Towards Flexible Device Participation in Federated Learning

Yichen Ruan Carnegie Mellon University Xiaoxi Zhang

Sun Yat-Sen University

Shu-Che Liang

Carlee Joe-Wong

Carnegie Mellon University Carnegie Mellon University

Abstract

Traditional federated learning algorithms impose strict requirements on the participation rates of devices, which limit the potential reach of federated learning. This paper extends the current learning paradigm to include devices that may become inactive, compute incomplete updates, and depart or arrive in the middle of training. We derive analytical results to illustrate how allowing more flexible device participation can affect the learning convergence when data is not independently and identically distributed (non-IID). We then propose a new federated aggregation scheme that converges even when devices may be inactive or return incomplete updates. We also study how the learning process can adapt to early departures or late arrivals, and analyze their impacts on the convergence.

1 Introduction

Federated learning is a cutting-edge learning framework that allows distributed devices to train a shared machine learning model cooperatively without sharing the raw data. In recent years, federated learning has exhibited remarkable performance in many applications such as next word suggestion, fault detection, and learning on private medical data (Li et al., 2020a). Generic federated learning involves a coordinator and a collection of devices. The training procedure consists of multiple rounds, each of which includes the following three steps: 1) Synchronization: the coordinator synchronizes the latest global model with all devices. 2) Local updates: each device trains a local model for a few local epochs, using samples from its local dataset. 3) Aggregation: the coordinator aggregates some, or all, of the local models to produce the next global model.

Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS) 2021, San Diego, California, USA. PMLR: Volume 130. Copyright 2021 by the author(s).

Our work focuses on cross-device federated learning (Kairouz et al., 2019), where participating entities are mostly mobile devices such as smart phones and tablets. These devices generally have limited computing and communication resources, e.g., due to battery limitations, and have different training data distributions, i.e., data is not independently and identically distributed (non-IID) among devices (Li et al., 2020b). To relieve the computation and communication burden, in the last step of the training procedure, the federated learning coordinator may only aggregate a subset of local models. However, only a few device selection policies ensure convergence in the non-IID setting, and the selection must be independent of the hardware status of devices (Li et al., 2020b). In other words, for the training to converge successfully, all selected devices must be able to train their local models and upload the results whenever they are selected. This is why the traditional federated learning paradigm requires participating devices to be dedicated to the training during the entire federated learning period, e.g., the popular FedAvg algorithm assumes mobile users will participate only when their phones are currently plugged-in, and have unlimited WI-FI access (McMahan et al., 2016).

Considering that federated learning typically takes thousands of communication rounds to converge, it is difficult to ensure that all devices will be available during the entire training in practice. Moreover, there are typically multiple apps running simultaneously on user devices, competing for already highly constrained hardware resources. As such, it cannot be guaranteed that devices will complete their assigned training tasks in every training round as expected. A similar challenge also arises in cloud based distributed learning due to the increasingly popular usage of preemptive cloud services, where the user process can be interrupted unexpectedly (Zhang et al., 2020).

While many methods have been proposed to mitigate the workload of individual devices, such as weight compression and federated dropout (Caldas et al., 2018) (Konečný et al., 2016), they cannot completely remove the possibility that devices are unable to fulfill their training responsibilities, e.g., due to poor wireless

connectivity. Thus, in large scale federated learning, many resource-constrained devices have to be excluded from joining federated learning in the first place, which restricts the potential availability of training datasets, and weakens the applicability of federated learning. Furthermore, existing work does not specify how to react when confronting unexpected device behaviors, and also does not analyze the (negative) effects of such behaviors on the training progress.

In this paper, we relax these restrictions and allow devices to follow more flexible participation patterns. Specifically, the paper incorporates four situations that are not yet well discussed in the literature: 1) Incompleteness: devices might submit only partially completed work in a round. 2) Inactivity: furthermore, devices might not complete any updates, or respond to the coordinator at all. 3) Early departures: in the extreme case, existing devices might quit the training without finishing all training rounds. 4) Late arrivals: apart from existing devices, new devices might join after the training has already started.

The difference between inactivity and departure is that inactive devices will temporarily disconnect with the coordinator, but are expected to come back in the near future. In contrast, departing devices will inform the coordinator that they do not plan to rejoin the training. For example, if a user quits the app running federated learning, a message can be sent to the coordinator; the coordinator thus knows who is departing. In the meanwhile, although devices' arriving and departing seem symmetric, they affect the model training differently, and thus require distinct treatments. The key difference is that arriving devices offer extra information about the data distribution, which can be utilized to accelerate the training, while departing devices reduce our available knowledge, thus degrading the applicability of the trained model.

Our approach to improve the flexibility of device participation comprises the following components that supplement the existing FedAvg algorithm and handle the challenges brought by flexible device participation.

• Debiasing for partial model updates. FedAvg aggregates device updates as a weighted sum, with weights that are proportional to the sizes of the local datasets. This choice of aggregation coefficients yields an unbiased gradient as in the centralized setting only when all data points from all devices are equally likely to join the learning (Li et al., 2020b). However, it in general fails to guarantee convergence to the globally optimal point in the presence of partial aggregation from incomplete and inactive devices. We show that by adapting the aggregation coefficients, the bias can be reduced and the convergence to a global optimum can still be established. Further-

- more, our analysis shows the bias originates from the heterogeneity in device participation, as well as from the degree to which local datasets are not IID.
- Fast-rebooting for device arrivals. Arriving devices interrupt the training by forcing the model to re-orient to the new device's data, thus slowing the convergence process. In this paper, we propose to rapidly reboot the training in the case of device arrivals by applying extra updates from the new devices. Intuitively, since an arriving device misses all previous epochs, the model training should emphasize more on its updates to compensate. We will rigorously prove this method indeed expedites learning convergence under certain conditions.
- Redefining model applicability for device departures. A model successfully trained by federated learning is expected to be applicable to the data from all participating devices. However, when a device withdraws itself from the learning, due to the lack of its future updates, we may no longer require the trained model to perform well on its data. It is then important to redefine the model's applicability. Namely, one can either keep the departing device as a part of the global learning objective, or exclude it to focus only on the remaining devices. The decision depends on which definition yields smaller training loss. We will show the key to this determination lies in the remaining training time.

In Section 2 we review relevant literature. In Section 3 we give a convergence analysis that incorporates flexible device participation. Based on this analysis, we detail our contributions, as outlined above, in Section 4 and we experimentally verify our theoretical results in Section 5. Finally we conclude in Section 6.

2 Related Works

The celebrated federated learning algorithm named FedAvg runs the stochastic gradient descent (SGD) algorithm in parallel on each device in the system and periodically averages the updated parameters from a small set of end devices. However, its performance degrades when the local data is non-IID (Hsieh et al., 2019) (Zhao et al., 2018). A few recent works provide theoretical results for the non-IID data case. For instance, Li et al. (2020b) analyze the convergence of FedAvg on non-IID data and establish an $O(\frac{1}{T})$ convergence rate for strongly convex and smooth optimization problems, where T is the number of rounds of local SGD updates. These works either simplify the heterogeneity of the devices, e.g., ignoring cases where some devices may partially finish some aggregation rounds or quit forever during the training (Li et al., 2020b), or consider alternative objective functions for the SGD

algorithm to optimize (Li et al.) 2018). Alternatively, some recent papers propose to combine federated learning with the multi-task learning paradigm (Smith et al.) 2017) (Corinzia and Buhmann, 2019) where multiple models are trained simultaneously, but they also entail dedicated device participation throughout the training.

The FedAvq algorithm with non-IID data across devices has also been modified in specific edge computing scenarios to reduce the communication overhead (Liu et al.) (Sattler et al., 2019) (Bonawitz et al., 2019) or maintain a good training convergence under a resource budget constraint (Wang et al., 2019). However, these works do not consider the possibility that the edge devices can be unavailable during the training process or join at different times, which are the main challenges of this work. An online learning framework (Chen et al., 2019) (Han et al., 2020) (Damaskinos et al., 2020) is a possible way to enable flexible device participation in the federated learning scenario. For instance, Chen et al. (2019) propose an asynchronous federated learning algorithm to handle unbalanced data that arrives in an online fashion onto different devices. Although the asynchronous aggregation in their proposed algorithm can be naturally applied to randomly inactive devices, the authors do not analyze how their algorithm's convergence is affected by the device inactivity or incompleteness and the data heterogeneity.

In recent years, some attempts have been made to relax the strict training requirements on the participating devices. For example, Tu et al. (2020) study federated learning in a fog network topology with possible data sharing among devices; Yang et al. (2020) incorporate heterogeneity of devices into the design of the learning systems; Nishio and Yonetani (2019) propose a client selection policy that adapts to the change of devices' hardware status. However, these works do not show how the variations in the devices could affect the convergence of training, nor do they incorporate the heterogeneity of user data into the algorithm design.

In (Rizk et al., 2020) and (Wang et al., 2020), the authors reveal that incomplete devices can block the convergence, but they consider neither other dynamic participation patterns such as inactivity, arrivals and departures, nor probabilistic models for uncertain device participation. To relieve the impact of incomplete devices, these works propose similar strategies as our paper by reweighting the contribution of local models. However, they focus mostly on removing the additional bias term originating from heterogeneous device updates, without looking into how this bias is related to the participation frequency of devices and the divergence among them. They also do not compare the proposed methods with alternative extensions of Fe-dAvg. In this work, we model the device participation

as random variables and incorporate them into the convergence analysis, and we compare the convergence rates for three reasonable aggregation schemes.

3 Convergence Analysis

In this section, we establish a convergence bound for federated learning with flexible device participation patterns. Our analysis generalizes the standard FedAvg to incorporate arbitrary aggregation coefficients. In the aggregation step, all devices are counted even if they cannot finish all local epochs. The analysis considers a non-IID data distribution and heterogeneous devices, i.e., some devices can be more stable than the others. We first derive the convergence bound with incomplete and inactive devices in Sections 3.1 to 3.2 and then discuss arrivals and departures in Section 3.3

3.1 Algorithm Description

Suppose there are N devices, where each device k defines a local objective function $F_k(w)$. Here w represents the parameters of the machine learning model to be optimized, and $F_k(w)$ may be defined as the average empirical loss over all data points at device k, as in typical federated learning frameworks (McMahan et al., 2016). The global objective is to minimize $F(w) = \sum_{k=1}^{N} p^k F_k(w)$, where $p^k = \frac{n_k}{n}$, n_k is the number of data points device k owns, and $n = \sum_{k=1}^{N} n_k$. Let w^* be the minimizer of F, and denote by F_k^* the minimum value of F_k . We quantify the degree to which data at each device k is distributed differently than that at other devices as $\Gamma_k = F_k(w^*) - F_k^*$ to capture that data distributions at different devices are non-IID, and let $\Gamma = \sum_{k=1}^{N} p^k \Gamma_k$ as in (Li et al., 2020b).

We consider discrete time steps $t=0,1,\ldots$ Model weights are synchronized when t is a multiple of E, i.e., each round consists of E time steps. Assume there are at most T rounds. For each round (say the τ th round), the following three steps are executed:

- Synchronization: the coordinator broadcasts the latest global weight $w_{\tau E}^{\mathcal{G}}$ to all devices. Each device updates its local weight so that: $w_{\tau E}^{k} = w_{\tau E}^{\mathcal{G}}$
- Local updates: each device runs stochastic gradient descent (SGD) on F_k for $i = 0, \dots, s_{\tau}^k 1$

$$w_{\tau E+i+1}^{k} = w_{\tau E+i}^{k} - \eta_{\tau} g_{\tau E+i}^{k} \tag{1}$$

Here η_{τ} is a staircase learning rate that decays with τ , $0 \leq s_{\tau}^{k} \leq E$ represents the number of local updates this device completes in this round, $g_{t}^{k} = \nabla F_{k}(w_{t}^{k}, \xi_{t}^{k})$ is the stochastic gradient at device

¹While some papers define local epochs and local updates separately, we use them interchangeably in this paper. Both refer to the times (1) is conducted in a global round.

k, and ξ_t^k is a mini-batch sampled from device k's local dataset. We also define $\bar{g}_t^k = \nabla F_k(w_t^k)$ as the full batch gradient at device k, hence $\bar{g}_t^k = \mathbb{E}_{\xi_t^k}[g_t^k]$.

Aggregation: the coordinator aggregates the gradients and generates the next global weight as

$$w_{(\tau+1)E}^{\mathcal{G}} = w_{\tau E}^{\mathcal{G}} + \sum_{k=1}^{N} p_{\tau}^{k} (w_{\tau E + s_{\tau}^{k}} - w_{\tau E}^{\mathcal{G}})$$

$$= w_{\tau E}^{\mathcal{G}} - \sum_{k=1}^{N} p_{\tau}^{k} \sum_{i=0}^{s_{\tau}^{k}} \eta_{\tau} g_{\tau E + i}^{k}$$
(2)

We define that a device k is inactive in round τ if $s_{\tau}^{k}=0$ (i.e., it completes no local updates), and say it is incomplete if $0 < s_{\tau}^{k} < E$. We treat each s_{τ}^{k} as a random variable that can follow an arbitrary distribution. Devices are heterogeneous if they have different distributions of s_{τ}^{k} , and otherwise they are homogeneous. We allow the aggregation coefficients p_{τ}^{k} to vary with τ . In Section \P , we will discuss different schemes of choosing p_{τ}^{k} and their impacts on the convergence.

As a special case, traditional FedAvg assumes all selected devices can complete all E local epochs, so that $s_{\tau}^{k} \equiv E$. Also, FedAvg with full device participation uses fixed aggregation coefficients $p_{\tau}^{k} \equiv p^{k}$, so that the right hand side of (2) can be written as $\sum_{k=1}^{N} p^{k} w_{\tau E}^{k}$, i.e., aggregating gradients is equivalent to aggregating the model parameters directly.

3.2 General Convergence Bound

The analysis relies on the following five assumptions. The first four are standard (Li et al.) 2020b). The last assumption ensures bounded aggregation coefficients and is satisfied by all schemes discussed in Section 4. In Section 5 we experimentally show that our proposed learning algorithm performs well even when some assumptions (like strong convexity) are violated.

Assumption 3.1. F_1, \ldots, F_N are all L-smooth, so that F is also L-smooth.

Assumption 3.2. F_1, \ldots, F_N are all μ -strongly convex, so that F is also μ -strongly convex.

Assumption 3.3. The variance of the stochastic gradients is bounded: $\mathbb{E}_{\xi} ||g_t^k - \bar{g}_t^k||^2 \leq \sigma_k^2, \ \forall k, t.$

Assumption 3.4. The expected squared norm of the stochastic gradients at each local device is uniformly bounded: $\mathbb{E}_{\xi} ||g_t^k||^2 \leq G^2$ for all k and t.

Assumption 3.5. There exists an upper bound $\theta > 0$ for the aggregation coefficient: $p_{\tau}^k/p^k \leq \theta, \forall k$.

Assume the following expectations exist and do not vary with time: $\mathbb{E}[p_{\tau}^k]$, $\mathbb{E}[p_{\tau}^k s_{\tau}^k]$, $\mathbb{E}[(p_{\tau}^k)^2 s_{\tau}^k]$, $\mathbb{E}[(\sum_{k=1}^N p_{\tau}^k - 2)_+(\sum_{k=1}^N p_{\tau}^k s_{\tau}^k)]$ for all rounds τ and devices k, and assume $\mathbb{E}[\sum_{k=1}^N p_{\tau}^k s_{\tau}^k] \neq 0$. Intuitively, this last assumption ensures that some updates are aggregated in each round, otherwise this round can be simply omitted.

Generally, p_{τ}^k 's are functions of s_{τ}^k , and these expectations can be estimated from device histories. Let $z_{\tau} \in \{0,1\}$ indicate the event that the ratio $\mathbb{E}[p_{\tau}^k s_{\tau}^k]/p^k$ does not take the same value for all k. We can obtain the following convergence bound for general p_{τ}^k :

Theorem 3.1. By choosing the learning rate $\eta_{\tau} = \frac{16E}{\mu \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}]} \frac{1}{\tau E + \gamma}$, we can obtain

$$\mathbb{E}\|w_{\tau E}^{\mathcal{G}} - w^*\|^2 \le \frac{M_{\tau}D + V}{\tau E + \gamma} \tag{3}$$

$$\begin{array}{lll} \textit{Here we define } \gamma = \max \left\{ \frac{32E(1+\theta)L}{\mu \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}]}, \frac{4E^{2}\theta}{\mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}]} \right\}, \\ M_{\tau} &= \sum_{t=0}^{\tau-1} \mathbb{E}[z_{t}], \quad D &= \frac{64E\sum_{k=1}^{N} \mathbb{E}[p_{\tau}^{k} s_{\tau}^{k}] \Gamma_{k}}{\mu \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}] \Gamma_{k}}, \\ V &= \max \left\{ \gamma^{2} \mathbb{E} \| w_{0}^{\mathcal{G}} - w^{*} \|^{2}, \left(\frac{16E}{\mu \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}]} \right)^{2} \frac{\mathbb{E}[B_{\tau}]}{E} \right\}, \\ B_{\tau} &= 2(2+\theta)L\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} \Gamma_{k} + \left(2 + \frac{\mu}{2(1+\theta)L} \right) E(E - 1)G^{2} \left(\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} + \theta(\sum_{k=1}^{N} p_{\tau}^{k} - 2) + \sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} \right) \\ &+ 2EG^{2}\sum_{k=1}^{N} \frac{(p_{\tau}^{k})^{2}}{p^{k}} s_{\tau}^{k} + \sum_{k=1}^{N} (p_{\tau}^{k})^{2} s_{\tau}^{k} \sigma_{k}^{2} \end{array}$$

Theorem 3.1 shows that the convergence rate is affected by the aggregation coefficients p_{τ}^{k} 's as they determine M_{τ} , D, and V. From (3), $w_{\tau E}^{\mathcal{G}}$ will eventually converge to a globally optimal solution only if M_{τ} increases sub-linearly with τ . In the original full-participation FedAvg, $p_{\tau}^k s_{\tau}^k \equiv p^k E$, so $z_{\tau} \equiv 0$ and $M_{\tau} \equiv 0$ as per the definitions. Thus, full-participation FedAvg converges according to (3), which is consistent with (Li et al.) 2020b). However, when considering flexible device participation, M_{τ} may increase with τ , which can cause FedAvg to converge to an arbitrary suboptimal point. The magnitude of M_{τ} is determined by the degree of heterogeneity in the device participation, and D is bounded by the non-IID metric Γ_k of local datasets. If M_{τ} increases linearly with τ (e.g., due to device departures), the model will converge to a suboptimal point with the loss bounded by $\frac{\overline{D}}{E}$. As we will see in Section 4.1, by smartly choosing the aggregation coefficients p_{τ}^k , the increase of M_{τ} can be controlled and a convergence to the global optimum can still be established.

While we only show results for s_{τ}^{k} whose distributions are static with time, Theorem 3.1 can be easily extended to time-varying distributed s_{τ}^{k} by replacing the corresponding expectations of $p_{\tau}^{k}s_{\tau}^{k}$ and $(p_{\tau}^{k})^{2}s_{\tau}^{k}$ with their minimum or maximum expectations over τ .

3.3 Shifts in the Global Objective

Recall the global objective is $F(w) = \sum_{k \in \mathcal{C}} p^k F_k(w)$, i.e., an average of local objectives for participating devices \mathcal{C} . A well trained model w^* is expected to perform well on all data points generated by devices in \mathcal{C} . In the presence of departing and arriving devices,

 $\mathcal C$ may shrink or expand dynamically during the training. The global objective thus varies accordingly. For example, after admitting an incoming device l with n_l data points: $\tilde{\mathcal C} \leftarrow \mathcal C + \{l\}$, the global objective becomes $\tilde{F}(w) = \tilde{p}^l F_l(w) + \sum_{k \in \mathcal C} \tilde{p}^k F_k(w)$, where $\tilde{p}^k = \frac{n_k}{n_C + n_l}$. The model \tilde{w}^* fully trained with this objective is then applicable to the new data from device l. We formally define objective shift as the process of changing the global objective, and the applicability of the trained model, by adding or removing devices from $\mathcal C$.

The following theorem bounds the offset between the global optima due to the objective shift. As we can intuitively expect, the difference reduces when the data becomes more IID $(\Gamma_l \to 0)$, and when the departing/arriving device owns fewer data points $(n_l \to 0)$:

Theorem 3.2. Suppose a device l arrives/departs, and let n be the total number of data points originally. Consider the objective shift $F \to \tilde{F}$, $w^* \to \tilde{w}^*$. Let $\tilde{\Gamma}_k = F_k(\tilde{w}^*) - F_k^*$ quantify the degree of non-IID with respect to the new objective. Then in the arrival case

$$\|w^* - \tilde{w}^*\| \le \frac{2\sqrt{2L}}{\mu} \frac{n_l}{n+n_l} \sqrt{\Gamma_l} \tag{4}$$

and in the departure case

$$\|w^* - \tilde{w}^*\| \le \frac{2\sqrt{2L}}{\mu} \frac{n_l}{n} \sqrt{\tilde{\Gamma}_l} \tag{5}$$

Objective shift is mandatory when a new device (say device l) arrives: Unless $F_l \equiv F$ (which is highly unlikely), incorporating updates from l will always move F(w) away from F^* . The best strategy without objective shift is then not to aggregate updates from l, and thus not to admit l into the learning process in the first place. In contrast, objective shift is optional when devices depart: we can keep the original objective F even if we will no longer receive updates from a departing device, if doing so yields smaller training loss.

Suppose an objective shift occurs at τ_0 . The remainder of the training is then equivalent to starting over from $w_{\tau_0 E}^{\mathcal{G}}$ but converging towards the new objective \tilde{w}^* . Combining Theorems 3.1 and 3.2 we can obtain the following convergence bound after the objective shifts:

Corollary 3.2.1. Assume the objective shifts at τ_0 with $\mathbb{E}\|w_{\tau_0 E}^{\mathcal{G}} - w^*\|^2 \leq \Delta_{\tau_0}$. By increasing the learning rate back to $\eta_{\tau} = \frac{16E}{\mu \mathbb{E}[\sum_{k=1}^N p_{\tau}^k s_{\tau}^k]} \frac{1}{(\tau - \tau_0)E + \gamma}$ for $\tau > \tau_0$, the convergence to the new objective can be bounded by

$$\mathbb{E}\|w_{\tau E}^{\mathcal{G}} - \tilde{w}^*\|^2 \le \frac{\tilde{M}_{\tau}\tilde{D} + \tilde{V}}{(\tau - \tau_0)E + \tilde{\gamma}}$$
 (6)

Here $\tilde{M}_{\tau}, \tilde{D}, \tilde{V}, \tilde{\gamma}$ are defined analogously to M_{τ}, D , V, γ but they respectively include/exclude the arriving/departing device. The first term in \tilde{V} equals $\tilde{\gamma}^2(\sqrt{\Delta_{\tau_0}} + \|w^* - \tilde{w}^*\|)^2 = O\left(\frac{V}{\tau_0 E + \gamma} + \Gamma_l\right)$.

The increase of the learning rate after the objective shift is necessary. Intuitively, if the shift happens at a large time τ_0 when $w_{\tau_0 E}^{\mathcal{G}}$ is close to the old optimal w^* and η_{τ_0} is close to zero, the learning rate used in Theorem 3.1 will be too small to steer the model to the new optimum, since $\|w_{\tau_0 E}^{\mathcal{G}} - \tilde{w}^*\| \approx \|w^* - \tilde{w}^*\|$.

Comparing (3) and (6), an objective shift yields an one-time increase in the loss, which forces us to take actions when confronting departures and arrivals. In the case of device departure, it is possible that retaining the old objective can result in a smaller training loss compared to doing a shift. In this situation, the trained model is still applicable to data of the departing device. In the arrival case, though objective shift is mandatory, we can still accelerate the training by a "fast-reboot", applying extra gradient updates from the arriving device.

We will discuss in Section 4.2 the fast-reboot method for the arrival case, and in Section 4.3 the decision of model applicability for the departure case.

4 Main Results

Based on the convergence analysis in Section 3 in this section, we present corollaries that can guide operators in reacting to flexible device participation.

4.1 Debiasing on Incomplete Aggregation

According to Theorem 3.1 the convergence bound is controlled by the expectation of p_{τ}^{k} and its functions. Below we discuss three plausible schemes of choosing p_{τ}^{k} , and compare their convergence rates in Table 1.

- Scheme A: Only aggregate parameters from devices that complete all E local epochs, with aggregation coefficient p^k_τ = Np^k/K_τ q^k_τ, where K_τ is the number of complete devices, q^k_τ ∈ {0,1} denotes if client k is complete. If K_τ = 0, this round is discarded.
- Scheme B: Allow clients to upload incomplete work (with $s_{\tau}^k < E$ updates), with fixed aggregation coefficient $p_{\tau}^k = p^k$.
- Scheme C: Accept incomplete works as in Scheme B, with adaptive $p_{\tau}^k = \frac{E}{s^k} p^k$, or $p_{\tau}^k = 0$ if $s_{\tau}^k = 0$.

Schemes A and B are natural extensions of *FedAvg*. Scheme C assigns a greater aggregation coefficient to devices that complete fewer local epochs. Though this idea seems counter-intuitive, as fewer local updates might lead to less optimal parameters (cf. Table 1), it turns out to be the only scheme that guarantees convergence when device participation is heterogeneous.

Corollary 4.0.1. Let K_{τ} be the number of devices that run all E epochs, I_{τ} indicate the appearance of any inactive devices in round τ , and write $\bar{\sigma}_N^2 \equiv \sum_k^N (p^k \sigma_k)^2$.

Table I gives the convergence rates of Schemes A, B, C when device updates may be incomplete and inactive.

Table 1: Convergence rates with incomplete and inactive devices. The bound for Scheme A assumes there is at least one complete device $(K_{\tau} \neq 0)$, and those for Schemes B, C assume s_{τ}^k is not trivially zero $(\mathbb{E}[s_{\tau}^k] \neq 0)$. While the three schemes have similar performance in the homogeneous setting, Schemes A and B fail to converge to the global optimum even assuming all devices are active. Scheme C works if inactive devices do not occur in every round $(\sum_t I_t < O(\tau))$.

	Homogeneous	Heterogeneous
A	$O\left(\frac{\mathbb{E}\left[\frac{N^2}{K_\tau}\right] + \bar{\sigma}_N^2 + \Gamma}{\tau}\right)$	$\leq \frac{D}{E}$
В	$O\left(\frac{\bar{\sigma}_N^2 + \Gamma}{\tau \mathbb{E}[s_{\tau}]}\right)$	$\leq \frac{D}{E}$
C	$O\left(\frac{\bar{\sigma}_N^2 + \Gamma}{\tau(\mathbb{E}[1/s_\tau])^{-1}}\right)$	$O\left(\frac{\sum_{t=0}^{\tau-1} I_t D + \sum_{k}^{N} (p^k \sigma_k)^2 \mathbb{E}\left[\frac{1}{s_{\tau}^k}\right] + \Gamma}{\tau}\right)$

The reason for enlarging the aggregation coefficients in Scheme C can be understood by observing from (2) that increasing p_{τ}^{k} is equivalent to increasing the learning rate of device k. Thus, by assigning devices that run fewer epochs a greater aggregation coefficient, these devices effectively run further in each local step, compensating for the additional epochs other devices completed. As shown in Figure 1. Scheme C ensures an unbiased gradient after aggregation, while Schemes A and B will favor devices that run more epochs. Ideally, allowing devices to adapt learning rates by themselves would effectively lead to the same result. However, when a device is running local updates, it may not yet know or be able to estimate the number of local epochs it will complete. In contrast, centralized intervention can make accurate adjustments a posterior.

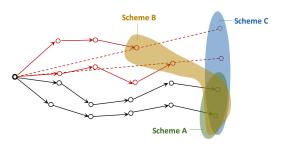


Figure 1: Snapshot of one aggregation round. The bottom two devices completed all E=5 local epochs, while the top two completed only 3 and 4 epochs. Scheme C enlarges the incomplete gradients by respectively 5/3 and 5/4 and produces unbiased aggregation results. Aggregations with Schemes A and B are biased towards devices that run more epochs.

Table I also reveals how the following system and

statistical factors affect the convergence asymptotically:

- The non-IID metric Γ is the major obstacle of convergence in the homogeneous case. In the heterogeneous setting, the D term (which grows with Γ) dominates the training loss. It controls the maximum non-diminishing loss D/E of Scheme A and B, and decelerates the training of Scheme C in the presence of inactive devices.
- Devices' activeness s_{τ}^{k} and K_{τ} contribute inversely to the training loss: The more devices participate, the faster the loss decays. When inactivity occurs frequently, Scheme C cannot converge either. E.g., if a device never responds to the coordinator (so $I_{t} \equiv 1$), its training loss can never converge to zero.
- The variance $\bar{\sigma}_N, \sigma_k$ in the stochastic gradient descent algorithm slows down the training as expected.

4.2 Fast-rebooting on Arrivals

Intuitively, when a device l arrives, \tilde{w}^* will be "dragged" towards its local optimum w_l^* . The gradients from device l may thus encode more information about the new optimum \tilde{w}^* compared to those from the other devices. Thus, by adding an extra update $-\delta^l \nabla F_l(w^{\mathcal{G}}), \delta^l > 0$ to the gradient aggregation, it is likely that w can move closer to \tilde{w}^* , allowing the training to fast-reboot from the point of arrivals. However, as shown in Figure 2 this intuition may not hold: it is also possible that $-\delta^l \nabla F_l(w^{\mathcal{G}})$ ends up driving $w^{\mathcal{G}}$ away from \tilde{w}^* . In fact, the success of this method is determined by the distance $b = ||w^{\mathcal{G}} - w^*||$. When b is small, it is highly likely the extra update can rapidly reboot the training. We formalize this statement in Corollary 4.0.2

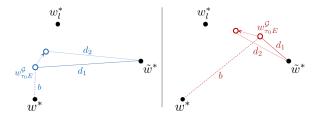


Figure 2: Left: when the distance to the old optimum $b = ||w^{\mathcal{G}} - w^*||$ is small, applying an extra update to $w^{\mathcal{G}}$ following the direction $-\nabla F_l(w^{\mathcal{G}})$ moves it closer to \tilde{w}^* $(d_2 < d_1)$. Right: for a large b, the extra update may on the contrary enlarge this distance $(d_2 > d_1)$.

Corollary 4.0.2. Assume $\nabla F(w)$ is continuous, and $0 < \|\nabla F(w)\|_2, \|\nabla^2 F(w)\|_2 \le W$ for any w (The latter is the induced l_2 norm for matrices). Let $w' = w - \delta^l \nabla F_l(w)$, then there exists a $\delta^l > 0$ such that $\|w' - \tilde{w}^*\| < \|w - \tilde{w}^*\|$ if w satisfies

$$||w - w^*|| < \frac{\tilde{F}(w^*) - \tilde{F}(\tilde{w}^*)}{\left(\frac{2\sqrt{2L}}{\mu}\tilde{p}^l\sqrt{\Gamma_l} + 1\right)\tilde{p}^lW}$$
 (7)

(7) defines a sphere around the original global optimum w^* within which the extra update helps fast-reboot. The radius of the sphere depends on the divergence between the new (arriving) and old data points. Generally, the longer the training has elapsed, the closer the global model is to w^* . Thus, the extra updating works best for devices that arrive late in the training.

When applied in practice, the extra updating can be conducted on-the-fly, by augmenting the aggregation coefficient of the arriving device so that $p_{\tau}^{l} = p^{l} + \delta^{l}$. Furthermore, the distance b can be estimated by the gradient norm with respect to the original objective.

As the name suggests, fast-reboot only accelerates the training for a certain duration after the device arrives. In fact, if there are no future interrupts, models with or without fast-rebooting eventually converge to the same global optimum. Nevertheless, fast-reboot is still beneficial if there is insufficient training time remaining (e.g., a device arrives near the end of the training).

4.3 Redefining Applicability on Departures

As is discussed in Section 3.3 when a device leaves, we need to redefine the applicability of the trained model. Namely, one can decide to either exclude this departing device and shift the objective, or keep including it and stick to the old objective. The decision depends on the time at which the device leaves. When including the device as a part of the global objective, from 3, since $M_{\tau} = \tau - \tau_0$ from then on, the training loss will always exceed a structural bias D/E. In contrast, if the device is excluded and the model is trained with a shifted global objective, there will be an immediate increase in the convergence bound as in Theorem 3.2 But afterwards, the bound will decrease and eventually the parameters will converge to the new global optimum.

Assume a device leaves at $\tau_0 < T$ and there are no subsequent arrivals/departures. Let $f_0(\tau)$ be the convergence bound if we include the device, and $f_1(\tau)$ be the bound if it is excluded. We can obtain $f_0(\tau) = \frac{(\tau - \tau_0)D + V}{\tau E + \gamma}$, $f_1(\tau) = \frac{\tilde{V}}{(\tau - \tau_0)E + \tilde{\gamma}}$. Here \tilde{M}_{τ} , \tilde{V}_{τ} , $\tilde{\gamma}$ are defined analogously to M_{τ} , V_{τ} , γ but they exclude the departing device. A device is excluded if by doing so, a smaller training loss can be obtained at the deadline T, which is summarized in the following corollary:

Corollary 4.0.3. Excluding a device that departs at τ_0 leads to smaller training loss if

$$\min_{\tau \ge \tau_0} f_0(\tau) \ge f_1(T) \tag{8}$$

Further assume $\tilde{\gamma} = \gamma$, and \tilde{V} is dominated by its first term so that $\tilde{V} = \frac{V}{\tau_0 E + \gamma} + \Gamma_l$. (8) then becomes

$$T - \tau_0 \ge O\left(\sqrt{\Gamma_l \tau_0}\right) \tag{9}$$

From (9), when the remaining training time $T - \tau_0$ is at least $O(\sqrt{\Gamma_l \tau_0})$, applying the trained model to the departing device becomes less promising. It is thus better to exclude it and shift the objective. As we can expect, the bound grows with Γ_l , since the non-IID contribution from the departing device increases the initial \tilde{V} . As τ_0 increases, the learning rate without shift gets smaller, mitigating the increase of the training loss from departing devices.

5 Experiments

In this section, we experimentally evaluate Section 4s results. Due to the limitations on hardware resources, the training process is performed in computer simulations. To ensure the simulation is consistent with the real learning environment, we use real-world traces to represent the participation patterns of simulated devices. We present our experiment setup in Section 5.1, and verify our theory results in Sections 5.2 - 5.4.

5.1 Experiment Setup

We create various data traces to represent the heterogeneous participation patterns of local devices. We set up a simple federated learning experiment with five Raspberry PIs as workers, and a desktop server as the coordinator. Each PI has a training process that runs the original FedAvg algorithm, and a competitor process doing CPU-intensive work simultaneously. We manually tune the workload of the competitor process so that it takes up 0%, 30%, 50%, 70%, 90% of the PI's CPU resources, simulating different device configurations in federated learning. Under the five settings, for each round, we record the percentage of required epochs the PI ends up submitting before a preset, fixed deadline. Due to the default load-balancing behavior of the operating system's CPU scheduler, these traces do not contain zero epochs (i.e. inactive cases). To generate inactive device participation patterns, we create another set of three traces with respectively low, medium and high bandwidth. Devices can thus be inactive due to weak transmission. Table 2 shows the mean and standard deviation of the percentage of epochs completed for each trace. In the following experiments, each simulated device is randomly assigned a trace. For each aggregation round τ , it randomly samples from its trace to obtain the number of local epochs s_{τ}^{k} .

Three datasets are used in this paper: MNIST (Le-Cun et al., 1998), EMNIST (Cohen et al., 2017) and SYNTHETIC(α, β) (Li et al., 2018). We build a two-layer MLP model and a two-convolution-layer CNN model respectively for MNIST and EMNIST, both models are defined by McMahan et al. (2016). For SYNTHETIC(α, β), we use an ordinary logistic re-

Table 2: The means and standard deviations for the percentage of required local epochs actually submitted to the coordinator during the federated training. The first five traces do not contain inactive cases.

Name								
Mean								
Stdev	0	14.8	11.3	11.7	14.8	23.3	22.3	18.3

gression model. All models use the vanilla SGD as local optimizers, with batch sizes of 10 for MNIST and EMNIST, and 20 for SYNTHETIC. When generating non-IID data, we sort the MNIST and EMNIST data by labels so that each device is assigned data from one label chosen uniformly at random. For SYNTHETIC(α, β), we vary the parameters α, β from 0 to 1. The larger α, β are, the less IID the dataset becomes. We use the staircase learning rate $\eta_{\tau} = \eta_0/\tau$ as adopted in our convergence analysis. The initial η_0 is 2e-3 for MNIST, 5e-4 for EMNIST, and 1 for SYNTHETIC(α, β). Unless otherwise noted, the number of samples at each device follows the Type-I Pareto distribution with the Pareto index of 0.5.

5.2 Comparison of Aggregation Schemes

We first examine the effects of the device heterogeneity and the non-IID data distributions on the convergence for each aggregation scheme. We conduct eight sets of experiments where we incrementally increase the number of participation traces to reflect the increasing heterogeneity in device participation. For SYNTHETIC, we use $\alpha = \beta = 0$ for the IID case, and $\alpha = \beta = 1$ for the non-IID case. We train on 100 devices for MNIST, 62 devices for EMNIST (by merge), and 50 devices for SYNTHETIC(α, β). Table 3 records the differences in the test accuracies between different aggregation schemes after 200 global epochs. The typical convergence process is depicted in Figure 3

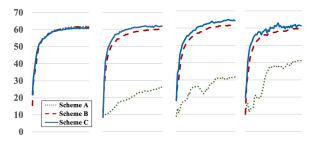


Figure 3: Test accuracy for non-IID EMNIST. Plots from left to right correspond to $|\mathcal{T}| = 1, 3, 5, 8$. (increasing device heterogeneity)

As we can see, Scheme C yields the best test accuracy on average. Compared to Schemes A and B, it achieves higher accuracy when devices get more heterogeneous and less IID. This is consistent with our loss bounds in Table $\[\]$ since Schemes A and B fail to converge to the global optimum in the heterogeneous case with non-IID data. On the other hand, Scheme A performs extremely badly with large $|\mathcal{T}|$. This is because the last few traces contain very few complete rounds, significantly increasing $\mathbb{E}[1/K_{\tau}]$. Noteworthily, Scheme C is no different from, or even worse than Scheme B in more homogeneous settings, this is consistent with Table $\[\]$ since $\frac{1}{\mathbb{E}[s_{\tau}]} \leq \mathbb{E}[\frac{1}{s_{\tau}}]$. When the traces contain inactive devices $(|\mathcal{T}| \geq 6)$, Scheme C becomes less stable due to the variance introduced by I_t in Corollary $\[\]$ 4.0.1

Table 3: The % improvement in the test accuracies of Scheme B w.r.t. Schemes A(left numbers) and Scheme C w.r.t. Scheme B (right numbers). $|\mathcal{T}| = j$ represents using the first j traces in Table 2.

(a) MNIST Data

$ \mathcal{T} $	1	2	3	4	
IID	-0.6 0.3	1.9 0.1	5.6 0.1	8.3 0.7	
NIID	0.2 -0.3	9.5 1.8	19.3 1.6	33.8 3.3	
17		-	1	8	
/	5	6	7	8	
IID	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{ c c c c c }\hline 7 \\ 5.8 & 1.9 \\ \hline \end{array}$	8 11.8 2.4	

(b) EMNIST Data

$ \mathcal{T} $	1	2	3	4	
IID	0.7 -0.6	0.9 0.1	4.2 0.7	4.8 1.0	
NIID	-0.1 -0.7	17.0 -2.0	34.2 1.8	37.9 4.8	
	1			8	
$ \mathcal{T} $	5	6	7	8	
$\frac{ \mathcal{T} }{ \mathbf{IID} }$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{ c c c c c }\hline 7 & & & \\\hline 4.0 & 1.5 & & \\\hline \end{array}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	

(c) SYNTHETIC Data

$ \mathcal{T} $	1	2	3	4	
IID	-0.6 0.5	2.1 0.1	6.6 0.0	9.0 0.7	
NIID	0.1 -0.4	9.6 1.5	22.2 1.8	38.2 3.2	
$\overline{ \mathcal{T} }$	5	6	7	8	
$\frac{ \mathcal{T} }{ ext{IID}}$	5 11.6 3.0	6 16.4 2.5	7 6.3 1.9	8 14.4 2.8	

5.3 Effectiveness of Fast-Reboot

We now investigate the effectiveness of the fast-reboot method described in Section 4.2. The experiments involve N-1 existing devices, and the arriving device joins at τ_0 . As is discussed in Section 4.2, the method makes no difference when data distribution is IID. We thus only consider non-IID cases. We set N=10 for MNIST and EMNIST (balanced) and N=30 for

SYNTHETIC(1,1). To avoid the interference brought by inactive devices, for this experiment we only use the first five traces in Table 2 and we adopt Scheme C as the aggregation method. All devices are given the same number of samples for fair comparison.

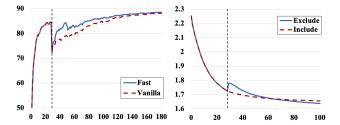


Figure 4: Evolution of the test accuracy (left) and loss (right) under device arrival (left, MNIST) and departure (right, SYNTHETIC) cases. The dashed vertical lines indicate the arriving (departing) time τ_0 . After τ_0 , except for the "include" option, models are tested with new datasets that include (exclude) holdout data from the arriving (departing) device.

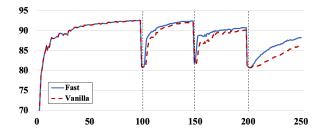


Figure 5: Test accuracy with and without fast-reboot for multiple arrivals for non-IID MNIST. The test dataset is updated every time a new device arrives to include its holdout data. The vertical dashed lines indicate the time the device arrives.

Table 4: The number of global epochs after the arriving time τ_0 until the test accuracy bounces back to that at $\tau_0 - 1$. Left: fast reboot. Right: vanilla reboot.

$ au_0$	10	30	50	70	
MNIST	4 4	22 27	55 63	59 66	
EMNIST	4 3	11 12	14 19	21 24	
SYNTHETIC	1 1	4 6	7 12	3 8	

When the device arrives, we increase the learning rate to $\eta_0/(\tau-\tau_0)$. The aggregation coefficient of the arriving device l is boosted to $p_{\tau}^l=3p^l$ initially, and decays to p^l by $O(\tau^{-2})$. Table 4 records the number of global epochs it takes to recover to the accuracy level before the arrival. Fast-reboot consistently achieves faster rebound, and works better for late arrivals as we expect. EMNIST-CNN enjoys less improvement from fast-reboot because CNN models converge more slowly

than MLP and logistic regression models. Thus, at the moment new devices arrive, EMNIST models have not fully converged to the old optima, degrading the effectiveness of fast-reboot as per Corollary 4.0.2 The typical fast-reboot process is shown in Figure 4.

Next we study the situation when multiple devices arrive in a row. Figure 5 shows the training process for MNIST data. Every time a device arrives, we increase the learning rate as per Corollary 3.2.1 Initially, seven devices are in the training. After 100 global epochs, the remaining three devices arrive at 50 epoch intervals, without waiting for the model to fully converge. From Figure 5 the fast-reboot trick accelerates the convergence for every device arrival.

5.4 Model Applicability upon Departures

The right plot in Figure 4 shows the typical change of the test loss after the device departs. We use the same setting as in Section 5.3. As is predicted in Section 4.3 an objective shift ('exclude') initially increases the test loss. But eventually, the two curves cross and excluding the device becomes more beneficial.

Table 5 summarizes the number of global epochs it takes for the curves to cross with SYNTHETIC(α, β). As we can see, the values increase with τ_0 and the non-IID metric (α, β), confirming Corollary 4.0.3.

Table 5: The number of global epochs after the departing time τ_0 until the test losses coincide for including and excluding options. The rows correspond to three choices of parameters (α, β) in SYNTHETIC (α, β) .

$\overline{ au_0}$	10	15	20	25	30	35	40	45	50
(.1,.1)	2	5	3	3	9	3	10	26	40
(.5,.5)	1	3	9	14	13	7	12	36	34
(1., 1.)	10	9	27	18	34	17	28	62	77

6 Conclusion and Future Work

This paper extends the federated learning paradigm to incorporate more flexible device participation. The analysis shows that incomplete local device updates can be utilized by scaling the corresponding aggregation coefficients, and a mild degree of device inactivity will not impact the convergence. Further investigation reveals how the convergence relates to heterogeneity in both the data and the device participation. The paper also proposes techniques to fast-reboot the training after new devices arrive, and provides an analytical criterion on when to exclude a departing device. In the future work, we will analyze groups of arrivals or departures, and investigate the possibility for users to dynamically update their datasets during the training.

Acknowledgements

This research was partially supported by the CMU CyLab IoT Initiative, and NSF CNS-1909306, CNS-1751075.

References

- Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir Ivanov, Chloe Kiddon, Jakub Konečný, Stefano Mazzocchi, H Brendan McMahan, et al. Towards federated learning at scale: System design. arXiv preprint arXiv:1902.01046, 2019.
- Sebastian Caldas, Jakub Konečny, H Brendan McMahan, and Ameet Talwalkar. Expanding the reach of federated learning by reducing client resource requirements. arXiv preprint arXiv:1812.07210, 2018.
- Yujing Chen, Yue Ning, and Huzefa Rangwala. Asynchronous online federated learning for edge devices. arXiv preprint arXiv:1911.02134, 2019.
- Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. In 2017 International Joint Conference on Neural Networks (IJCNN), pages 2921–2926. IEEE, 2017.
- Luca Corinzia and Joachim M Buhmann. Variational federated multi-task learning. arXiv preprint arXiv:1906.06268, 2019.
- Georgios Damaskinos, Rachid Guerraoui, Anne-Marie Kermarrec, Vlad Nitu, Rhicheek Patra, and Francois Taiani. Fleet: Online federated learning via staleness awareness and performance prediction. arXiv preprint arXiv:2006.07273, 2020.
- Pengchao Han, Shiqiang Wang, and Kin K Leung. Adaptive gradient sparsification for efficient federated learning: An online learning approach. arXiv preprint arXiv:2001.04756, 2020.
- Kevin Hsieh, Amar Phanishayee, Onur Mutlu, and Phillip B Gibbons. The non-iid data quagmire of decentralized machine learning. arXiv preprint arXiv:1910.00189, 2019.
- Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Keith Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. arXiv preprint arXiv:1912.04977, 2019.
- Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. arXiv preprint arXiv:1610.05492, 2016.

- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11): 2278–2324, 1998.
- Li Li, Yuxi Fan, Mike Tse, and Kuo-Yi Lin. A review of applications in federated learning. *Computers & Industrial Engineering*, page 106854, 2020a.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. arXiv preprint arXiv:1812.06127, 2018.
- Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedayg on non-iid data. In *In Proc. of ICLR*, 2020b.
- Lumin Liu, Jun Zhang, S. H. Song, and Khaled B. Letaief. Client-edge-cloud hierarchical federated learning. https://arxiv.org/abs/1905.06641.
- H Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, et al. Communication-efficient learning of deep networks from decentralized data. arXiv preprint arXiv:1602.05629, 2016.
- Takayuki Nishio and Ryo Yonetani. Client selection for federated learning with heterogeneous resources in mobile edge. In *ICC 2019-2019 IEEE International Conference on Communications (ICC)*, pages 1–7. IEEE, 2019.
- Elsa Rizk, Stefan Vlaski, and Ali H Sayed. Dynamic federated learning. arXiv preprint arXiv:2002.08782, 2020.
- Felix Sattler, Simon Wiedemann, Klaus-Robert Müller, and Wojciech Samek. Robust and communication-efficient federated learning from non-iid data. *IEEE transactions on neural networks and learning systems*, 2019.
- Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S Talwalkar. Federated multi-task learning. In Advances in Neural Information Processing Systems, pages 4424–4434, 2017.
- Yuwei Tu, Yichen Ruan, Su Wang, Satyavrat Wagle, Christopher G Brinton, and Carlee Joe-Wang. Network-aware optimization of distributed learning for fog computing. arXiv preprint arXiv:2004.08488, 2020.
- Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. arXiv preprint arXiv:2007.07481, 2020.
- Shiqiang Wang, Tiffany Tuor, Theodoros Salonidis, Kin K. Leung, Christian Makaya, Ting He, and Kevin Chan. Adaptive federated learning in resource constrained edge computing systems. https://arxiv.org/abs/1804.05271, 2019.

- Chengxu Yang, QiPeng Wang, Mengwei Xu, Shangguang Wang, Kaigui Bian, and Xuanzhe Liu. Heterogeneity-aware federated learning. arXiv preprint arXiv:2006.06983, 2020.
- Xiaoxi Zhang, Jianyu Wang, Gauri Joshi, and Carlee Joe-Wong. Machine learning on volatile instances. arXiv preprint arXiv:2003.05649, 2020.
- Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning on non-iid data. https://arxiv.org/abs/1806_00582, 2018.

A Proof of Theorems and Corollaries

A.1 Proof of Theorem 3.1

A.1.1 Equivalent View

For ease of the analysis, we introduce for each client k and each global round τ a sequence of virtual variables $\alpha_{\tau E}^k, \alpha_{\tau E+1}^k, \dots, \alpha_{(\tau+1)E-1}^k$. Here each $\alpha_t^k \in \{0,1\}$ and $\sum_{i=0}^E \alpha_{\tau E+i}^k = s_{\tau}^k$. Since s_{τ}^k is a random variable, α_t^k 's are also random variables, and the distributions of α_t^k 's determine the distribution of s_{τ}^k . For example, if $\alpha_t^k \stackrel{iid}{\sim} \text{Bernoulli}(p)$, then $s_{\tau}^k \sim \text{Bin}(E,p)$. In general, we do not make any assumption on the distributions and correlations of α_t^k 's. Our results are thus valid for any realization of s_{τ}^k .

With the definition of α_t^k 's, we can rewrite (1) (2) as:

$$w_{\tau E+i+1}^{k} = w_{\tau E+i}^{k} - \eta_{\tau} g_{\tau E+i}^{k} \alpha_{\tau E+i}^{k}$$
(10)

$$w_{(\tau+1)E}^{\mathcal{G}} = w_{\tau E}^{\mathcal{G}} - \sum_{k=1}^{N} p_{\tau}^{k} \sum_{i=0}^{E} \eta_{\tau} g_{\tau E+i}^{k} \alpha_{\tau E+i}^{k}$$

$$\tag{11}$$

Note that $w_t^{\mathcal{G}}$ is visible only when t is a multiple of E. To generalize it to arbitrary t, we define \bar{w}_t such that $\bar{w}_0 = w_0^{\mathcal{G}}$, and

$$\bar{w}_{\tau E+i+1} = \bar{w}_{\tau E+i} - \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} g_{\tau E+i}^{k} \alpha_{\tau E+i}^{k}$$
(12)

Note that $\bar{w}_{\tau E+i} = \sum_{k=1}^{N} p_{\tau}^{k} w_{\tau E+i}^{k}$ only if $\sum_{k=1}^{N} p_{\tau}^{k} = 1$, which generally does not hold.

Lemma A.1. For any τ , $\bar{w}_{\tau E} = w_{\tau E}^{\mathcal{G}}$.

Proof. We will prove by induction. By definition, $\bar{w}_0 = w_0^{\mathcal{G}}$. Suppose $\bar{w}_{\tau E} = w_{\tau E}^{\mathcal{G}}$, then

$$\bar{w}_{(\tau+1)E} = \bar{w}_{(\tau+1)E-1} - \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} g_{(\tau+1)E-1}^{k} \alpha_{(\tau+1)E-1}^{k}$$

$$= \dots = \bar{w}_{\tau E} - \sum_{i=0}^{E-1} \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} g_{\tau E+i}^{k} \alpha_{\tau E+i}^{k}$$

$$= w_{\tau E}^{\mathcal{G}} - \sum_{k=1}^{N} p_{\tau}^{k} \sum_{i=0}^{E-1} \eta_{\tau} g_{\tau E+i}^{k} \alpha_{\tau E+i}^{k} = w_{(\tau+1)E}^{\mathcal{G}}$$

$$(13)$$

Thus, in the following analysis we will just use \bar{w}_t to denote the global weight.

A.1.2 Key Lemmas

We first present a couple of important lemmas:

Lemma A.2.

$$\mathbb{E}_{\xi} \| \sum_{k=1}^{N} p_{\tau}^{k} (g_{t}^{k} - \bar{g}_{t}^{k}) \|^{2} \le \sum_{k=1}^{N} (p_{\tau}^{k})^{2} \sigma_{k}^{2}$$
(14)

Proof.

$$\|\sum_{k=1}^{N} p_{\tau}^{k} (g_{t}^{k} - \bar{g}_{t}^{k})\|^{2} = \sum_{k=1}^{N} \|p_{\tau}^{k} (g_{t}^{k} - \bar{g}_{t}^{k})\|^{2} + \sum_{j \neq k} p_{\tau}^{k} p_{\tau}^{j} \langle g_{t}^{k} - \bar{g}_{t}^{k}, g_{t}^{j} - \bar{g}_{t}^{j} \rangle$$

$$(15)$$

Since each client is running independently, the covariance

$$\mathbb{E}_{\xi} \langle g_t^k - \bar{g}_t^k, g_t^j - \bar{g}_t^j \rangle = 0 \tag{16}$$

Thus,

$$\mathbb{E}_{\xi} \| \sum_{k=1}^{N} p_{\tau}^{k} (g_{t}^{k} - \bar{g}_{t}^{k}) \|^{2} = \sum_{k=1}^{N} \mathbb{E}_{\xi} \| p_{\tau}^{k} (g_{t}^{k} - \bar{g}_{t}^{k}) \|^{2} \le \sum_{k=1}^{N} (p_{\tau}^{k})^{2} \sigma_{k}^{2}$$

$$(17)$$

Lemma A.3. For $i = 0, \dots, E-1$ and all τ, k

$$\mathbb{E}_{\xi}\left[\sum_{k=1}^{N} p_{\tau}^{k} \|\bar{w}_{\tau E+i} - w_{\tau E+i}^{k}\|^{2}\right] \leq (E-1)G^{2} \eta_{\tau}^{2} \left(\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} + (\sum_{k=1}^{N} p_{\tau}^{k} - 2)_{+} \sum_{k=1}^{N} \frac{(p_{\tau}^{k})^{2}}{p^{k}} s_{\tau}^{k}\right)$$
(18)

Proof. Note that $w_{\tau E}^k = \bar{w}_{\tau E}$ for all k.

$$\|\bar{w}_{\tau E+i} - w_{\tau E+i}^{k}\|^{2} = \|(\bar{w}_{\tau E+i} - \bar{w}_{\tau E}) - (w_{\tau E+i}^{k} - \bar{w}_{\tau E})\|^{2}$$

$$= \|\bar{w}_{\tau E+i} - \bar{w}_{\tau E}\|^{2} - 2\langle \bar{w}_{\tau E+i} - \bar{w}_{\tau E}, w_{\tau E+i}^{k} - \bar{w}_{\tau E} \rangle + \|w_{\tau E+i}^{k} - \bar{w}_{\tau E}\|^{2}$$
(19)

From (10) (12),

$$\sum_{k=1}^{N} p_{\tau}^{k} w_{\tau E+i}^{k} = \sum_{k=1}^{N} p_{\tau}^{k} w_{\tau E+i-1}^{k} - \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} g_{\tau E+i-1}^{k} \alpha_{\tau E+i-1}^{k}$$

$$= \sum_{k=1}^{N} p_{\tau}^{k} w_{\tau E+i-1}^{k} + \bar{w}_{\tau E+i} - \bar{w}_{\tau E+i-1}$$

$$= \cdots = \sum_{k=1}^{N} p_{\tau}^{k} w_{\tau E}^{k} + \bar{w}_{\tau E+i} - \bar{w}_{\tau E}$$

$$(20)$$

Thus,

$$-2\sum_{k=1}^{N} p_{\tau}^{k} \langle \bar{w}_{\tau E+i} - \bar{w}_{\tau E}, w_{\tau E+i}^{k} - \bar{w}_{\tau E} \rangle$$

$$= -2 \langle \bar{w}_{\tau E+i} - \bar{w}_{\tau E}, \sum_{k=1}^{N} p_{\tau}^{k} w_{\tau E}^{k} + \bar{w}_{\tau E+i} - \bar{w}_{\tau E} - \sum_{k=1}^{N} p_{\tau}^{k} \bar{w}_{\tau E} \rangle$$

$$= -2 \|\bar{w}_{\tau E+i} - \bar{w}_{\tau E}\|^{2}$$
(21)

$$\sum_{k=1}^{N} p_{\tau}^{k} \|\bar{w}_{\tau E+i} - w_{\tau E+i}^{k}\|^{2} = (\sum_{k=1}^{N} p_{\tau}^{k} - 2) \|\bar{w}_{\tau E+i} - \bar{w}_{\tau E}\|^{2} + \sum_{k=1}^{N} p_{\tau}^{k} \|w_{\tau E+i}^{k} - \bar{w}_{\tau E}\|^{2}$$
(22)

$$\|\bar{w}_{\tau E+i} - \bar{w}_{\tau E}\|^{2} = \|\sum_{j=0}^{i-1} \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} g_{\tau E+j}^{k} \alpha_{\tau E+j}^{k}\|^{2}$$

$$= \|\eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \left(\sum_{j=0}^{i-1} g_{\tau E+j}^{k} \alpha_{\tau E+j}^{k}\right)\|^{2} = \eta_{\tau}^{2} \|\sum_{k=1}^{N} p^{k} \left(\frac{p_{\tau}^{k}}{p^{k}} \sum_{j=0}^{i-1} g_{\tau E+j}^{k} \alpha_{\tau E+j}^{k}\right)\|^{2}$$

$$\leq \eta_{\tau}^{2} \sum_{k=1}^{N} \frac{(p_{\tau}^{k})^{2}}{p^{k}} \|\sum_{j=0}^{i-1} g_{\tau E+j}^{k} \alpha_{\tau E+j}^{k}\|^{2}$$

$$(23)$$

Here

$$\begin{split} \| \sum_{j=0}^{i-1} g_{\tau E+j}^k \alpha_{\tau E+j}^k \|^2 &= \sum_{j=0}^{i-1} \| g_{\tau E+j}^k \alpha_{\tau E+j}^k \|^2 + 2 \sum_{p < q} \langle g_{\tau E+p}^k \alpha_{\tau E+p}^k, g_{\tau E+q}^k \alpha_{\tau E+q}^k \rangle \\ &\leq \sum_{j=0}^{i-1} \| g_{\tau E+j}^k \alpha_{\tau E+j}^k \|^2 + 2 \sum_{p < q} \| g_{\tau E+p}^k \alpha_{\tau E+p}^k \| \| g_{\tau E+q}^k \alpha_{\tau E+q}^k \| \\ &\leq \sum_{j=0}^{i-1} \| g_{\tau E+j}^k \alpha_{\tau E+j}^k \|^2 + \sum_{p < q} \left(\| g_{\tau E+p}^k \alpha_{\tau E+p}^k \|^2 + \| g_{\tau E+q}^k \alpha_{\tau E+q}^k \|^2 \right) \\ &= i \sum_{j=0}^{i-1} \| g_{\tau E+j}^k \alpha_{\tau E+j}^k \|^2 \end{split}$$

So

$$\mathbb{E}_{\xi} \| \sum_{j=0}^{i-1} g_{\tau E+j}^k \alpha_{\tau E+j}^k \|^2 \le iG^2 \sum_{j=0}^{i-1} \alpha_{\tau E+j}^k \le (E-1)G^2 s_{\tau}^k$$
 (25)

Plug (25) to (23) we have

$$\mathbb{E}_{\xi} \|\bar{w}_{\tau E+i} - \bar{w}_{\tau E}\|^{2} \le (E-1)G^{2} \eta_{\tau}^{2} \sum_{k=1}^{N} \frac{(p_{\tau}^{k})^{2}}{p^{k}} s_{\tau}^{k}$$
(26)

Similarly

$$\mathbb{E}_{\xi} \sum_{k=1}^{N} p_{\tau}^{k} \| w_{\tau E+i}^{k} - \bar{w}_{\tau E} \|^{2} = \mathbb{E}_{\xi} \sum_{k=1}^{N} p_{\tau}^{k} \| \eta_{\tau} \sum_{j=0}^{i-1} g_{\tau E+j}^{k} \alpha_{\tau E+j}^{k} \|^{2} \le (E-1)G^{2} \eta_{\tau}^{2} \sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}$$

$$(27)$$

Plug (26) (27) to (22) we have

$$\mathbb{E}_{\xi}\left[\sum_{k=1}^{N} p_{\tau}^{k} \|\bar{w}_{\tau E+i} - w_{\tau E+i}^{k}\|^{2}\right] \leq (E-1)G^{2}\eta_{\tau}^{2} \left(\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} + (\sum_{k=1}^{N} p_{\tau}^{k} - 2)_{+} \sum_{k=1}^{N} \frac{(p_{\tau}^{k})^{2}}{p^{k}} s_{\tau}^{k}\right)$$
(28)

A.1.3 Bounding $\|\bar{w}_{\tau E+i+1} - w^*\|^2$

$$\|\bar{w}_{\tau E+i+1} - w^*\|^2 = \|\bar{w}_{\tau E+i} - \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k g_{\tau E+i}^k - w^* - \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k \bar{g}_{\tau E+i}^k + \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k \bar{g}_{\tau E+i}^k \|^2$$

$$= \|\bar{w}_{\tau E+i} - w^* - \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k \bar{g}_{\tau E+i}^k \|^2 + \eta_{\tau}^2 \|\sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k (\bar{g}_{\tau E+i}^k - g_{\tau E+i}^k) \|^2$$

$$+ 2\eta_{\tau} \langle \bar{w}_{\tau E+i} - w^* - \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k \bar{g}_{\tau E+i}^k, \sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k (\bar{g}_{\tau E+i}^k - g_{\tau E+i}^k) \rangle$$

$$A_2$$

$$(29)$$

Since $\mathbb{E}_{\xi}[g_{\tau E+i}^k] = \bar{g}_{\tau E+i}^k$, we have $\mathbb{E}_{\xi}[A_2] = 0$. We then bound A_1 .

$$A_{1} = \|\bar{w}_{\tau E+i} - w^{*} - \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k}\|^{2}$$

$$= \|\bar{w}_{\tau E+i} - w^{*}\|^{2} \underbrace{-2\eta_{\tau} \langle \bar{w}_{\tau E+i} - w^{*}, \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k} \rangle}_{B_{1}} + \underbrace{\eta_{\tau}^{2} \|\sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k}\|^{2}}_{B_{2}}$$

$$(30)$$

Since F_k is L-smooth,

$$\|\alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k}\|^{2} \le 2L(F_{k}(w_{\tau E+i}^{k}) - F_{k}^{*})\alpha_{\tau E+i}^{k}$$
(31)

By the convexity of l_2 norm

$$B_{2} = \eta_{\tau}^{2} \| \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k} \|^{2} = \eta_{\tau}^{2} \| \sum_{k=1}^{N} p^{k} (\frac{p_{\tau}^{k}}{p^{k}} \alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k}) \|^{2}$$

$$\leq \eta_{\tau}^{2} \sum_{k=1}^{N} \frac{(p_{\tau}^{k})^{2}}{p^{k}} \| \alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k} \|^{2} \leq 2L\theta \eta_{\tau}^{2} \sum_{k=1}^{N} p_{\tau}^{k} (F_{k}(w_{\tau E+i}^{k}) - F_{k}^{*}) \alpha_{\tau E+i}^{k}$$

$$(32)$$

$$B_{1} = -2\eta_{\tau} \langle \bar{w}_{\tau E+i} - w^{*}, \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k} \rangle = -2\eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \langle \bar{w}_{\tau E+i} - w^{*}, \alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k} \rangle$$

$$= -2\eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \langle \bar{w}_{\tau E+i} - w_{\tau E+i}^{k}, \alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k} \rangle - 2\eta_{\tau} \sum_{k=1}^{N} p_{k} \langle w_{\tau E+i}^{k} - w^{*}, \alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k} \rangle$$

$$(33)$$

Here

$$-2\langle \bar{w}_{\tau E+i} - w_{\tau E+i}^{k}, \alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k} \rangle \leq 2|\langle \bar{w}_{\tau E+i} - w_{\tau E+i}^{k}, \alpha_{\tau E+i}^{k} \bar{g}_{\tau E+i}^{k} \rangle|$$

$$\leq 2\alpha_{\tau E+i}^{k} ||\bar{w}_{\tau E+i} - w_{\tau E+i}^{k}||||\bar{g}_{\tau E+i}^{k}|| \leq \left(\frac{1}{\eta_{\tau}} ||\bar{w}_{\tau E+i} - w_{\tau E+i}^{k}||^{2} + \eta_{\tau} ||\bar{g}_{\tau E+i}^{k}||^{2}\right) \alpha_{\tau E+i}^{k}$$
(34)

Since F_k is μ -strongly convex

$$\langle w_{\tau E+i}^k - w^*, \alpha_{\tau E+i}^k \bar{g}_{\tau E+i}^k \rangle \ge \left((F_k(w_{\tau E+i}^k) - F_k(w^*)) + \frac{\mu}{2} \|w_{\tau E+i}^k - w^*\|^2 \right) \alpha_{\tau E+i}^k$$
 (35)

Plug (34) (35) to (33)

$$B_{1} \leq \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \left(\|\bar{w}_{\tau E+i} - w_{\tau E+i}^{k}\|^{2} + \eta_{\tau}^{2} \|\bar{g}_{\tau E+i}^{k}\|^{2} - 2\eta_{\tau} \left((F_{k}(w_{\tau E+i}^{k}) - F_{k}(w^{*})) + \frac{\mu}{2} \|w_{\tau E+i}^{k} - w^{*}\|^{2} \right) \right)$$
(36)

Plug (32) (36) to (30)

$$A_{1} \leq \|\bar{w}_{\tau E+i} - w^{*}\|^{2} + 2L\theta \eta_{\tau}^{2} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(w_{\tau E+i}^{k}) - F_{k}^{*})$$

$$+ \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \left(\|\bar{w}_{\tau E+i} - w_{\tau E+i}^{k}\|^{2} + \underbrace{\eta_{\tau}^{2} \|\bar{g}_{\tau E+i}^{k}\|^{2}}_{\leq 2\eta_{\tau}^{2} L(F_{k}(w_{\tau E+i}^{k}) - F_{k}^{*})} - 2\eta_{\tau} \left((F_{k}(w_{\tau E+i}^{k}) - F_{k}(w^{*})) + \frac{\mu}{2} \|w_{\tau E+i}^{k} - w^{*}\|^{2} \right) \right)$$

$$\leq \|\bar{w}_{\tau E+i} - w^{*}\|^{2} - \mu \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \|w_{\tau E+i}^{k} - w^{*}\|^{2} + \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \|\bar{w}_{\tau E+i} - w_{\tau E+i}^{k}\|^{2}$$

$$+ 2(1+\theta)L\eta_{\tau}^{2} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(w_{\tau E+i}^{k}) - F_{k}^{*}) - 2\eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(w_{\tau E+i}^{k}) - F_{k}(w^{*}))$$

$$C \qquad (37)$$

$$||w_{\tau E+i}^{k} - w^{*}||^{2} = ||w_{\tau E+i}^{k} - \bar{w}_{\tau E+i} + \bar{w}_{\tau E+i} - w^{*}||^{2}$$

$$= ||w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}||^{2} + ||\bar{w}_{\tau E+i} - w^{*}||^{2} + 2\langle w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}, \bar{w}_{\tau E+i} - w^{*}\rangle$$

$$\geq ||w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}||^{2} + ||\bar{w}_{\tau E+i} - w^{*}||^{2} - 2||w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}|||\bar{w}_{\tau E+i} - w^{*}||$$

$$\geq ||w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}||^{2} + ||\bar{w}_{\tau E+i} - w^{*}||^{2} - (2||w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}||^{2} + \frac{1}{2}||\bar{w}_{\tau E+i} - w^{*}||^{2})$$

$$= \frac{1}{2}||\bar{w}_{\tau E+i} - w^{*}||^{2} - ||w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}||^{2}$$

$$(38)$$

Thus,

$$A_{1} \leq \left(1 - \frac{1}{2}\mu\eta_{\tau}\sum_{k=1}^{N}p_{\tau}^{k}\alpha_{\tau E+i}^{k}\right)\|\bar{w}_{\tau E+i} - w^{*}\|^{2} + \left(1 + \mu\eta_{\tau}\right)\sum_{k=1}^{N}p_{\tau}^{k}\alpha_{\tau E+i}^{k}\|\bar{w}_{\tau E+i} - w_{\tau E+i}^{k}\|^{2} + C$$
(39)

Let $\gamma_{\tau} = 2\eta_{\tau}(1 - (1 + \theta)L\eta_{\tau})$. Assume $\eta_{\tau} \leq \frac{1}{2(1+\theta)L}$, hence $\eta_{\tau} \leq \gamma_{\tau} \leq 2\eta_{\tau}$.

$$C = -2\eta_{\tau}(1 - (1 + \theta)L\eta_{\tau}) \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(w_{\tau E+i}^{k}) - F_{k}^{*}) + 2\eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(w^{*}) - F_{k}^{*})$$

$$= -\gamma_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(w_{\tau E+i}^{k}) - F_{k}^{*} + F_{k}(w^{*}) - F_{k}(w^{*})) + 2\eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(w^{*}) - F_{k}^{*})$$

$$= -\gamma_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(w_{\tau E+i}^{k}) - F_{k}(w^{*})) + (2\eta_{\tau} - \gamma_{\tau}) \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(w^{*}) - F_{k}^{*})$$

$$\leq -\gamma_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(w_{\tau E+i}^{k}) - F_{k}(w^{*})) + 2(1 + \theta)L\eta_{\tau}^{2} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \Gamma_{k}$$

$$(40)$$

Next we bound D

$$\sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(w_{\tau E+i}^{k}) - F_{k}(w^{*})) = \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(w_{\tau E+i}^{k}) - F_{k}(\bar{w}_{\tau E+i})) + \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(\bar{w}_{\tau E+i}) - F_{k}(w^{*}))$$

$$\geq \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (\nabla F_{k}(\bar{w}_{\tau E+i}), w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}) + \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(\bar{w}_{\tau E+i}) - F_{k}(w^{*}))$$

$$\geq -\sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \|\nabla F_{k}(\bar{w}_{\tau E+i})\| \|w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}\| + \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(\bar{w}_{\tau E+i}) - F_{k}(w^{*}))$$

$$\geq -\frac{1}{2} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (\eta_{\tau} \underbrace{\|\nabla F_{k}(\bar{w}_{\tau E+i})\|^{2}}_{\leq 2L(F_{k}(\bar{w}_{\tau E+i}) - F_{k}^{*})} + \frac{1}{\eta_{\tau}} \|w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}\|^{2}) + \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(\bar{w}_{\tau E+i}) - F_{k}(w^{*}))$$

$$\geq -\sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (\eta_{\tau} L(F_{k}(\bar{w}_{\tau E+i}) - F_{k}^{*}) + \frac{1}{2\eta_{\tau}} \|w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}\|^{2}) + \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(\bar{w}_{\tau E+i}) - F_{k}(w^{*}))$$

$$\geq -\sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (\eta_{\tau} L(F_{k}(\bar{w}_{\tau E+i}) - F_{k}^{*}) + \frac{1}{2\eta_{\tau}} \|w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}\|^{2}) + \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(\bar{w}_{\tau E+i}) - F_{k}(w^{*}))$$

$$\geq -\sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (\eta_{\tau} L(F_{k}(\bar{w}_{\tau E+i}) - F_{k}^{*}) + \frac{1}{2\eta_{\tau}} \|w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}\|^{2}) + \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(\bar{w}_{\tau E+i}) - F_{k}(w^{*}))$$

Thus,

$$C \leq \gamma_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (\eta_{\tau} L \underbrace{(F_{k}(\bar{w}_{\tau E+i}) - F_{k}^{*})}_{F_{k}(\bar{w}_{\tau E+i}) - F_{k}(w^{*}) + F_{k}(w^{*}) - F_{k}^{*}} + \frac{1}{2\eta_{\tau}} \|w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}\|^{2})$$

$$- \gamma_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(\bar{w}_{\tau E+i}) - F_{k}(w^{*})) + 2(1+\theta) L \eta_{\tau}^{2} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i} \Gamma_{k}$$

$$= \gamma_{\tau} (\eta_{\tau} L - 1) \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(\bar{w}_{\tau E+i}) - F_{k}(w^{*})) + \underbrace{\frac{\gamma_{\tau}}{2\eta_{\tau}}}_{\leq 1} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \|w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}\|^{2}$$

$$+ 2(1+\theta) L \eta_{\tau}^{2} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i} \Gamma_{k} + \underbrace{\gamma_{\tau}}_{\leq 2\eta_{\tau}} \eta_{\tau} L \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \Gamma_{k}$$

$$\leq \gamma_{\tau} (\eta_{\tau} L - 1) \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(\bar{w}_{\tau E+i}) - F_{k}(w^{*})) + \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \|w_{\tau E+i}^{k} - \bar{w}_{\tau E+i}\|^{2}$$

$$+ 2(2+\theta) L \eta_{\tau}^{2} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \Gamma_{k}$$

$$(42)$$

Plug (42) to (39) we have

$$A_{1} \leq \|\bar{w}_{\tau E+i} - w^{*}\|^{2} - \mu \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \|w_{\tau E+i}^{k} - w^{*}\|^{2} + 2 \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \|\bar{w}_{\tau E+i} - w_{\tau E+i}^{k}\|^{2}$$

$$+ 2(2+\theta) L \eta_{\tau}^{2} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \Gamma_{k} + \gamma_{\tau} (\eta_{\tau} L - 1) \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (F_{k}(\bar{w}_{\tau E+i}) - F_{k}(w^{*}))$$

$$(43)$$

Plug (43) to (29),

$$\|\bar{w}_{\tau E+i+1} - w^*\|^2 \le \left(1 - \frac{1}{2}\mu\eta_{\tau} \sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k\right) \|\bar{w}_{\tau E+i} - w^*\|^2$$

$$+ \eta_{\tau}^2 \|\sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k (\bar{g}_{\tau E+i}^k - g_{\tau E+i}^k)\|^2 + \underbrace{\left(2 + \mu\eta_{\tau}\right)}_{\le 2 + \frac{\mu}{2(1+\theta)L}} \sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k \|\bar{w}_{\tau E+i} - w_{\tau E+i}^k\|^2$$

$$+ 2(2+\theta)L\eta_{\tau}^2 \sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k \Gamma_k + \underbrace{\gamma_{\tau} (1 - \eta_{\tau} L)}_{\le 2\eta_{\tau}} \sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k (F_k(w^*) - F_k(\bar{w}_{\tau E+i}))$$

$$(44)$$

Define

$$C_{\tau E+i} = \left(2 + \frac{\mu}{2(1+\theta)L}\right) \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \|\bar{w}_{\tau E+i} - w_{\tau E+i}^{k}\|^{2} + \|\sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} (\bar{g}_{\tau E+i}^{k} - g_{\tau E+i}^{k})\|^{2} + 2(2+\theta)L \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \Gamma_{k}$$

$$(45)$$

Thus,

$$\|\bar{w}_{\tau E+i+1} - w^*\|^2 \le \left(1 - \frac{1}{2}\mu\eta_{\tau} \sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k\right) \|\bar{w}_{\tau E+i} - w^*\|^2 + \eta_{\tau}^2 B_{\tau E+i}$$

$$+2\eta_{\tau} \sum_{k=1}^{N} p_{\tau}^k \alpha_{\tau E+i}^k (F_k(w^*) - F_k(\bar{w}_{\tau E+i}))$$

$$(46)$$

Apply the lemmas we have

$$\mathbb{E}_{\xi}[C_{\tau E+i}] \leq \sum_{k=1}^{N} (p_{\tau}^{k})^{2} \alpha_{\tau E+i}^{k} \sigma_{k}^{2} + 2(2+\theta) L \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \Gamma_{k}
+ (2 + \frac{\mu}{2(1+\theta)L}) (E-1) G^{2} \left(\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} + (\sum_{k=1}^{N} p_{\tau}^{k} - 2)_{+} \sum_{k=1}^{N} \frac{(p_{\tau}^{k})^{2}}{p^{k}} s_{\tau}^{k} \right)$$
(47)

For convenience we write $\Delta_{\tau E+i} = \|\bar{w}_{\tau E+i} - w^*\|^2$, and $\bar{\Delta}_{\tau E+i} = \mathbb{E}[\Delta_{\tau E+i}]$, where the expectation is taken over all random variables up to $\tau E + i$.

A.1.4 Bounding $\|\bar{w}_{\tau E} - w^*\|$

Summing from τE to $(\tau + 1)E$ we have

$$\sum_{i=1}^{E} \Delta_{\tau E+i} \le \sum_{i=0}^{E-1} \left(1 - \frac{1}{2} \mu \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k}\right) \Delta_{\tau E+i} + \eta_{\tau}^{2} C_{\tau} + 2 \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} (F_{k}(w^{*}) - F_{k}(\bar{w}_{\tau E+l}))$$
(48)

where $C_{\tau} = \sum_{i=0}^{E-1} C_{\tau E+i}$, and $\bar{w}_{\tau E+l} = \operatorname{argmin}_{\bar{w}_{\tau E+i}} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} F_{k}(\bar{w}_{\tau E+i})$.

Reorganize it we can get

$$\Delta_{(\tau+1)E} \le \Delta_{\tau E} - \frac{1}{2}\mu\eta_{\tau} \sum_{i=0}^{E-1} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} \Delta_{\tau E+i} + \eta_{\tau}^{2} C_{\tau} + 2\eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} (F_{k}(w^{*}) - F_{k}(\bar{w}_{\tau E+l}))$$
(49)

We then seek to find a lower bound for $\Delta_{\tau E+i}$.

$$\sqrt{\Delta_{\tau E+i+1}} = \|\bar{w}_{\tau E+i+1} - w^*\| = \|\bar{w}_{\tau E+i+1} - \bar{w}_{\tau E+i} + \bar{w}_{\tau E+i} - w^*\|
\leq \|\bar{w}_{\tau E+i+1} - \bar{w}_{\tau E+i}\| + \sqrt{\Delta_{\tau E+i}}
= \|\eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} g_{\tau E+i}^{k} \| + \sqrt{\Delta_{\tau E+i}}$$
(50)

Define $h_{\tau E+i} = \| \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} g_{\tau E+i}^{k} \|$.

Thus,

$$\sqrt{\Delta_{(\tau+1)E}} \leq \sqrt{\Delta_{(\tau+1)E-1}} + \eta_{\tau} h_{(\tau+1)E-1}$$

$$\leq \dots \leq \sqrt{\Delta_{\tau E+i}} + \sum_{j=i}^{E-1} \eta_{\tau} h_{\tau E+j}$$
(51)

$$\Delta_{(\tau+1)E} \leq \Delta_{\tau E+i} + 2\sqrt{\Delta_{\tau E+i}} (\sum_{j=i}^{E-1} \eta_{\tau} h_{\tau E+j}) + (\sum_{j=i}^{E-1} \eta_{\tau} h_{\tau E+j})^{2}
\leq 2\Delta_{\tau E+i} + 2(\sum_{j=i}^{E-1} \eta_{\tau} h_{\tau E+j})^{2}$$
(52)

$$\Delta_{\tau E+i} \ge \frac{1}{2} \Delta_{(\tau+1)E} - (\sum_{j=i}^{E-1} \eta_{\tau} h_{\tau E+j})^2 \ge \frac{1}{2} \Delta_{(\tau+1)E} - (\sum_{j=0}^{E-1} \eta_{\tau} h_{\tau E+j})^2$$
(53)

Plug (53) to (49) we can get

$$(1 + \frac{1}{4}\mu\eta_{\tau}\sum_{k=1}^{N}p_{\tau}^{k}s_{\tau}^{k})\Delta_{(\tau+1)E} \leq \Delta_{\tau E} + \frac{1}{2}\mu\eta_{\tau}^{3}\sum_{k=1}^{N}p_{\tau}^{k}s_{\tau}^{k}(\sum_{i=0}^{E-1}h_{\tau E+i})^{2} + \eta_{\tau}^{2}C_{\tau} + 2\eta_{\tau}\sum_{k=1}^{N}p_{\tau}^{k}s_{\tau}^{k}(F_{k}(w^{*}) - F_{k}(\bar{w}_{\tau E+l}))$$

$$(54)$$

Define $H_{\tau} = (\sum_{i=0}^{E-1} h_{\tau E+i})^2$. Apply Lemma A.2, Lemma A.3 and Assumption 3.4, we have

$$\mathbb{E}_{\xi}[h_{\tau E+i}^{2}] = \mathbb{E}_{\xi} \| \sum_{k=1}^{N} p_{\tau}^{k} \alpha_{\tau E+i}^{k} g_{\tau E+i}^{k} \|^{2} \\
\leq \sum_{k=1}^{N} \frac{(p_{\tau}^{k})^{2}}{p^{k}} \mathbb{E}_{\xi} \| \alpha_{\tau E+i}^{k} g_{\tau E+i}^{k} \|^{2} \leq \sum_{k=1}^{N} \frac{(p_{\tau}^{k})^{2}}{p^{k}} G^{2} \alpha_{\tau E+i}^{k} \tag{55}$$

$$\mathbb{E}_{\xi}[H_{\tau}] = \mathbb{E}_{\xi}[(\sum_{i=0}^{E-1} h_{\tau E+i})^{2}] \le \mathbb{E}_{\xi}[E\sum_{i=0}^{E-1} h_{\tau E+i}^{2}] \le EG^{2}\sum_{k=1}^{N} \frac{(p_{\tau}^{k})^{2}}{p^{k}} s_{\tau}^{k}$$
(56)

$$\mathbb{E}_{\xi}[C_{\tau}] = \sum_{i=0}^{E-1} \mathbb{E}_{\xi}[C_{\tau E+i}] = \sum_{k=1}^{N} (p_{\tau}^{k})^{2} s_{\tau}^{k} \sigma_{k}^{2} + 2(2+\theta) L \sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} \Gamma_{k}
+ (2 + \frac{\mu}{2(1+\theta)L}) E(E-1) G^{2} \Big(\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} + \theta (\sum_{k=1}^{N} p_{\tau}^{k} - 2)_{+} \sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} \Big)$$
(57)

Write $\bar{\Delta}_{\tau E+i} = \mathbb{E}_{\xi}[\Delta_{\tau E+i}], \bar{C}_{\tau} = \mathbb{E}_{\xi}[C_{\tau}], \ \bar{H}_{\tau} = \mathbb{E}_{\xi}[(\sum_{i=0}^{E-1} h_{\tau E+i})^2],$ then

$$(1 + \frac{1}{4}\mu\eta_{\tau}\sum_{k=1}^{N}p_{\tau}^{k}s_{\tau}^{k})\bar{\Delta}_{(\tau+1)E} \leq \bar{\Delta}_{\tau E} + \frac{1}{2}\mu\eta_{\tau}^{3}\sum_{k=1}^{N}p_{\tau}^{k}s_{\tau}^{k}\bar{H}_{\tau} + \eta_{\tau}^{2}\bar{C}_{\tau} + 2\eta_{\tau}\mathbb{E}_{\xi}\sum_{k=1}^{N}p_{\tau}^{k}s_{\tau}^{k}(F_{k}(w^{*}) - F_{k}(\bar{w}_{\tau E+l}))$$

$$(58)$$

Let $z_{\tau}=0$ indicate the event that for all k, $\mathbb{E}[p_{\tau}^k s_{\tau}^k]=c_{\tau}p^k$ for come constant c_{τ} that does not depend on k, otherwise $z_{\tau}^k=1$. Note that if $z_{\tau}=0$, then $\sum_{k=1}^N p_{\tau}^k s_{\tau}^k (F_k(w^*)-F_k(\bar{w}_{\tau E+l}))=c_{\tau}(F(w^*)-F(\bar{w}_{\tau E+l}))\leq 0$. Otherwise, we have

$$\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} (F_{k}(w^{*}) - F_{k}(\bar{w}_{\tau E+l})) = \sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} (\underbrace{F_{k}(w^{*}) - F_{k}^{*}}_{\Gamma_{k}} + \underbrace{F_{k}^{*} - F_{k}(\bar{w}_{\tau E+l})}_{\leq 0}) \\
\leq \sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} \Gamma_{k} \tag{59}$$

Put it together

$$\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} (F_{k}(w^{*}) - F_{k}(\bar{w}_{\tau E+l})) \leq z_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k} \Gamma_{k}$$

$$(60)$$

Assume $\eta_{\tau} \leq \frac{4}{\mu E \theta} \leq \frac{4}{\mu \sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}}$, divide both sides with $1 + \frac{1}{4} \mu \eta_{\tau} \sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}$ in (58) we can get

$$\bar{\Delta}_{(\tau+1)E} \leq \left(1 - \frac{\frac{1}{4}\mu\eta_{\tau}\sum_{k=1}^{N}p_{\tau}^{k}s_{\tau}^{k}}{1 + \frac{1}{4}\mu\eta_{\tau}\sum_{k=1}^{N}p_{\tau}^{k}s_{\tau}^{k}}\right)\bar{\Delta}_{\tau E} + 2\eta_{\tau}^{2}\bar{H}_{\tau} + \eta_{\tau}^{2}\bar{C}_{\tau}
+ 2\eta_{\tau}z_{\tau}\sum_{k=1}^{N}p_{\tau}^{k}s_{\tau}^{k}\Gamma_{k}
\leq \left(1 - \frac{1}{8}\mu\eta_{\tau}\sum_{k=1}^{N}p_{\tau}^{k}s_{\tau}^{k}\right)\bar{\Delta}_{\tau E} + \eta_{\tau}^{2}B_{\tau} + 2\eta_{\tau}z_{\tau}\sum_{k=1}^{N}p_{\tau}^{k}s_{\tau}^{k}\Gamma_{k}$$
(61)

Note that p_{τ}^k , s_{τ}^k are independent with $\bar{\Delta}_{\tau E}$. Taking expectation over p_{τ}^k and s_{τ}^k we get

$$\mathbb{E}[\bar{\Delta}_{(\tau+1)E}] \le \left(1 - \frac{1}{8}\mu\eta_{\tau}\mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}]\right)\bar{\Delta}_{\tau E} + \eta_{\tau}^{2}\mathbb{E}[B_{\tau}] + 2\eta_{\tau} z_{\tau} \sum_{k=1}^{N} \mathbb{E}[p_{\tau}^{k} s_{\tau}^{k}]\Gamma_{k}$$
(62)

A.2 Proof of Theorem 3.1

When the distributions of s_{τ}^k do not change with time, we have $B_{\tau}=B$. We prove the convergence by induction. Let $\eta_{\tau}=\frac{8}{\mu\mathbb{E}[\sum_{k=1}^N p_{\tau}^k s_{\tau}^k]}\frac{2E}{\tau E+\gamma}$. Initially, $\frac{V_0}{\gamma^2}\geq \mathbb{E}[\bar{\Delta}_0]$. Suppose $\mathbb{E}[\bar{\Delta}_{\tau E}]\leq \frac{M_{\tau}D+V}{\tau E+\gamma}$, then

$$\mathbb{E}[\bar{\Delta}_{(\tau+1)E}] \leq \frac{\tau E + \gamma - 2E}{\tau E + \gamma} \frac{M_{\tau}D + V}{\tau E + \gamma} + \left(\frac{16E}{\mu \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}]}\right)^{2} \frac{B}{(\tau E + \gamma)^{2}} + \frac{\frac{1}{2}z_{\tau}D}{\tau E + \gamma} \\
\leq \frac{\tau E + \gamma - E}{(\tau E + \gamma)^{2}} (M_{\tau}D + V) + \frac{\frac{1}{2}z_{\tau}D}{\tau E + \gamma} + \underbrace{\left(\frac{16E}{\mu \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}]}\right)^{2} \frac{B}{(\tau E + \gamma)^{2}} - \frac{E(M_{\tau}D + V)}{(\tau E + \gamma)^{2}}}_{\leq 0} \\
\leq \frac{M_{\tau}D + V}{(\tau + 1)E + \gamma} + \frac{\frac{1}{2}\frac{\tau E + \gamma + E}{\tau E + \gamma} z_{\tau}D}{(\tau + 1)E + \gamma} \leq \frac{M_{\tau+1}D + V}{(\tau + 1)E + \gamma} \tag{63}$$

Thus $\bar{\Delta}_{(\tau+1)E} \leq \frac{M_{\tau+1}D+V}{(\tau+1)E+\gamma}$.

We can check it satisfies the previous assumptions regarding η_{τ} :

$$\eta_{\tau} \leq \eta_{0} = \frac{16E/(\mu \mathbb{E}\left[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}\right])}{E + \gamma} \\
\leq \frac{16E/(\mu \mathbb{E}\left[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}\right])}{32E(1+\theta)L/(\mu \mathbb{E}\left[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}\right])} = \frac{1}{2(1+\theta)L}$$
(64)

$$\eta_{\tau} \leq \eta_{0} = \frac{16E/(\mu \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}])}{E + \gamma} \leq \frac{16E/(\mu \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}])}{4E^{2}\theta/(\mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}])} = \frac{4}{\mu E \theta}$$
(65)

A.2.1 Extension to Time-Varying Distributions

When the distribution of s_{τ}^{k} vary with time, we can still establish a convergence with slightly different definitions.

$$\text{Redefine } \gamma = \max \left\{ \frac{32E(1+\theta)L}{\mu \min_{\tau} \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^k s_{\tau}^k]}, \frac{4E^2\theta}{\min_{\tau} \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^k s_{\tau}^k]} \right\}, V_{\tau} = \max \left\{ \gamma^2 \mathbb{E} \|w_0^{\mathcal{G}} - w^*\|^2, \left(\frac{16E}{\mu}\right)^2 \sum_{t=0}^{\tau-1} \frac{\mathbb{E}[B_t]}{\left(\mathbb{E}[\sum_{k=1}^{N} p_{\tau}^k s_{\tau}^k]\right)^2} \right\}$$

We now prove by induction that with this definition, we can obtain

$$\mathbb{E}[\bar{\Delta}_{\tau E}] \le \frac{M_{\tau} D}{\tau E + \gamma} + \frac{V_{\tau}}{(\tau E + \gamma)^2} \tag{66}$$

Let $\eta_{\tau} = \frac{8}{\mu \mathbb{E}[\sum_{k=1}^{K} p_{\tau}^{k} s_{\tau}^{k}]} \frac{2E}{(\tau+1)E+\gamma}$. Initially, $\frac{V_{0}}{\gamma^{2}} \geq \mathbb{E}[\bar{\Delta}_{0}]$. Suppose $\mathbb{E}[\bar{\Delta}_{\tau E}] \leq \frac{M_{\tau}D}{\tau E+\gamma} + \frac{V_{\tau}}{(\tau E+\gamma)^{2}}$, then

$$\mathbb{E}[\bar{\Delta}_{(\tau+1)E}] \leq \frac{\tau E + \gamma - E}{(\tau+1)E + \gamma} \left(\frac{M_{\tau}D}{\tau E + \gamma} + \frac{V_{\tau}}{(\tau E + \gamma)^{2}} \right) + \frac{(16E)^{2}\mathbb{E}[\bar{B}_{\tau} + 2\bar{H}_{\tau}]}{(\mu \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}])^{2} ((\tau+1)E + \gamma)^{2}} + \frac{z_{\tau}D}{(\tau+1)E + \gamma} \\
\leq \frac{(\tau E + \gamma - E)M_{\tau}D}{(\tau E + \gamma)^{2} - E^{2}} + \frac{\tau E + \gamma - E}{(\tau E + \gamma)^{2} - E^{2}} \frac{V_{\tau}}{(\tau+1)E + \gamma} + \frac{(16E)^{2}\mathbb{E}[\bar{B}_{\tau} + 2\bar{H}_{\tau}]}{(\mu \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}])^{2} ((\tau+1)E + \gamma)^{2}} + \frac{z_{\tau}D}{(\tau+1)E + \gamma} \\
\leq \frac{M_{\tau}D}{(\tau+1)E + \gamma} + \frac{V_{\tau}}{((\tau+1)E + \gamma)^{2}} + \frac{(16E)^{2}\mathbb{E}[\bar{B}_{\tau} + 2\bar{H}_{\tau}]}{(\mu \mathbb{E}[\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}])^{2} ((\tau+1)E + \gamma)^{2}} + \frac{z_{\tau}D}{(\tau+1)E + \gamma} \\
= \frac{M_{\tau+1}D}{(\tau+1)E + \gamma} + \frac{V_{\tau+1}}{((\tau+1)E + \gamma)^{2}} \tag{67}$$

Thus $\bar{\Delta}_{(\tau+1)E} \leq \frac{M_{\tau+1}D}{(\tau+1)E+\gamma} + \frac{V_{\tau+1}}{((\tau+1)E+\gamma)^2}$

Easy to check previous assumptions regarding η_{τ} are all satisfied.

A.3 Proof of Theorem 3.2

• Departure Case: $\tilde{n} = n - n_l$

$$\begin{split} &\|\tilde{w}^* - w^*\| \le \frac{2}{\mu} \|\nabla F(\tilde{w}^*)\| = \frac{2}{\mu} \left\| \nabla F(\tilde{w}^*) - \underbrace{\nabla \tilde{F}(\tilde{w}^*)}_{=0} \right\| \\ &= \frac{2}{\mu} \left\| \sum_{k \ne l} (p^k - \tilde{p}^k) \nabla F_k(\tilde{w}^*) + p^l \nabla F_l(\tilde{w}^*) \right\| \\ &= \frac{2}{\mu} \left\| \sum_{k \ne l} \left(\frac{n_k}{n} - \frac{n_k}{n - n_l} \right) \nabla F_k(\tilde{w}^*) + p^l \nabla F_l(\tilde{w}^*) \right\| \\ &= \frac{2}{\mu} \left\| - \sum_{k \ne l} \left(\frac{n_l n_k}{n(n - n_l)} \right) \nabla F_k(\tilde{w}^*) + p^l \nabla F_l(\tilde{w}^*) \right\| \\ &= \frac{2}{\mu} \left\| - p^l \underbrace{\sum_{k \ne l} \tilde{p}^k \nabla F_k(\tilde{w}^*) + p^l \nabla F_l(\tilde{w}^*)}_{=\nabla \tilde{F}(\tilde{w}^*) = 0} \right\| \\ &= \frac{2p^l}{\mu} \left\| \nabla F_l(\tilde{w}^*) \right\| \le \frac{2p^l}{\mu} \sqrt{2L \left(F_l(\tilde{w}^*) - F_l^* \right)} = \frac{2\sqrt{2L}}{\mu} p^l \sqrt{\tilde{\Gamma}_l} \end{split}$$

• Arrival Case: $\tilde{n} = n + n_l$

$$\|\tilde{w}^* - w^*\| = \|w^* - \tilde{w}^*\| \le \frac{2}{\mu} \|\nabla \tilde{F}(w^*)\| = \frac{2}{\mu} \left\| \nabla \tilde{F}(w^*) - \underbrace{\nabla F(w^*)}_{=0} \right\|$$

$$= \dots = \frac{2}{\mu} \left\| -\tilde{p}^l \underbrace{\sum_{k \neq l} p^k \nabla F_k(w^*) + \tilde{p}^l \nabla F_l(w^*)}_{=\nabla F(w^*) = 0} \right\|$$

$$= \frac{2\tilde{p}^l}{\mu} \|\nabla F_l(w^*)\| = \frac{2\sqrt{2L}}{\mu} \tilde{p}^l \sqrt{\Gamma_l}$$

A.4 Proof of Corollary 4.0.1

A.4.1 Scheme A

In Scheme A, we only consider devices whose $s_{\tau}^k = E$. Let q_{τ}^k be an indicator denoting if client k is complete in round τ . Thus, $K_{\tau} = \sum_{k=1}^{N} q_{\tau}^k$.

Homogeneous participation. Obviously q_{τ}^k 's are homogeneous when s_{τ}^k 's are homogeneous. Thus, $\mathbb{E}[q_{\tau}^k] = q_{\tau}$, where $q_{\tau} = \mathbb{P}(s_{\tau} = E)$. We then have $\mathbb{P}(K_{\tau} = 0) = (1 - q_{\tau})^N$. When choosing $p_{\tau}^k = \frac{Np^k}{K_{\tau}}q_{\tau}^k$, $\theta = N$. Note that by the definition of q_{τ}^k , we have $q_{\tau}^k s_{\tau}^k = Eq_{\tau}^k$, so $\mathbb{E}[p_{\tau}^k s_{\tau}^k] = E\mathbb{E}[p_{\tau}^k]$. Similarly, we can replace all s_{τ}^k terms with E. Next we calculate $\mathbb{E}[p_{\tau}^k]$:

$$\mathbb{E}_{q}[p_{\tau}^{k}|K_{\tau} \neq 0] = Np^{k}\mathbb{E}_{q}\left[\frac{q_{\tau}^{k}}{\sum_{i=1}^{N}q_{\tau}^{i}}\Big|K_{\tau} \neq 0\right] \\
=Np^{k}\sum_{i=0}^{N-1}\frac{1}{1+i}\binom{N-1}{i}\frac{(q_{\tau})^{i}(1-q_{\tau})^{N-1-i}q_{\tau}}{1-(1-q_{\tau})^{N}} \\
=Np^{k}\sum_{i=1}^{N}\frac{1}{i}\binom{N-1}{i-1}\frac{(q_{\tau})^{i}(1-q_{\tau})^{N-i}}{1-(1-q_{\tau})^{N}} = Np^{k}\sum_{i=1}^{N}\frac{1}{i}\frac{(N-1)!}{(i-1)!(N-i)!}\frac{(q_{\tau})^{i}(1-q_{\tau})^{N-i}}{1-(1-q_{\tau})^{N}} \\
=p^{k}\sum_{i=1}^{N}\binom{N}{i}\frac{(q_{\tau})^{i}(1-q_{\tau})^{N-i}}{1-(1-q_{\tau})^{N}} = p^{k}\frac{1-\binom{N}{0}(q_{\tau})^{0}(1-q_{\tau})^{N}}{1-(1-q_{\tau})^{N}} = p^{k}$$
(68)

Similarly,

$$\mathbb{E}_{q}[(p_{\tau}^{k})^{2}|K_{\tau} \neq 0] = (Np^{k})^{2}\mathbb{E}_{q}\left[\frac{q_{\tau}^{k}}{(\sum_{i=1}^{N}q_{\tau}^{i})^{2}}\Big|K_{\tau} \neq 0\right] \\
= (Np^{k})^{2}\sum_{i=0}^{N-1}\frac{1}{(1+i)^{2}}\binom{N-1}{i}\frac{(q_{\tau})^{i}(1-q_{\tau})^{N-1-i}q_{\tau}}{1-(1-q_{\tau})^{N}} \\
= (Np^{k})^{2}\sum_{i=1}^{N}\frac{1}{i^{2}}\binom{N-1}{i-1}\frac{(q_{\tau})^{i}(1-q_{\tau})^{N-i}}{1-(1-q_{\tau})^{N}} = N(p^{k})^{2}\sum_{i=1}^{N}\frac{1}{i}\binom{N}{i}\frac{(q_{\tau})^{i}(1-q_{\tau})^{N-i}}{1-(1-q_{\tau})^{N}} \\
= N(p^{k})^{2}\mathbb{E}\left[\frac{1}{K_{\tau}}|K_{\tau} \neq 0\right] \tag{69}$$

It is possible that $\sum_{k=1}^{N} p_{\tau}^{k} > 2$, so we need to calculate $\mathbb{E}[p_{\tau}^{k} p_{\tau}^{l} | K_{\tau} \neq 0]$

$$\mathbb{E}_{q}[p_{\tau}^{k}p_{\tau}^{l}|K_{\tau} \neq 0] = N^{2}p^{k}p^{l}\mathbb{E}\left[\frac{q_{\tau}^{k}q_{\tau}^{l}}{(\sum_{i=1}^{N}q_{\tau}^{i})^{2}}|K_{\tau} \neq 0\right]$$

$$=N^{2}p^{k}p^{l}\sum_{i=0}^{N-2}\frac{1}{(2+i)^{2}}\binom{N-2}{i}\frac{(q_{\tau})^{i}(1-q_{\tau})^{N-2-i}(q_{\tau})^{2}}{1-(1-q_{\tau})^{N}}$$

$$=\frac{N}{N-1}p^{k}p^{l}\sum_{i=2}^{N}\frac{i-1}{i}\binom{N}{i}\frac{(q_{\tau})^{i}(1-q_{\tau})^{N-i}}{1-(1-q_{\tau})^{N}} = \frac{N}{N-1}p^{k}p^{l}\mathbb{E}\left[1-\frac{1}{K_{\tau}}|K_{\tau} \neq 0\right]$$
(70)

For all k and τ , $\mathbb{E}[p_{\tau}^k s_{\tau}^k | K_{\tau} \neq 0] = Ep^k$, thus $z_{\tau} = 0, M_{\tau} = 0$ for all k, τ .

Therefore, $\mathbb{E}[B] = O(N^2 \mathbb{E}[\frac{1}{K_{\tau}} | K_{\tau} \neq 0] + \sum_{k=1}^{N} (p^k \sigma_k)^2 + \Gamma), \gamma = O(N)$, hence $V = O(N^2 \mathbb{E}[\frac{1}{K_{\tau}} | K_{\tau} \neq 0] + \sum_{k=1}^{N} (p^k \sigma_k)^2 + \Gamma)$. Plug them into Theorem 3.1, we can get an asymptotic rate of $O\left(\frac{\mathbb{E}[\frac{N^2}{K_{\tau}}] + \bar{\sigma}_N^2 + \Gamma}{\tau}\right)$.

Heterogeneous Participation. When s_{τ}^{k} 's (i.e., q_{τ}^{k} 's) are heterogeneous, generally $\mathbb{E}[p_{\tau}^{k}] \neq p^{k}$, furthermore, we may have $z_{\tau} = 1$ for all τ . To see this, consider an example where a device k_{0} has $q_{\tau}^{k_{0}} = 1$, i.e. $\mathbb{P}(s_{\tau}^{k} = E) = 1$, whereas all the rest devices have $\mathbb{E}[q_{\tau}^{k}] = q_{\tau}$, then we can show that

$$\mathbb{E}_{q}[p_{\tau}^{k_{0}}|K_{\tau} \neq 0] = \mathbb{E}_{q}[p_{\tau}^{k_{0}}] = Np^{k_{0}}\mathbb{E}_{q}\left[\frac{q_{\tau}^{k_{0}}}{(\sum_{i=1}^{N} q_{\tau}^{i})^{2}}\right] = p^{k_{0}}\frac{1 - (1 - q_{\tau})^{N}}{q_{\tau}}$$
(71)

and for $k \neq k_0$

$$\mathbb{E}_{q}[p_{\tau}^{k}|K_{\tau} \neq 0] = \mathbb{E}_{q}[p_{\tau}^{k}] = Np^{k}\mathbb{E}_{q}\left[\frac{q_{\tau}^{k}}{(\sum_{i=1}^{N}q_{\tau}^{i})^{2}}\right] \\
= \frac{p^{k}}{(N-1)q_{\tau}} \sum_{i=2}^{N} (i-1)\binom{N}{i} (q_{\tau})^{i} (1-q_{\tau})^{N-i} \\
= \frac{p^{k}}{(N-1)q_{\tau}} \left(Nq_{\tau} - Nq_{\tau}(1-q_{\tau})^{N-1} - (1-(1-q_{\tau})^{N} - Nq_{\tau}(1-q_{\tau})^{N-1})\right) \\
= p^{k} \frac{Nq_{\tau} + (1-q_{\tau})^{N} - 1}{(N-1)q_{\tau}} \tag{72}$$

Thus, different k will have different ratio of $\mathbb{E}[p_{\tau}^k s_{\tau}^k/p^k] = E\mathbb{E}[p_{\tau}^k/p^k]$, which indicates $z_{\tau} = 1$. Since this is true for all τ , we have $M_{\tau} = \tau$. Thus according to Theorem [3.1] the learning will not converge to the global optimal, and the remainder loss is bounded by D/E.

A.4.2 Scheme B

In Scheme B, $p_{\tau}^k = p^k$ is a fixed number, so we only need to take expectation over s_{τ}^k , and $c_p = 1$. Since $\sum_{k=1}^N p^k = 1 < 2$, we can bound $\mathbb{E}[(\sum_{k=1}^N p_{\tau}^k - 2)_+ (\sum_{k=1}^N p_{\tau}^k s_{\tau}^k)] < 0$.

Homogeneous Participation. When s_{τ}^k 's are homogeneous, i.e. $s_{\tau}^k \stackrel{iid}{\sim} s_{\tau}$, then $\mathbb{E}[p_{\tau}^k s_{\tau}^k]/p^k = \mathbb{E}[s_{\tau}]$. This is the same for all k, thus $z_{\tau} = 0$ and $M_{\tau} = 0$. Moreover, we have $\mathbb{E}[B] = O(\mathbb{E}[s_{\tau}](\bar{\sigma}_N^2 + \Gamma))$, $\gamma = O(1/\mathbb{E}[s_{\tau}])$, $V = O\left(\left(\bar{\sigma}_N^2 + \Gamma\right) \frac{1}{\mathbb{E}[s_{\tau}]}\right)$, which yields an asymptotic convergence rate of $O\left(\frac{\bar{\sigma}_N^2 + \Gamma}{\tau \mathbb{E}[s_{\tau}]}\right)$.

Heterogeneous Participation. When s_{τ}^{k} 's are heterogeneous, $\mathbb{E}[p_{\tau}^{k}s_{\tau}^{k}]/p^{k} = \mathbb{E}[s_{\tau}^{k}]$ varies with k. Thus, $z_{\tau} = 1$ and $M_{\tau} = \tau$. Therefore, the algorithm will not converge to the global optimum according to Theorem [3.1]

A.4.3 Scheme C

In Scheme C, $p_{\tau}^k = \frac{Ep^k}{s_{\tau}^k}$, so $\theta = E$. It is possible that $\sum_{k=1}^K p_{\tau}^k > 2$, so we need to calculate $\mathbb{E}\left[(\sum_{k=1}^N p_{\tau}^k)(\sum_{k=1}^N p_{\tau}^k s_{\tau}^k)\right]$.

Homogeneous Participation. When s_{τ}^{k} 's are homogeneous, $\mathbb{E}[p_{\tau}^{k}s_{\tau}^{k}]/p^{k} = E$ for all k. Thus, $z_{\tau} = 0$, $M_{\tau} = 0$. Moreover, we have

$$\mathbb{E}\left[\sum_{k=1}^{N} p_{\tau}^{k}\right] = E\mathbb{E}\left[\frac{1}{s_{\tau}}\right] \tag{73}$$

$$\mathbb{E}[\sum_{k=1}^{N} (p_{\tau}^{k})^{2}] = \left(E\mathbb{E}\left[\frac{1}{s_{\tau}}\right]\right)^{2} \sum_{k=1}^{N} (p^{k})^{2}$$
(74)

$$\mathbb{E}\left[\sum_{k=1}^{N} (p_{\tau}^{k})^{2} s_{\tau}^{k}\right] = E^{2} \mathbb{E}\left[\frac{1}{s_{\tau}}\right] \sum_{k=1}^{N} (p^{k})^{2}$$
(75)

$$\mathbb{E}\left[\left(\sum_{k=1}^{N} p_{\tau}^{k}\right)\left(\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}\right)\right] = E^{2} \mathbb{E}\left[\frac{1}{s_{\tau}}\right]$$
(76)

Therefore, we have $\mathbb{E}[B] = O\left(\mathbb{E}\left[\frac{1}{s_{\tau}}\right](\bar{\sigma}_N + \Gamma)\right) = V$, which yields a convergence rate of $O\left(\frac{\bar{\sigma}_N^2 + \Gamma}{\tau(\mathbb{E}[1/s_{\tau}])^{-1}}\right)$.

Heterogeneous Participation. Even when s_{τ}^k 's are heterogeneous, we still have $\mathbb{E}[p_{\tau}^k s_{\tau}^k]/p^k = E$ for active all

k. Thus, $z_{\tau}=0$ only if $I_{\tau}=1$. Thus, $M_{\tau}=\sum_{t=0}^{\tau-1}I_{t}$. Moreover,

$$\mathbb{E}\left[\sum_{k=1}^{N} p_{\tau}^{k}\right] = E \sum_{k=1}^{N} p^{k} \mathbb{E}\left[\frac{1}{s_{\tau}^{k}}\right] \tag{77}$$

$$\mathbb{E}[\sum_{k=1}^{N} (p_{\tau}^{k})^{2}] = E^{2} \sum_{k=1}^{N} \left(p^{k} \mathbb{E}\left[\frac{1}{s_{\tau}^{k}}\right] \right)^{2}$$
(78)

$$\mathbb{E}\left[\sum_{k=1}^{N} (p_{\tau}^{k})^{2} s_{\tau}^{k}\right] = E^{2} \sum_{k=1}^{N} (p^{k})^{2} \mathbb{E}\left[\frac{1}{s_{\tau}^{k}}\right]$$
(79)

$$\mathbb{E}\Big[\left(\sum_{k=1}^{N} p_{\tau}^{k}\right)\left(\sum_{k=1}^{N} p_{\tau}^{k} s_{\tau}^{k}\right)\Big] = E^{2} \sum_{k=1}^{N} p^{k} \mathbb{E}\Big[\frac{1}{s_{\tau}^{k}}\Big]$$
(80)

Thus,
$$\mathbb{E}[B] = O\left(\sum_{k=1}^{N} (p^k \sigma_k)^2 \mathbb{E}\left[\frac{1}{s_{\tau}^k}\right] + \Gamma\right) = V$$
, and the convergence rate is $O\left(\frac{\sum\limits_{t=0}^{\tau-1} I_t D + \sum\limits_{k}^{N} (p^k \sigma_k)^2 \mathbb{E}\left[\frac{1}{s_{\tau}^k}\right] + \Gamma}{\tau}\right)$

A.5 Proof of Corollary 4.0.2

We first introduce the following lemma:

Lemma A.4. Suppose device l arrives, then for any w, we have

$$F_l(w) = \frac{1}{\tilde{p}^l} \left(\tilde{F}(w) - \frac{n}{\tilde{n}} F(w) \right)$$
(81)

Proof. We expand the right hand side expression and show it equals $F_l(w)$:

$$\frac{1}{\tilde{p}^{l}} \left(\tilde{F}(w) - \frac{n}{\tilde{n}} F(w) \right) = \frac{1}{\tilde{p}^{l}} \left(\sum_{k=1}^{N} \tilde{p}^{k} F_{k}(w) + \tilde{p}^{l} F_{l}(w) - \sum_{k=1}^{N} \frac{n}{\tilde{n}} p^{k} F_{k}(w) \right)
= \frac{1}{\tilde{p}^{l}} \left(\sum_{k=1}^{N} \tilde{p}^{k} F_{k}(w) + \tilde{p}^{l} F_{l}(w) - \sum_{k=1}^{N} \tilde{p}^{k} F_{k}(w) \right) = F_{l}(w)$$
(82)

Next we investigate the effect of applying additional update from l. Suppose the current global weight is $w_{\tau E}^{\mathcal{G}} = w$, and assume we perform full batch gradient for the additional update. After the update, it becomes

$$w' = w - \eta_{\tau} \delta^{l} \nabla F_{l}(w) \tag{83}$$

We are interested in the distance between w' and the new global optimum \tilde{w}^* :

$$\|w' - \tilde{w}^*\|^2 = \|w - \eta_\tau \delta^l \nabla F_l(w) - \tilde{w}^*\|^2$$

$$= \|w - \tilde{w}^*\|^2 \underbrace{-2\eta_\tau \delta^l \langle w - \tilde{w}^*, \nabla F_l(w) \rangle + (\eta_\tau \delta^l)^2 \|\nabla F_l(w)\|^2}_{A(w, \delta^l)}$$
(84)

Obviously, the additional update helps fast-reboot if $A(w, \delta^l) < 0$.

Applying Lemma A.4 we can get

$$A(w, \delta^l) = -2\frac{\eta_\tau \delta^l}{\tilde{p}^l} \langle w - \tilde{w}^*, \nabla \tilde{F}(w) - \frac{n}{\tilde{n}} \nabla F(w) \rangle + \left(\eta_\tau \delta^l\right)^2 \|\nabla F_l(w)\|^2$$
(85)

Write $b = w - w^*$, and use the mean value theorem we have

$$-\langle w - \tilde{w}^*, \nabla F_l(w) \rangle = -\langle b + w^* - \tilde{w}^*, \nabla F_l(w^*) + \nabla^2 F_l(\xi)b \rangle$$

$$= -\langle w^* - \tilde{w}^*, \nabla F_l(w^*) \rangle - \langle w^* - \tilde{w}^*, \nabla^2 F_l(\xi)b \rangle - \langle b, \nabla F_l(w) \rangle$$

$$\leq -\langle w^* - \tilde{w}^*, \nabla F_l(w^*) \rangle + \|w^* - \tilde{w}^*\| \|\nabla^2 F_l(\xi)\|_2 \|b\| + \|\nabla F_l(w)\| \|b\|$$

$$\leq -\langle w^* - \tilde{w}^*, \nabla F_l(w^*) \rangle + (\|w^* - \tilde{w}^*\| + 1)W \|b\|$$

$$\leq -\frac{1}{\tilde{p}^l} \langle w^* - \tilde{w}^*, \nabla \tilde{F}(w^*) - \frac{n}{\tilde{n}} \underbrace{\nabla F(w^*)}_{=0} \rangle + \left(\frac{2\sqrt{2L}}{\mu} \tilde{p}^l \sqrt{\Gamma_l} + 1\right) W \|b\|$$

$$\leq -\frac{1}{\tilde{p}^l} \left(\tilde{F}(w^*) - \tilde{F}(\tilde{w}^*)\right) + \left(\frac{2\sqrt{2L}}{\mu} \tilde{p}^l \sqrt{\Gamma_l} + 1\right) W \|b\|$$

$$(86)$$

Therefore,

$$A(w, \delta^l) \le 2\frac{\eta_\tau \delta^l}{\tilde{p}^l} \left(\left(\frac{2\sqrt{2L}}{\mu} \tilde{p}^l \sqrt{\Gamma_l} + 1 \right) W \|b\| - \left(\tilde{F}(w^*) - \tilde{F}(\tilde{w}^*) \right) \right) + (\eta_\tau \delta^l)^2 W^2$$

$$\tag{87}$$

For $\delta^l > 0$, the right hand side can be negative if and only if $||b|| < \frac{\tilde{F}(w^*) - \tilde{F}(\tilde{w}^*)}{\left(\frac{2\sqrt{2L}}{\mu}\tilde{p}^l\sqrt{\Gamma_l} + 1\right)\tilde{p}^lW}$.

A.6 Proof of Corollary 4.0.3

The loss bound without objective shift is $f_0(\tau) = \frac{(\tau - \tau_0)D + V}{\tau E + \gamma}$, and the bound with shift is $f_1(\tau) = \frac{\frac{V}{\tau_0 E + \gamma} + \Gamma_l}{(\tau - \tau_0)E + \gamma}$. Note that $f_0(\tau)$ is a monotonic function. When it is increasing, we just need $f_0(\tau_0) = f_1(\tau)$, which yields

$$\tau - \tau_0 = \frac{1 - \gamma}{E} + \frac{\Gamma_l(\tau_0 E + \gamma)}{EV} = O\left(\frac{\Gamma_l \tau_0}{V}\right)$$
(88)

Now we consider monotonically decreasing $f_0(\tau)$, which is more commonly observed in experiments. Let $C_1 = DE$, $C_2 = \gamma D + VE - E\Gamma_l$, $C_3 = V(\gamma - 1)$, the only possible root for the quadratic equation $f_0(\tau) = f_1(\tau)$ is

$$\tau - \tau_0 = \frac{EV}{\tau_0 E + \gamma} - C_2 + \sqrt{4C_1 \Gamma_l(\tau_0 E + \gamma) + \left(\frac{EV}{\tau_0 E + \gamma}\right)^2 - \frac{2C_2 EV}{\tau_0 E + \gamma} + (C_2^2 - 4C_1 C_3)}$$

$$= O(\sqrt{\tau_0 \Gamma_l})$$
(89)