

We describe the proposed FLARE framework with the joint device selection and bandwidth allocation strategy for WFL, as outlined in Algorithm 2. Specifically, each round of FL training starts with all devices sending their local computation and channel information to the BS. After the BS determines the value of  $\bar{\tau}_r$ , it calls Alg. 1 or its linear programming-based low-complexity alternative to decide the selected device subset  $\mathcal{M}_r$ , and sends the global model to the selected devices in  $\mathcal{M}_r$ . Then, the selected devices adjust their

local learning rate  $\eta_l$  based on  $\bar{\tau}_r$ . After local computation, with the allocated bandwidths, the selected devices transfer their updated accumulated gradients to the BS for global model aggregation with the global learning rate  $\eta_g$ .

## V. EXPERIMENT RESULTS

In this section, we empirically evaluate the proposed FLARE framework and scheduling strategy for WFL.

### A. Experiment Setup

We carry out experiments on the MNIST [35] and CIFAR-10 [36] datasets. In the case of MNIST, we employ a CNN for handwritten digit recognition. In the case of CIFAR-10, we use another CNN for image classification. We consider both i.i.d. and non-i.i.d. data distributions for both datasets. In the non-i.i.d. case, the training data samples are sorted by label and distributed randomly among the devices. Unless specified otherwise, we set the global learning rate  $\eta_g = 1$ , and the local learning rate is set to  $\eta_l = 0.005$  for MNIST and  $\eta_l = 0.01$  for CIFAR-10. The batch size is  $D = 40$ .

Consider a WFL edge network with  $K = 40$  devices uniformly distributed between 100 m and 500 m away from the BS. The path-loss exponent is 3.76. The total bandwidth  $B$  is 10 MHz. The maximum transmit power of a device is  $p_{i,\max} = 20$  dBm. The power spectrum density of the AWGN is  $N_0 = -114$  dBm/MHz at the BS. The model size  $S$  is  $1 \times 10^7$  or  $6.4 \times 10^7$  bits for MNIST and CIFAR-10, respectively. The CPU frequency of a device follows a uniform distribution between 2 GHz and 4 GHz. According to our experimental tests, the required number of CPU cycles for computing a data sample is specified empirically to be 110 cycles/bit for MNIST and 85 cycles/bit for CIFAR-10. To simulate device heterogeneity, we assume that the number of local updates  $\tau_{r,i}$  of each device  $i$  follows an exponential distribution with the mean value  $\tau$  [37]. All experiments are conducted on a computer with an Intel 13700K processor and 64 GB memory, running Python 3.9, Numpy 1.23.3, and PyTorch 1.13.0, installed on a Windows 10 operating system.

### B. Evaluation of FLARE

Fig. 3 demonstrates the benefit of the local learning rate adjustments in FLARE in the presence of device and data heterogeneity, where we assume that the wireless bandwidth is sufficient. The following baselines are considered.

- Fixed Aggregation: The same number of local training iterations with persistent learning rates is performed at the selected devices per round.
- Flexible Aggregation: Different numbers of local training iterations with persistent learning rates can be performed at the selected devices per round.

We see that FLARE dramatically outperforms the other two settings, attributed to the adaptively adjusted local learning rates of the selected devices.

Under FLARE, we further consider the following strategies that the BS can take to determine  $\bar{\tau}_r$  in each round  $r$ .

- Max strategy (MaS):  $\bar{\tau}_r = \max_{i \in \mathcal{M}_r} \{\tau_{r,i}\}$  for  $r \geq 1$ ;

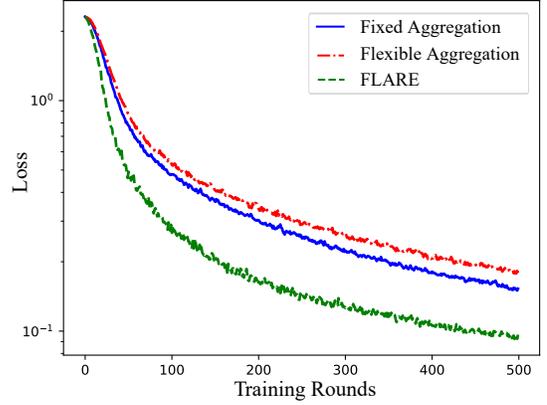


Fig. 3. Comparison of training loss on non-i.i.d. MNIST dataset with  $K = 60$ ,  $M_r = 20$ , and uniform sampling. For the fixed aggregation, we have  $\tau_{r,i} = 7, \forall r \in \mathcal{M}_r$ . For the flexible aggregation, every 20 devices among the  $K$  devices perform 12, 6, and 3 local updates, respectively. For FLARE,  $\bar{\tau}_r$  is set as the maximal local updates among the selected devices in each round.

- Mean strategy (MeS):  $\bar{\tau}_r = \text{mean}_{i \in \mathcal{M}_r} \{\tau_{r,i}\}$  for  $r \geq 1$ ;
- Fixed max strategy (FMaS):  $\bar{\tau}_r = \max_{i \in \mathcal{M}_r} \{\tau_{1,i}\}$  for  $r > 1$ ;
- Fixed mean strategy (FMeS):  $\bar{\tau}_r = \text{mean}_{i \in \mathcal{M}_r} \{\tau_{1,i}\}$  for  $r > 1$ .

Here, FMaS and FMeS suggest that the BS directly utilizes  $\{\bar{\tau}_{1,i}\}$  from the first training round in various ways.

Fig. 4 shows the convergence of FLARE with  $\tau = 3$ . A random, uniform selection of 10 out of  $K = 40$  devices and a non-uniform device selection are considered. As for the non-uniform device selection, we divide all devices into four groups before training. The probabilities of devices being selected in the four groups are 0.05, 0.15, 0.2, and 0.6 per round. MaS and MeS outperform FedAvg upon convergence, by up to 4.5% and 9.4% under the i.i.d. and non-i.i.d. settings, respectively. MaS converges faster than MeS because  $\bar{\tau}_r$  is larger in most rounds and hence the test accuracy is better under MaS. In contrast, MeS allows for a stable accuracy growth, as it reduces  $\phi_1$  and  $\phi_2$  in (16). Moreover, FLARE tolerates non-uniform device selection with limited accuracy loss, e.g., only 2.4% under MaS, as revealed by comparing Figs. 4(b) and 4(d).

Under FMaS and FMeS,  $\bar{\tau}_r$  is adjusted. FMeS tends to offer greater robustness since no excessively large  $\bar{\tau}_r$  is produced. However, the gain of FMeS is marginal, compared to FedAvg. Additionally, FMaS and FMeS only introduce negligible additional communications between the BS and devices in the first training round, making them slightly more communication-efficient than MaS and MeS.

### C. Comparison with State-Of-The-Art

We proceed to evaluate the FLARE framework for resource-constrained WFL, by considering the following benchmarks.

- Pre-tuned Scheduling (PS) [27]: The BS randomly selects  $M_r$  devices in each round  $r$ . The optimal value of  $M_r$  is pre-tuned experimentally.
- Channel Proportion (CP) [24]: The BS selects the device subset  $\mathcal{M}_r$  with the best instantaneous channel qualities

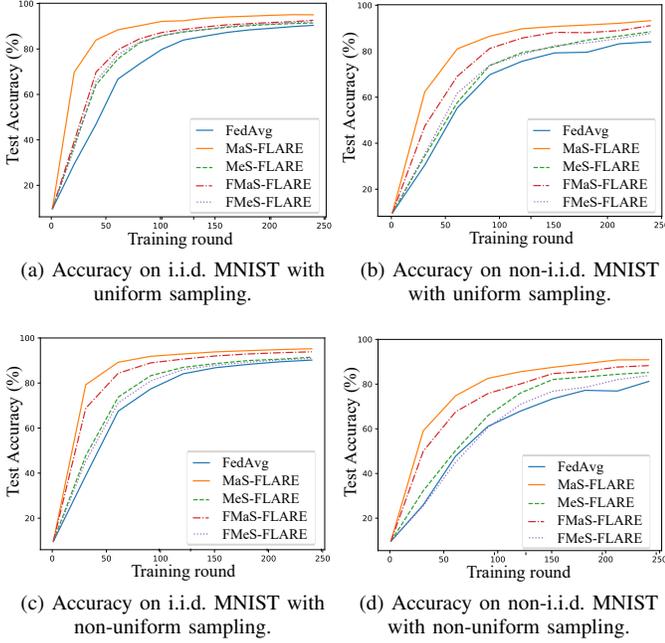


Fig. 4. The performance of different  $\bar{\tau}_r$ -strategies of FLARE on MNIST with both uniform and non-uniform sampling.

iteratively until the delay constraint (11d) is satisfied; that is, based on the channel states of the devices, the BS selects devices one after another and utilizes **Theorem 3** to allocate the available bandwidth until  $t_{\text{thr}}$  is reached.

- Device Maximization (DM) [23]: The BS iteratively selects the devices that minimize the growth of the total delay  $t_r$ , and then distributes the available bandwidth  $B$  evenly among the selected devices.
- Computation Minimization (CM) [38]: The BS ranks the expected computing times of all devices, followed by a binary search for the largest device subset satisfying all constraints.

For fair comparisons, the optimal bandwidth allocation policy in **Theorem 3** is also applied in PS, CP, and CM. We consider two ways to implement the proposed scheduling policy under FLARE. The one utilizing the MaS strategy to determine  $\bar{\tau}_r$  round-by-round is dubbed SPF. The one without the learning rate adjustments is called SP.

The convergence performances of the different policies and datasets are plotted in Figs. 5 and 6. For MNIST, two settings are considered:  $t_{\text{thr}} = 0.4$  and 1. For CIFAR-10,  $t_{\text{thr}} = 1.5$  and 5. In both figures, SPF and SP consistently outperform the benchmarks under both i.i.d. and non-i.i.d. datasets. CM performs the worst, followed by DM in all cases. This is due to their tendency to favor devices with insufficient local updates, thus hindering convergence, even though DM considers channel states. Similarly, CP focuses solely on the channel strengths of the devices and overlooks their different computational abilities and unequal contributions to model updates, while PS lacks adaptability to time-varying environments due to a persistent number of devices selected per round.

By comparing SPF and SP, it is evident that the inclusion of

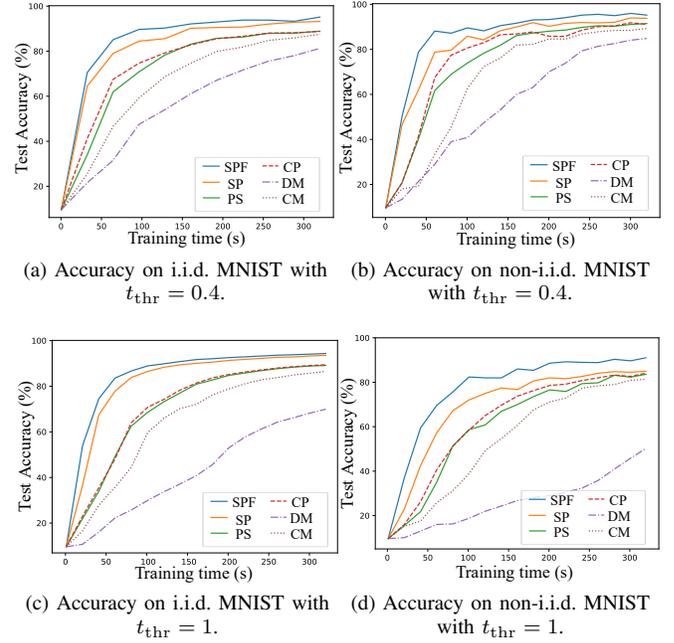


Fig. 5. Convergence performance of different policies on MNIST.

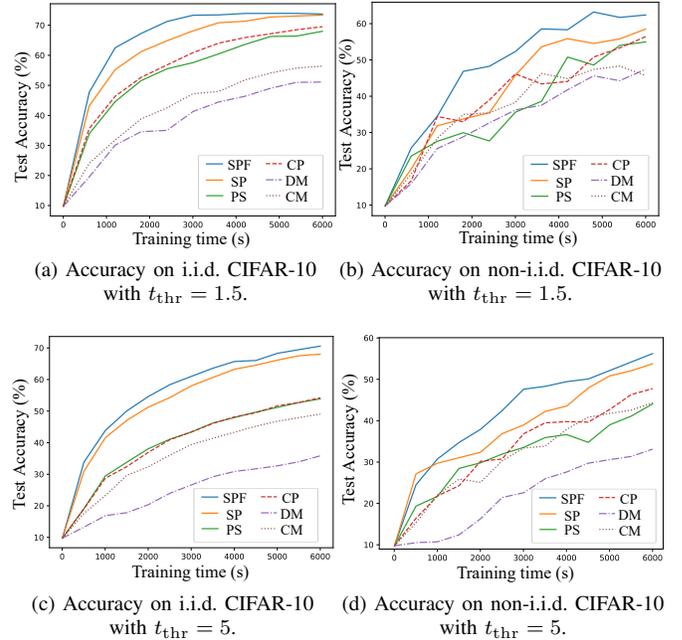


Fig. 6. Convergence performance of different policies on CIFAR-10.

FLARE consistently improves performance, especially under the non-i.i.d. settings. This can be attributed to the exacerbation of differences between the local models induced by heterogeneous updates in non-i.i.d. distributions, leading to an increased performance gap between SP and SPF.

We examine the average number of participating devices per round and their local update numbers under different policies in Fig. 7. It is noticed that SPF/SP demonstrates consistently the highest number of local updates, and maintains a sufficient number of participating devices in both the cases depicted

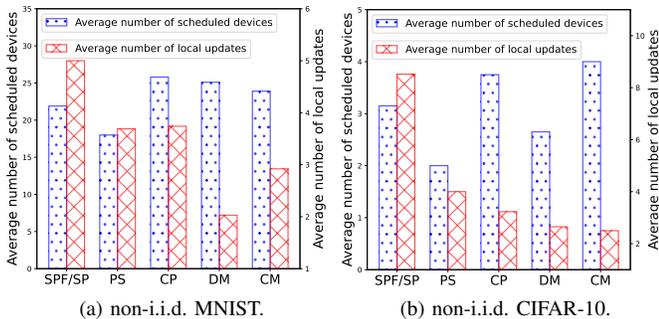


Fig. 7. Average participated devices per round and local updates per device of different policies with  $t_{\text{thr}} = 2$ . For SPF/SP, the participated device subset is determined by calling Alg. 1 at the BS side in each round.

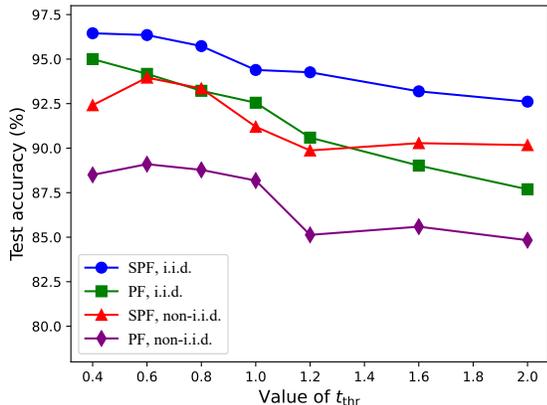


Fig. 8. The maximum achievable accuracy on i.i.d. and non-i.i.d. MNIST dataset. The total training time  $t$  is 300s.

in Figs. 7(a) and 7(b). This indicates that it is crucial to select devices with a larger number of local updates (and hence higher contributions) while ensuring adequate device participation, to accelerate the convergence.

Last but not least, we investigate the impact of  $t_{\text{thr}}$  on the achievable accuracy of SPF and SP within a fixed training time. Specifically, we explore the values of  $t_{\text{thr}}$  from  $\{0.4, 0.6, 0.8, 1, 1.2, 1.6, 2\}$ . As  $t_{\text{thr}}$  increases, the system becomes more tolerant of larger  $t_{r,i}^{\text{comp}}$  and  $t_{r,i}^{\text{comm}}$ , resulting in an increased number of devices selected according to Alg. 1. Meanwhile, the number of training rounds decreases. As shown in Fig. 8, in the i.i.d. case where data heterogeneity does not exist, the performance degenerates with  $t_{\text{thr}}$  since the enlarging subset of selected devices cannot fully compensate for the reduction in the number of training rounds under SPF and PF. This contrasts with the non-i.i.d. case, where the optimal performance is acquired at  $t_{\text{thr}} = 0.6$ , indicating a trade-off between training time and training rounds. Consistent accuracy is also attained when  $t_{\text{thr}} > 0.6$ . This stability stems from the fact that the increasing number of devices can counteract effectively the adverse effect of data heterogeneity.

## VI. CONCLUSION

In this paper, we presented the FLARE framework specifically designed to address data, device, and channel heterogeneity and expedite WFL training. Our contribution included

a comprehensive convergence analysis of FLARE with non-convex loss functions under data and device heterogeneity and arbitrary device scheduling policy. We optimized the device selection and bandwidth allocation of FLARE in the face of resource constraints by minimizing the convergence upper bound. Efficient solutions were developed by revealing a nested optimization structure of the intended problem, as well as a much simpler linear structure when the neural network models have large Lipschitz constants. Experiments corroborated the efficacy of FLARE in mitigating heterogeneity across diverse system settings, and that the proposed scheduling policy consistently outperformed the state of the art in test accuracy.

## APPENDIX A PROOF OF THEOREM 1

We begin with the following general assumption:

$$\mathbb{E}_{\xi} [\|\nabla f(\mathbf{x}; \xi) - \nabla f(\mathbf{x})\|^2] \leq \sigma^2 \text{ and } \|\nabla f(\mathbf{x})\| \leq g. \quad (26)$$

Setting the local learning rate  $\tilde{\eta}_l = \eta_l \bar{\tau}_r / \tau_r$  in round  $r$  and denoting  $\xi^i$  as the local data used for the  $i$ -th local iteration in an SGD round, we have

$$\begin{aligned} \|\tilde{\mathbf{w}} - \mathbf{w}\|^2 &= \left\| \tilde{\eta}_l \sum_{i=0}^{\tau_r-1} \nabla f(\tilde{\mathbf{w}}^i, \xi^i) - \eta_l \sum_{j=0}^{\bar{\tau}_r-1} \nabla f(\mathbf{w}^j, \xi^j) \right\|^2 \\ &\stackrel{(a)}{\leq} 2\tau_r \tilde{\eta}_l^2 \sum_{i=0}^{\tau_r-1} \|\nabla f(\tilde{\mathbf{w}}^i, \xi^i) - \nabla f(\tilde{\mathbf{w}}^i) + \nabla f(\tilde{\mathbf{w}}^i)\|^2 \\ &\quad + 2\bar{\tau}_r \eta_l^2 \sum_{j=0}^{\bar{\tau}_r-1} \|\nabla f(\mathbf{w}^j, \xi^j) - \nabla f(\mathbf{w}^j) + \nabla f(\mathbf{w}^j)\|^2 \\ &\stackrel{(b)}{\leq} 4\tau_r \tilde{\eta}_l^2 \sum_{i=0}^{\tau_r-1} \left( \|\nabla f(\tilde{\mathbf{w}}^i, \xi^i) - \nabla f(\tilde{\mathbf{w}}^i)\|^2 + \|\nabla f(\tilde{\mathbf{w}}^i)\|^2 \right) \\ &\quad + 4\bar{\tau}_r \eta_l^2 \sum_{j=0}^{\bar{\tau}_r-1} \left( \|\nabla f(\mathbf{w}^j, \xi^j) - \nabla f(\mathbf{w}^j)\|^2 + \|\nabla f(\mathbf{w}^j)\|^2 \right) \\ &\stackrel{(c)}{\leq} 2 \cdot 4\bar{\tau}_r^2 \eta_l^2 (g + \sigma^2) = 8\eta_l^2 \bar{\tau}_r^2 (g + \sigma^2). \end{aligned} \quad (27)$$

where (a) and (b) are obtained because of the Cauchy-Schwartz inequality, and (c) is obtained by substituting (26).

## APPENDIX B PROOF OF LEMMA 1

Since  $f_i, \forall i \in \mathcal{K}$  is  $L$ -smooth,  $F(\mathbf{w})$  is  $L$ -smooth. Since  $\nabla F$  is also Lipschitz continuous, it readily leads to

$$\mathbb{E}[F(\mathbf{w}_{r+1}) - F(\mathbf{w}_r)] \leq \underbrace{\mathbb{E}[\langle \nabla F(\mathbf{w}_r), \eta_g \Delta_r \rangle]}_{A_1} + \underbrace{\frac{L}{2} \mathbb{E}[\|\eta_g \Delta_r\|^2]}_{A_2}, \quad (28)$$

where  $\mathbb{E}[\cdot]$  takes expectation over random minibatch sampling.

The upper bound of  $A_1$  is given by

$$A_1 = -\eta_g \eta_l \left\langle \nabla F(\mathbf{w}_r), \mathbb{E} \left[ \sum_{i \in \mathcal{M}_r} \frac{1}{M_r} \sum_{j=0}^{\tau_{r,i}-1} \frac{\bar{\tau}_r}{\tau_{r,i}} g_{r,i}^j \right] \right\rangle$$

$$\begin{aligned}
&= -\eta_g \eta \left\langle \sqrt{\bar{\tau}_r} \nabla F(\mathbf{w}_r), \frac{1}{M_r} \sum_{i \in \mathcal{M}_r} \sum_{j=0}^{\tau_{r,i}-1} \frac{\sqrt{\bar{\tau}_r}}{\tau_{r,i}} \nabla f_i(\mathbf{w}_{r,i}^j) \right\rangle \\
&\stackrel{(a)}{\leq} a_1 + \frac{\eta_g \eta}{2} \left\| \frac{1}{M_r} \sum_{i \in \mathcal{M}_r} \sum_{j=0}^{\tau_{r,i}-1} \frac{\sqrt{\bar{\tau}_r}}{\tau_{r,i}} (\nabla f_i(\mathbf{w}_{r,i}^j) - \nabla f_i(\mathbf{w}_r)) \right\|^2 \\
&\stackrel{(b)}{\leq} a_1 + \underbrace{\frac{L^2 \eta_g \eta}{2M_r} \sum_{i \in \mathcal{M}_r} \frac{\bar{\tau}_r}{\tau_{r,i}} \sum_{j=0}^{\tau_{r,i}-1} \mathbb{E} \left[ \|\mathbf{w}_{r,i}^j - \mathbf{w}_r\|^2 \right]}_{B_1}, \quad (29)
\end{aligned}$$

where  $a_1 = -\frac{\bar{\tau}_r \eta_g \eta}{2} \|\nabla F(\mathbf{w}_r)\|^2$ . Here, (a) is due to  $-2 \langle \mathbf{a}, \mathbf{b} \rangle \leq -\|\mathbf{a}\|^2 + \|\mathbf{a} - \mathbf{b}\|^2$  with  $\mathbf{a} = \sqrt{\bar{\tau}_r} \nabla F(\mathbf{w}_r)$  and  $\mathbf{b} = \frac{1}{M_r} \sum_{i \in \mathcal{M}_r} \sum_{j=0}^{\tau_{r,i}-1} \frac{\sqrt{\bar{\tau}_r}}{\tau_{r,i}} \nabla f_i(\mathbf{w}_{r,i}^j)$ ; and (b) is from **Assumption 1**.

To derive the upper bound of  $B_1$ , we start with the upper bound of  $\mathbb{E} \left[ \|\mathbf{w}_{r,i}^j - \mathbf{w}_r\|^2 \right] \forall i, j$ . Note that  $\mathbb{E} \left[ \|\mathbf{w}_{r,i}^j - \mathbf{w}_r\|^2 \right] = 0$  when  $j = 0$ , since  $\mathbf{w}_{r,i}^0 = \mathbf{w}_r$  by definition. When  $j \geq 1$ , we have

$$\begin{aligned}
\mathbb{E} \left[ \|\mathbf{w}_{r,i}^j - \mathbf{w}_r\|^2 \right] &= \mathbb{E} \left[ \left\| \mathbf{w}_{r,i}^{j-1} - \frac{\bar{\tau}_r}{\tau_{r,i}} \eta g_{r,i}^{j-1} - \mathbf{w}_r \right\|^2 \right] \\
&\stackrel{(a)}{=} \mathbb{E} \left[ \left\| \mathbf{w}_{r,i}^{j-1} - \frac{\bar{\tau}_r}{\tau_{r,i}} \eta \nabla f_i(\mathbf{w}_{r,i}^{j-1}) - \mathbf{w}_r \right\|^2 \right] + \frac{\bar{\tau}_r^2}{\tau_{r,i}^2} \sigma^2 \eta^2 \\
&\stackrel{(b)}{\leq} \left( 1 + \frac{1}{b\tau_{r,i} - 1} \right) \mathbb{E} \left[ \|\mathbf{w}_{r,i}^{j-1} - \mathbf{w}_r\|^2 \right] \\
&\quad + (1 + b\tau_{r,i} - 1) \frac{\bar{\tau}_r^2}{\tau_{r,i}^2} \eta^2 \|\nabla f_i(\mathbf{w}_{r,i}^{j-1})\|^2 + \frac{\bar{\tau}_r^2}{\tau_{r,i}^2} \sigma^2 \eta^2 \\
&\stackrel{(c)}{\leq} \left( 1 + \frac{1}{b\tau_{r,i} - 1} \right) \mathbb{E} \left[ \|\mathbf{w}_{r,i}^{j-1} - \mathbf{w}_r\|^2 \right] + \frac{\bar{\tau}_r^2}{\tau_{r,i}^2} \sigma^2 \eta^2 \\
&\quad + \frac{2b\bar{\tau}_r^2 \eta^2}{\tau_{r,i}} \|\nabla f_i(\mathbf{w}_{r,i}^{j-1}) - \nabla f_i(\mathbf{w}_r)\|^2 + \frac{2b\bar{\tau}_r^2 \eta^2}{\tau_{r,i}} \|\nabla f_i(\mathbf{w}_r)\|^2 \\
&\stackrel{(d)}{\leq} \left( 1 + \frac{1}{b\tau_{r,i} - 1} + 2b \frac{L^2 \bar{\tau}_r^2}{\tau_{r,i} \eta^2} \right) \mathbb{E} \left[ \|\mathbf{w}_{r,i}^{j-1} - \mathbf{w}_r\|^2 \right] \\
&\quad + \frac{2b\bar{\tau}_r^2 \eta^2}{\tau_{r,i}} \|\nabla f_i(\mathbf{w}_r)\|^2 + \frac{\bar{\tau}_r^2}{\tau_{r,i}^2} \sigma^2 \eta^2 \\
&\stackrel{(e)}{\leq} \frac{b\tau_{r,i} + 1}{b\tau_{r,i} - 1} \mathbb{E} \left[ \|\mathbf{w}_{r,i}^{j-1} - \mathbf{w}_r\|^2 \right] + \frac{\bar{\tau}_r^2 \sigma^2 \eta^2}{\tau_{r,i}^2} + \frac{2b\bar{\tau}_r^2 \eta^2}{\tau_{r,i}} \|\nabla f_i(\mathbf{w}_r)\|^2, \quad (30)
\end{aligned}$$

where (a) holds since  $\mathbb{E} [\|\mathbf{x}\|^2] = \mathbb{E} [\|\mathbf{x} - \mathbb{E}[\mathbf{x}]\|^2] + \|\mathbb{E}[\mathbf{x}]\|^2$ ; (b) is due to the relaxed triangle inequality  $\|\mathbf{x}_1 + \mathbf{x}_2\|^2 \leq (1 + \beta) \|\mathbf{x}_1\|^2 + \left(1 + \frac{1}{\beta}\right) \|\mathbf{x}_2\|^2, \forall \beta = b\tau_{r,i} - 1 \geq 0$ ; (c) comes from the Cauchy-Schwartz inequality; (d) is based on **Assumption 1**; (e) holds under the condition that  $2b \frac{L^2 \bar{\tau}_r^2}{\tau_{r,i} \eta^2} \leq \frac{1}{b\tau_{r,i} - 1}$ . Without loss of generality, we set  $b = 1.2$  to ensure  $\beta = b\tau_{r,i} - 1 > 0$  since  $\tau_{r,i} \geq 1$ . Moreover,  $\eta \leq \frac{0.58}{L\bar{\tau}_r}$  is required by resolving the condition  $2.4 \frac{L^2 \bar{\tau}_r^2}{\tau_{r,i} \eta^2} \leq \frac{1}{1.2\tau_{r,i} - 1}$ .

By unrolling the recursion of (30),  $B_1$  is upper bounded by

$$B_1 \leq \mathbb{E} \left[ \frac{L^2 \eta_g \eta}{2M_r} \sum_{i \in \mathcal{M}_r} \frac{\bar{\tau}_r e_{r,i}}{\tau_{r,i}} \right]$$

$$\begin{aligned}
&\stackrel{(a)}{\leq} \frac{L^2 \eta_g \eta}{2M_r} \mathbb{E} \left[ \sum_{i \in \mathcal{M}_r} 2.4 \bar{\tau}_r^3 \eta^2 \|\nabla f_i(\mathbf{w}_r)\|^2 + \frac{\bar{\tau}_r^3 \eta^2}{\tau_{r,i}} \sigma^2 \right] \\
&\stackrel{(b)}{\leq} 1.2 L^2 \bar{\tau}_r^3 \eta_g \eta^3 [G^2 + H^2 \|\nabla F(\mathbf{w}_r)\|^2] + \frac{L^2 \bar{\tau}_r^3 \eta_g \eta^3}{2M_r} \sum_{i \in \mathcal{M}_r} \frac{\sigma^2}{\tau_{r,i}}, \quad (31)
\end{aligned}$$

where  $e_{r,i} = \frac{1.2\tau_{r,i}-1}{2} \left[ \frac{1-3.2\tau_{r,i}}{2} + \frac{1.2\tau_{r,i}-1}{2} \left( \frac{1.2\tau_{r,i}+1}{1.2\tau_{r,i}-1} \right)^{\tau_{r,i}} \right] \cdot \left[ 2.4 \frac{\bar{\tau}_r^2}{\tau_{r,i}} \eta^2 \|\nabla f_i(\mathbf{w}_r)\|^2 + \frac{\bar{\tau}_r^2}{\tau_{r,i}^2} \sigma^2 \eta^2 \right]$ ; (a) holds due to  $\frac{1.2\tau_{r,i}-1}{2\tau_{r,i}^2} \left[ \frac{1-3.2\tau_{r,i}}{2} + \frac{1.2\tau_{r,i}-1}{2} \left( \frac{1.2\tau_{r,i}+1}{1.2\tau_{r,i}-1} \right)^{\tau_{r,i}} \right] < 1, \forall \tau_{r,i} \geq 1$ ; and (b) follows readily from **Assumption 3**.

Next, we derive the upper bound of  $A_2$  in (28), as given by

$$\begin{aligned}
A_2 &= \frac{L}{2} \eta_g^2 \eta^2 \mathbb{E} \left[ \left\| \sum_{i \in \mathcal{M}_r} \frac{1}{M_r} \sum_{j=0}^{\tau_{r,i}-1} \frac{\bar{\tau}_r}{\tau_{r,i}} g_{r,i}^j \right\|^2 \right] \\
&\stackrel{(a)}{=} \frac{L}{2} \eta_g^2 \eta^2 \mathbb{E} \left[ \left\| \frac{1}{M_r} \sum_{i \in \mathcal{M}_r} \sum_{j=0}^{\tau_{r,i}-1} \frac{\bar{\tau}_r}{\tau_{r,i}} (g_{r,i}^j - \nabla f_i(\mathbf{w}_{r,i}^j)) \right\|^2 \right] \\
&\quad + \frac{L}{2} \eta_g^2 \eta^2 \left\| \frac{1}{M_r} \sum_{i \in \mathcal{M}_r} \sum_{j=0}^{\tau_{r,i}-1} \frac{\bar{\tau}_r}{\tau_{r,i}} \nabla f_i(\mathbf{w}_{r,i}^j) \right\|^2 \\
&\stackrel{(b)}{\leq} \frac{L}{2} \eta_g^2 \eta^2 \frac{1}{M_r^2} \sum_{i \in \mathcal{M}_r} \frac{\bar{\tau}_r^2}{\tau_{r,i}^2} \sum_{j=0}^{\tau_{r,i}-1} \sigma^2 \\
&\quad + \frac{L}{2} \eta_g^2 \eta^2 \bar{\tau}_r^2 \left\| \frac{1}{M_r} \sum_{i \in \mathcal{M}_r} \sum_{j=0}^{\tau_{r,i}-1} \frac{1}{\tau_{r,i}} \nabla f_i(\mathbf{w}_{r,i}^j) \right\|^2 \\
&\stackrel{(c)}{\leq} \frac{L \eta_g^2 \eta^2 \bar{\tau}_r^2}{2M_r^2} \sigma^2 \sum_{i \in \mathcal{M}_r} \frac{1}{\tau_{r,i}} + \frac{L \bar{\tau}_r^2 \eta_g^2 \eta^2}{2\rho^2} \|\nabla F(\mathbf{w}_r)\|^2, \quad (32)
\end{aligned}$$

where (a) is due to  $\mathbb{E} [\|\mathbf{x}\|^2] = \mathbb{E} [\|\mathbf{x} - \mathbb{E}[\mathbf{x}]\|^2] + \|\mathbb{E}[\mathbf{x}]\|^2$ ; (b) is due to  $\mathbb{E} [\|x_1 + \dots + x_n\|^2] \leq n \mathbb{E} [\|x_1\|^2 + \dots + \|x_n\|^2]$ . We can further prove that there exists a constant  $\rho > 0$  such that

$$\left\| \frac{1}{M_r} \sum_{i \in \mathcal{M}_r} \frac{1}{\tau_{r,i}} \nabla f_i(\mathbf{w}_r) \right\|^2 \leq \left\| \frac{1}{\rho} \nabla F(\mathbf{w}_r) \right\|^2, \forall r, i. \quad (33)$$

Applying the Cauchy-Schwartz inequality and **Assumption 3** on the LHS of (33) yields  $\rho^2 \leq \frac{M_r \|\nabla F(\mathbf{w}_r)\|^2}{\sum_{i \in \mathcal{M}_r} \frac{1}{\tau_{r,i}} (G^2 + H^2 \|\nabla F(\mathbf{w}_r)\|^2)}$ . Setting  $\tau_{r,i} = 1, \forall r, i$ , the upper bound of  $\rho$  is  $\rho \leq \frac{1}{\sqrt{H^2 + G^2 / \|\nabla F(\mathbf{w}_r)\|^2}}$  with a bounded global gradient  $\|\nabla F(\mathbf{w}_r)\|$ , and (c) follows from (33).

With the upper bound of  $A_1$  based on (29) and (31) and the upper bound of  $A_2$  based on (32), (28) can be reorganized into (16), where we define a positive constant  $q_r$  that satisfies the condition  $0 < q_r < 1 - \frac{\bar{\tau}_r \eta_g \eta}{\rho^2} - 2.4 H^2 L^2 \bar{\tau}_r^2 \eta^2 < 1$ .

## APPENDIX C

### PROOF OF THEOREM 2

By adding  $F(\mathbf{w}^*)$  to both sides of (16) and then reorganizing and taking expectation, we have

$$\mathbb{E} [\|\nabla F(\mathbf{w}_r)\|^2] \leq a_2 + \frac{2.4}{q_r} G^2 L^2 \bar{\tau}_r^2 \eta^2 \quad (34)$$

$$+ \frac{L\bar{\tau}\eta_l}{q_r} \left( \frac{L\bar{\tau}\eta_l}{M_r} + \frac{\eta_g}{M_r^2} \right) \sigma^2 \sum_{i \in \mathcal{M}_r} \frac{1}{\tau_{r,i}},$$

where  $a_2 = \frac{2(\mathbb{E}[F(\mathbf{w}_r) - F(\mathbf{w}^*)] - \mathbb{E}[F(\mathbf{w}_{r+1}) - F(\mathbf{w}^*)])}{\bar{\tau}\eta_g\eta_l q_r}$ . To derive an upper bound of  $\sum_{i \in \mathcal{M}_r} \frac{1}{\tau_{r,i}}$  in (34), we define  $\kappa_r \triangleq \max_{r,i \in \mathcal{M}_r} \left\{ \frac{\bar{\tau}}{\tau_{r,i}} \right\}$ , and rewrite (34) as

$$\mathbb{E} \left[ \|\nabla F(\mathbf{w}_r)\|^2 \right] \leq a_2 + \frac{2.4}{q_r} G^2 L^2 \bar{\tau}^2 \eta_l^2 + \frac{L\eta_l}{q_r} \left( L\bar{\tau}\eta_l + \frac{\kappa_r \eta_g}{M_r} \right) \sigma^2. \quad (35)$$

By summing up (35) from  $r = 1, \dots, R$  and dividing both sides of the resulting inequality by  $R$ , we obtain

$$\begin{aligned} \frac{1}{R} \sum_{r=1}^R \mathbb{E} \left[ \|\nabla F(\mathbf{w}_r)\|^2 \right] &\leq \mathbb{E}_{\mathcal{M}} \left[ \frac{1}{R} \sum_{r=1}^R a_2 \right. \\ &\left. + \frac{1}{R} \sum_{r=1}^R \frac{2.4}{q_r} G^2 L^2 \bar{\tau}^2 \eta_l^2 + \frac{1}{R} \sum_{r=1}^R \frac{L\eta_l}{q_r} \left( L\bar{\tau}\eta_l + \frac{\kappa_r \eta_g}{M_r} \right) \sigma^2 \right], \end{aligned} \quad (36)$$

where  $\mathbb{E}_{\mathcal{M}}[\cdot]$  takes expectation over the past device scheduling decisions, i.e.,  $\mathcal{M} = \{\mathcal{M}_r | r = 1, \dots, R\}$ . Define  $q \triangleq \min_r \{q_r\}$ ,  $\kappa \triangleq \max_r \{\kappa_r\}$ , and  $M \triangleq \min_r \{M_r\}$  with  $M \geq 1$ , and further rescale the RHS of (36). We finally obtain

$$\begin{aligned} \frac{1}{R} \sum_{r=1}^R \mathbb{E} \left[ \|\nabla F(\mathbf{w}_r)\|^2 \right] &\leq \frac{2[F(\mathbf{w}_1) - F(\mathbf{w}^*)]}{\bar{\tau}\eta_g\eta_l q R} \\ &+ \frac{2.4}{q} G^2 L^2 \bar{\tau}^2 \eta_l^2 + \frac{L\eta_l}{q} \left( L\bar{\tau}\eta_l + \frac{\kappa \eta_g}{M} \right) \sigma^2, \end{aligned} \quad (37)$$

which completes this proof.

#### APPENDIX D PROOF OF THEOREM 3

For illustration convenience, we suppress the indexes of training rounds and devices, i.e.,  $r$  and  $i$ . The transmission rate  $u$  monotonically increases w.r.t  $b > 0$ , since

$$\frac{du}{db} = \frac{1}{\ln 2} \left( \ln \left( 1 + \frac{ph^2}{bN_0} \right) - \frac{ph^2}{bN_0 + ph^2} \right) > 0. \quad (38)$$

We can prove by contradiction that  $t_r^*$  must be equal for all selected devices under the optimal bandwidth allocation. To see it, we first assume that  $\mathbf{b}$  denotes the optimal bandwidth allocation, with  $t_{r,i}$  inconsistent across devices, and let the fastest and slowest devices be denoted by  $m, n$ , respectively. Due to synchronous aggregation, the total delay of each round is determined by the slowest device. Then according to (38), the bandwidth  $b_m$  allocated to Device  $m$  can be partially reassigned to Device  $n$  to increase its transmission rate and in turn decrease the total latency determined by Device  $n$ . Such an operation can be repeated continuously until a solution  $\mathbf{b}^*$  of Problem **P3** is obtained when the latency  $t_{r,i}$  of all the selected devices is the same. The solution  $\mathbf{b}^*$  clearly leads to a strictly smaller objective in (20a) than that of "optimal"  $\mathbf{b}'$ , leading to a contradiction. Having proved that  $t_r^*$  is equal, the

optimal bandwidth allocation can be obtained by solving the following system of equations:

$$\begin{cases} t_{r,i}^{\text{comp}} + \frac{S}{b_{r,i}^* \log_2 \left( 1 + \frac{p_i h_{r,i}^2}{b_{r,i}^* N_0} \right)} = t_r^*, \forall i \in \mathcal{M}_r \\ \sum_{i \in \mathcal{M}_r} b_{r,i}^* = B. \end{cases} \quad (39)$$

By leveraging the definition of Lambert-W function, the theorem readily follows.

#### APPENDIX E PROOF OF THEOREM 4

Consider the currently selected subset of devices, denoted as  $\mathcal{Q}_k$ , including  $Q$  devices, and the objective function (19a). If the introduction of another device  $i \in \mathcal{K} \setminus \mathcal{Q}_k$  into  $\mathcal{Q}_k$  allows for a further decrease of (19a), then we have

$$\frac{Q+1+\gamma}{(Q+1)^2} \left( \Delta + \frac{1}{\tau_{r,i}} \right) < \frac{Q+\gamma}{Q^2} \Delta, \quad (40)$$

where  $\Delta = \sum_{q \in \mathcal{Q}_k} \frac{1}{\tau_{r,q}}$ . This readily implies (24), the proof is thus complete.

#### REFERENCES

- [1] W. Jiang, B. Han, M. A. Habibi, and H. D. Schotten, "The road towards 6G: A comprehensive survey," *IEEE Open J. Commun. Soc.*, vol. 2, pp. 334–366, 2021.
- [2] S. Hu, X. Chen, W. Ni, E. Hossain, and X. Wang, "Distributed machine learning for wireless communication networks: Techniques, architectures, and applications," *IEEE Commun. Surv. Tutorials*, vol. 23, no. 3, pp. 1458–1493, 2021.
- [3] T. Chen, S. Barbarossa, X. Wang, G. B. Giannakis, and Z.-L. Zhang, "Learning and management for internet of things: Accounting for adaptivity and scalability," *Proc. IEEE*, vol. 107, no. 4, pp. 778–796, 2019.
- [4] J. Konečný, H. B. McMahan, F. X. Yu, A. T. Suresh, D. Bacon, and P. Richtárik, "Federated learning: Strategies for improving communication efficiency," in *Proc. Int. Conf. Learn. Representations (ICLR)*, 2018.
- [5] 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, pp. 1273–1282.
- [6] S. Niknam, H. S. Dhillon, and J. H. Reed, "Federated learning for wireless communications: Motivation, opportunities, and challenges," *IEEE Commun. Mag.*, vol. 58, no. 6, pp. 46–51, 2020.
- [7] H. Shiri, J. Park, and M. Bennis, "Communication-efficient massive UAV online path control: Federated learning meets mean-field game theory," *IEEE Trans. Commun.*, vol. 68, no. 11, pp. 6840–6857, 2020.
- [8] F. Yin, Z. Lin, Q. Kong, Y. Xu, D. Li, S. Theodoridis, and S. R. Cui, "Fedloc: Federated learning framework for data-driven cooperative localization and location data processing," *IEEE Open J. Signal Process.*, vol. 1, pp. 187–215, 2020.
- [9] M. Ammad-Ud-Din, E. Ivannikova, S. A. Khan, W. Oyomno, Q. Fu, K. E. Tan, and A. Flanagan, "Federated collaborative filtering for privacy-preserving personalized recommendation system," arXiv preprint arXiv:1901.09888, 2019.
- [10] A. Nilsson, S. Smith, G. Ulm, E. Gustavsson, and M. Jirstrand, "A performance evaluation of federated learning algorithms," in *Proc. Wksp. Dtrib. Infrastruct. Deep Learn. (DIDL Wksp)*, 2018, pp. 1–8.
- [11] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of FedAvg on non-iid data," arXiv preprint arXiv:1907.02189, 2019.
- [12] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 50–60, 2020.
- [13] B. Luo, X. Li, S. Wang, J. Huang, and L. Tassiulas, "Cost-effective federated learning in mobile edge networks," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 12, pp. 3606–3621, 2021.
- [14] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, "Federated learning with non-iid data," arXiv preprint arXiv:1806.00582, 2018.

- [15] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "Scaffold: Stochastic controlled averaging for federated learning," in *Proc. Int. Conf. Mach. Learn. (ICML)*. PMLR, 2020, pp. 5132–5143.
- [16] Q. Li, B. He, and D. Song, "Model-contrastive federated learning," in *Proc. IEEE/CVF Conf. Comp. Vision Pattern Recognit. (CVPR)*, 2021, pp. 10 713–10 722.
- [17] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," *Proc. Mach. Learn. Syst. (MLSys)*, vol. 2, pp. 429–450, 2020.
- [18] 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.
- [19] A. Mitra, R. Jaafar, G. J. Pappas, and H. Hassani, "Linear convergence in federated learning: Tackling client heterogeneity and sparse gradients," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, vol. 34, 2021, pp. 14 606–14 619.
- [20] B. Xiao, X. Yu, W. Ni, X. Wang, and H. V. Poor, "Over-the-air federated learning: Status quo, open challenges, and future directions," *Fundamental Research*, pp. 1–1, 2024.
- [21] X. Yu, B. Xiao, W. Ni, and X. Wang, "Optimal adaptive power control for over-the-air federated edge learning under fading channels," *IEEE Trans. Commun.*, vol. 71, no. 9, pp. 5199–5213, 2023.
- [22] N. H. Tran, W. Bao, A. Zomaya, M. N. H. Nguyen, and C. S. Hong, "Federated learning over wireless networks: Optimization model design and analysis," in *Proc. IEEE Conf. Comput. Commun. (INFOCOM)*, 2019, pp. 1387–1395.
- [23] T. Nishio and R. Yonetani, "Client selection for federated learning with heterogeneous resources in mobile edge," in *Proc. IEEE Int. Conf. Commun. (ICC)*, 2019, pp. 1–7.
- [24] H. H. Yang, Z. Liu, T. Q. S. Quek, and H. V. Poor, "Scheduling policies for federated learning in wireless networks," *IEEE Trans. Commun.*, vol. 68, no. 1, pp. 317–333, 2020.
- [25] M. M. Amiri, D. Gündüz, S. R. Kulkarni, and H. V. Poor, "Convergence of update aware device scheduling for federated learning at the wireless edge," *IEEE Trans. Wireless Commun.*, vol. 20, no. 6, pp. 3643–3658, 2021.
- [26] C.-H. Hu, Z. Chen, and E. G. Larsson, "Scheduling and aggregation design for asynchronous federated learning over wireless networks," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 4, pp. 874–886, 2023.
- [27] W. Shi, S. Zhou, Z. Niu, M. Jiang, and L. Geng, "Joint device scheduling and resource allocation for latency constrained wireless federated learning," *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 453–467, 2021.
- [28] M. Chen, Z. Yang, W. Saad, C. Yin, H. V. Poor, and S. Cui, "A joint learning and communications framework for federated learning over wireless networks," *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 269–283, 2021.
- [29] 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, 2020.
- [30] B. Luo, W. Xiao, S. Wang, J. Huang, and L. Tassiulas, "Tackling system and statistical heterogeneity for federated learning with adaptive client sampling," in *Proc. IEEE Conf. Comput. Commun. (INFOCOM)*, 2022, pp. 1739–1748.
- [31] E. Diao, J. Ding, and V. Tarokh, "HeteroFL: Computation and communication efficient federated learning for heterogeneous clients," in *Proc. Int. Conf. Learn. Representations (ICLR)*, 2021.
- [32] Z. Qu, X. Li, R. Duan, Y. Liu, B. Tang, and Z. Lu, "Generalized federated learning via sharpness aware minimization," in *International conference on machine learning*. PMLR, 2022, pp. 18 250–18 280.
- [33] J. Xu and H. Wang, "Client selection and bandwidth allocation in wireless federated learning networks: A long-term perspective," *IEEE Trans. Wireless Commun.*, vol. 20, no. 2, pp. 1188–1200, 2021.
- [34] H. Yang, M. Fang, and J. Liu, "Achieving linear speedup with partial worker participation in non-iid federated learning," in *Proc. Int. Conf. Learn. Representations (ICLR)*, 2021.
- [35] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [36] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [37] A. Reiszadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani, "Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization," in *Proc. Int. Conf. Artif. Intell. Stat. (AISTATS)*. PMLR, 2020, pp. 2021–2031.
- [38] C. Wang, X. Wei, and P. Zhou, "Optimize scheduling of federated learning on battery-powered mobile devices," in *Proc. IEEE Int. Parallel Distrib. Process. Symp. (IPDPS)*, 2020, pp. 212–221.