# Resource-Efficient and Delay-Aware Federated Learning Design under Edge Heterogeneity

David Nickel*, Frank Po-Chen Lin*, Seyyedali Hosseinalipour*, Nicolo Michelusi†, and Christopher G. Brinton*

*Electrical and Computer Engineering, Purdue University, West Lafayette, IN, USA

†Electrical, Computer, and Energy Engineering, Arizona State University, Tempe, AZ, USA

*{dnickel, lin1183, hosseina, cgb}@purdue.edu, †nicolo.michelusi@asu.edu

*Abstract*—**Federated learning (FL) has emerged as a popular technique for distributing machine learning across wireless edge devices. We examine FL under two salient properties of contemporary networks: device-server communication delays and device computation heterogeneity. Our proposed `StoFedDelAv` algorithm incorporates a local-global model combiner into the FL synchronization step. We theoretically characterize the convergence behavior of `StoFedDelAv` and obtain the optimal *combiner weights*, which consider the global model delay and expected local gradient error at each device. We then formulate a network-aware optimization problem which tunes the minibatch sizes of the devices to jointly minimize energy consumption and machine learning training loss, and solve the non-convex problem through a series of convex approximations. Our simulations reveal that `StoFedDelAv` outperforms the current art in FL, evidenced by the obtained improvements in optimization objective.**

## I. INTRODUCTION

Recent advancements in smart devices (e.g. cell phones, smart cars) have resulted in a paradigm shift for machine learning (ML) [1], aiming to migrate intelligence management from cloud datacenters to the network edge [2]. Federated learning (FL) has been promoted as one of the main frameworks for distributing ML over wireless networks [3], where model training is conducted without data exchange across devices.

Conventional FL operates in two iterative steps [4]: (i) local training, where edge devices update their local models using their own datasets; and (ii) global aggregation, where a cloud server computes the global model based on local models received from the edge devices, and synchronizes them [5]. Implementations of this process over the wireless edge are complicated by heterogeneity in communication and computation capabilities found across devices [6]. In this work, we augment FL to provide resilience to these factors.

### A. Related Works

Several works in FL have focused on techniques for improving device-to-server communication efficiency in the global aggregation step. Some have focused on reducing the number of uplink/downlink communication rounds by performing multiple iterations of local model updates between consecutive global aggregations [7], [8]. Works [9], [10] showed that device-server communication requirements in FL can be further reduced through direct device-to-device model synchronization.

Building upon this, there has been a recent trend towards control methodologies for optimizing device participation in FL. The authors of [11] proposed a joint optimization formulation considering learning, resource allocation, and device selection to minimize convergence time. In [12], the authors minimized the total energy consumption of the system under device heterogeneity constraints. In [13], the authors developed over-the-air FL for maximizing global model aggregation speed under proper device selection and beamforming design.

Such works have largely neglected the effect of *communication delay* on the performance of model training in FL. In [14], we took a step towards addressing this by establishing a delay-aware FL framework. Specifically, we introduced a mechanism for devices to combine local and global models during the synchronization step to account for communication delay. Nevertheless, [14] considers a scenario in which the edge devices train their models in the local straining step using full-batch gradient descent (GD). This can introduce large inefficiencies with respect to the energy consumed versus model convergence obtained in FL, especially when training models over heterogeneous wireless devices. In practice, an edge device can potentially store more data than it can process in a timely manner. An energy-efficient solution to this is using minibatch stochastic gradient descent (SGD), which on the other hand has the downside of introducing estimation noise [9]. In this paper, we address these challenges by coupling the selection of device minibatch sizes with the weighting of local and global model combiners based on heterogeneity conditions.

### B. Outline and Summary of Contributions

- We develop a delay-aware FL framework, `StoFedDelAv`, which incorporates a local-global model combiner to jointly optimize model training performance and network resource consumption in the presence of device-server communication delays and device computation heterogeneity.

- We theoretically characterize the convergence behavior of `StoFedDelAv` and optimize the local-global model combiner weight in the presence of communication delay. We further formulate a network-aware learning optimization problem which aims to tune the SGD minibatch sizes across the devices according to resource constraints. We demonstrate that the problem is a non-convex signomial program, and solve it using a series of convex approximations.

- Our experiments show that `StoFedDelAv` outperforms the current art in FL in terms of model convergence speed and network resource utilization when the minibatch size and local-global model combiner are carefully adjusted.

## II. SYSTEM MODEL AND ALGORITHM

### A. Network and Machine Learning Model

We consider a set $\mathcal{I} = \{1, \cdots, I\}$ of $I$ edge devices connected to a cloud server, which acts as a model aggregator (see Fig. 1). Each edge device $i$ is associated with a dataset $\mathcal{D}_i$, where each datapoint $(\mathbf{x}, y) \in \mathcal{D}_i$ comprises an $m$-dimensional feature vector, $\mathbf{x} \in \mathbb{R}^m$, and a label, $y \in \mathbb{R}$.

We let $f_i(\mathbf{x}, y; \mathbf{w})$ be the loss of the machine learning model associated with datapoint $(\mathbf{x}, y)$ and model parameter vector $\mathbf{w} \in \mathbb{R}^n$. The local loss function of device $i$ is given by

$$F_i(\mathbf{w}) = \sum_{(\mathbf{x},y) \in \mathcal{D}_i} f_i(\mathbf{x}, y; \mathbf{w})/N_i. \tag{1}$$

The global loss is subsequently defined as

$$F(\mathbf{w}) = \sum_{i \in \mathcal{I}} \rho_i F_i(\mathbf{w}), \tag{2}$$

where $\rho_i = N_i / \sum_{j \in \mathcal{I}} N_j$ is the weight associated with device $i$. $N_i = |\mathcal{D}_i|$ is the size of the local dataset. The goal of the ML training is to find the optimal parameter given by

$$\mathbf{w}^\star = \arg\min_{\mathbf{w}} F(\mathbf{w}). \tag{3}$$

To aid in convergence analysis of model training across the network, the following assumptions are made:

**Assumption 1.** *The loss functions are assumed to be L-Lipschitz and $\beta$-Smooth, i.e., $\|F_i(\boldsymbol{w}_1) - F_i(\boldsymbol{w}_2)\| \leq L\|\boldsymbol{w}_1 - \boldsymbol{w}_2\|, \forall i$, and $\|\nabla F_i(\boldsymbol{w}_1) - \nabla F_i(\boldsymbol{w}_2)\| \leq \beta\|\boldsymbol{w}_1 - \boldsymbol{w}_2\|, \forall i$.*

**Assumption 2.** *The local and global gradients are assumed to have a bounded dissimilarity, i.e. $\|\nabla F_i(\boldsymbol{w}) - \nabla F(\boldsymbol{w})\| \leq \delta_i, \forall \boldsymbol{w}, \forall i$, where $0 \leq \delta_i \leq 2L$. We let $\delta = \sum_i \rho_i \delta_i$.*
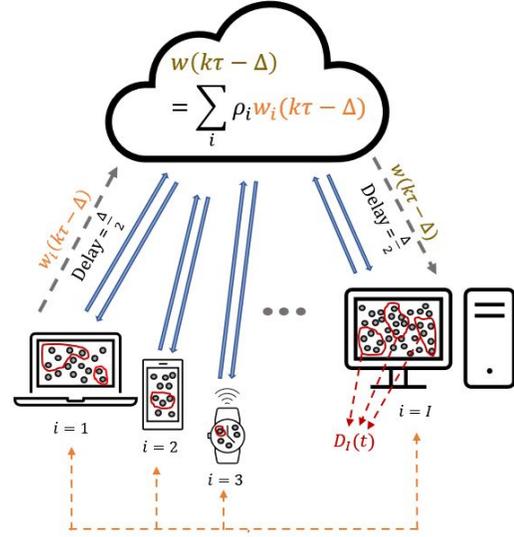
Note that a higher value of $\delta$ implies a larger statistical diversity across the local datasets of the edge devices.

### B. StoFedDelAv *Algorithm*

We propose the StoFedDelAv algorithm (see Alg. 1), considering the effect of the communication delay between the edge devices and the cloud server. We divide the full training cycle into discrete time-instances $t \in \{1, 2, ..., T\}$, where the training consists of $K = \frac{T}{\tau}$ rounds of aggregation. $\tau$ denotes the number of SGD steps taken by each device for each round of global aggregation indexed by $k \in \{0, 1, ..., K-1\}$, where each aggregation period spans the interval $\mathcal{T}_k = \{k\tau - \Delta + 1, ..., (k+1)\tau - \Delta\}$. The communication delay, i.e., the duration between when edge devices send their models to the server and the reception of the resulting global model is denoted by $\Delta$, where $\tau \geq \Delta \geq 0$. Without loss of generality, we assume the uplink and downlink communication delay to be symmetric, i.e., $\Delta/2$, for both upstream and downstream communications.

Let $\mathbf{w}_i(t)$ denote the local model trained at each device $i$ and $\mathbf{w}(t) = \sum_i \rho_i \mathbf{w}_i(t)$ be the global model at each time instance $t$. The model training starts with the cloud server initializing all the local models such that $\mathbf{w}_i(-\Delta) = \mathbf{w}(-\Delta), \forall i$.

Between two consecutive global aggregations, each device sends its local model $\mathbf{w}_i(t)$ to the server at $t \in \{k\tau - \Delta, \forall k \geq 0\}$, after waiting for the communication delay between edge and server, i.e., $\Delta/2$, and the global model $\mathbf{w}(t)$ is computed at the server at $t \in \{k\tau - \Delta/2, \forall k \geq 0\}$. Finally, the devices receive the global model at $k\tau$ to perform local model synchronization.



Fig. 1: System architecture and illustration of our proposed methodology for delay-aware federated learning.

**Distributed SGD:** At time $t$, the edge devices sample their datasets randomly and without replacement, obtaining minibatch $\mathcal{D}_i(t) \subseteq \mathcal{D}_i$, where $|\mathcal{D}_i(t)|$ is the number of datapoints selected and is the same for each $t \in \mathcal{T}_k$. Let $n_i(k) \triangleq |\mathcal{D}_i(t)|, \forall t \in \mathcal{T}_k$ be the minibatch size of device $i$ for the $k$-th aggregation period, each local device take an SGD step on their local model using unbiased gradient estimator as:

$$g_i(\mathbf{w}_i(t); \mathcal{D}_i(t)) = |\mathcal{D}_i(t)|^{-1} \sum_{(\mathbf{x},y) \in \mathcal{D}_i(t)} \nabla f_i(\mathbf{x}, y; \mathbf{w}_i(t)), \tag{4}$$

where

$$g_i(\mathbf{w}_i(t); \mathcal{D}_i(t)) = \nabla F_i(\mathbf{w}_i(t)) + \nu_i(t) \tag{5}$$

with $\nu_i(t)$ being a zero-mean noise.

At each time $t \in \mathcal{T}_k \setminus \{k\tau\}$, each edge device updates the local model Using the gradient estimate as:

$$\mathbf{w}_i(t) = \mathbf{w}_i(t-1) - \eta g_i(\mathbf{w}_i(t); \mathcal{D}_i(t)), \quad t \in \mathcal{T}_k. \tag{6}$$

**Model Synchronization:** At time $t = k\tau$, after receiving the delayed global model $\mathbf{w}(t - \Delta)$ from the cloud server, each edge device performs one additional local SGD update followed by synchronization. During synchronization, each edge device performs local update by replacing its local model with a combination of the global and local model with the global/local *combiner weight* $\alpha(k) \in (0, 1]$. The expression for the local model after synchronization is given by

$$\mathbf{w}_i(t) = \alpha_t(k)\mathbf{w}(t - \Delta) \\ + (1 - \alpha_t(k)) \left[ \mathbf{w}_i(t-1) - \eta g_i(\mathbf{w}_i(t-1); \mathcal{D}_i(t)) \right], \tag{7}$$

where $\alpha_t(k)$ is the weight assigned to the global model:

$$\alpha_t(k) = \begin{cases} \alpha(k), & t = k\tau, k \in \{0, 1, ..., K-1\} \\ 0, & \text{otherwise} \end{cases}. \tag{8}$$

Let $\widehat{\boldsymbol{\alpha}} = \{\alpha(0), ..., \alpha(K)\}$ be the set of *combiner weights* across the global aggregation instances. $\alpha = 1$ corresponds to standard FL.

At the $K$-th global aggregation, the server chooses the best $\mathbf{w}(t)$ it has found thus far. Since the server only has access to the global model at $t = k\tau - \Delta$, the model selected at $K$ is

$$\mathbf{w}^K = \min_{\mathbf{w} \in \mathcal{W}} F(\mathbf{w}), \tag{9}$$

with $\mathcal{W} \triangleq \{\mathbf{w}(k\tau - \Delta), k = 0, 1, ..., K-1\}$.

**Algorithm 1:** Stochastic Federated Delayed Averaging

---

**Input:** $\widehat{\alpha}, \tau, \mathcal{I}, T$
**Output:** $\mathbf{w}^K$
Initialize $\mathbf{w}_i(-\Delta), \ \forall i$;
**for** $k = 0 : K - 1$ **do**
    **for** $t = k\tau - \Delta + 1 : (k+1)\tau - \Delta$ **do**
        **for** $i \in \mathcal{I}$ **do**
            **if** $t = (k+1)\tau - \Delta$ **then**
                | Each device $i$ sends $\mathbf{w}_i$ to the server
            **else**
                | Device $i \in \mathcal{I}$ updates its model using (7)
        **end**
        **if** $t = (k+1)\tau - \Delta/2$ **then**
            // Procedure at the cloud server
            Compute $\mathbf{w}((k+1)\tau - \Delta)$ and send it to
            the edge for synchronization and update
            $\mathbf{w}^K$ with (9)
    **end**
**end**

---

## III. Convergence Analysis of StoFedDelAv

In this section, we explore the optimality gap between the model chosen at the latest global aggregation $K$ and the optimal model. We then obtain the optimal model combiner weight. *All the proofs can be found in our online technical report [15].*

**Definition 1.** *The local data variability of device $i$ is measured via $\Theta_i \geq 0, \forall i$, satisfying* $\|\nabla f_i(\boldsymbol{x}_1, y_1; \boldsymbol{w}) - \nabla f_i(\boldsymbol{x}_2, y_2; \boldsymbol{w})\| \leq \Theta_i \|\boldsymbol{x}_1 - \boldsymbol{x}_2\|, \ \forall (\boldsymbol{x}_1, y_1), (\boldsymbol{x}_2, y_2) \in \mathcal{D}_i$.

**Definition 2.** *For $k \in \{0, ..., K-1\}$, the centralized GD during $t \in \{k\tau - \Delta + 1, ..., (k+1)\tau - \Delta\}$ is defined as $\boldsymbol{c}_k(t) = \boldsymbol{c}_k(t-1) - \eta \nabla F(\boldsymbol{c}_k(t-1))$ initialized such that $\boldsymbol{c}_k(k\tau - \Delta) = \boldsymbol{w}(k\tau - \Delta)$.*

We now characterize the variance of SGD noise in (5):

**Lemma 1.** *Using Definition 1, the SGD noise at each local iteration $t$ at each device $i$ can be upper bounded as follows:*

$$\mathbb{E}\left[\|\nu_i(t)\|^2\right] \leq \left(1 - |\mathcal{D}_i(t)|/N_i\right) 2 \left(\Theta_i S_i\right)^2 / |\mathcal{D}_i(t)|, \quad (10)$$

*where $S_i^2$ is the sample variance of data at device $i$.*

Since the the minibatch size (i.e, $|\mathcal{D}_i(t)|, \forall i$ in the above definition) is fixed during each local training interval and only varies across global aggregations, with some abuse of notation, we replace $t$ with $k$ in the above definition and express the SGD noise during period $k$, using Jensen's inequality, as

$$\mathbb{E}\left[\|\nu_i(k)\|\right] \leq \Theta_i S_i \sqrt{2\left(1 - n_i(k)/N_i\right)/n_i(k)}. \quad (11)$$

In Theorem 1, we bound the loss gap, i.e., $F(\mathbf{w}^K) - F(\mathbf{w}^\star)$:

**Theorem 1.** *If $\eta < \frac{2}{\beta}$, then under Assumption 1, we have*

$$F(\mathbf{w}^K) - F(\mathbf{w}^\star) \leq \frac{1}{2\eta\phi T} + \sqrt{\frac{1}{4\eta^2\phi^2 T^2} + \frac{L\Psi(\widehat{\alpha})}{\eta\phi T}} + L\Psi(\widehat{\alpha})$$

$$\triangleq \mathcal{L}(\{n_i(k)\}_{i \in \mathcal{I}, 1 \leq k \leq K}), \quad (12)$$

$$\Psi(\widehat{\alpha}) \triangleq \sum_{k=1}^{K} \psi(\alpha(k), k), \quad (13)$$

$$\begin{aligned} \psi(\alpha(k),k) &= \mathbb{E}[\|\mathbf{w}((k+1)\tau - \Delta) - \mathbf{c}_k((k+1)\tau - \Delta)\| \\ &\leq (1 - \alpha(k))\epsilon(k)([1 + \eta\beta]^\tau - 1) \\ &\quad + (1 - \alpha(k))h(\tau, k) + \alpha(k)h(\tau - \Delta, k) \\ &\quad + \alpha(k)\eta\Delta L[1 + \eta\beta]^{\tau - \Delta} + \eta\sigma(k)[\tau - \alpha(k)\Delta], \end{aligned} \quad (14)$$

$$h(x,k) \triangleq ((\delta + \sigma(k))/\beta)[(1 + \eta\beta)^x - 1] - \eta(\delta + \sigma(k))x, \quad (15)$$

$$\epsilon(k) \triangleq (1 - (1 - \alpha(k))^k)\left[2\eta(L + \sigma(k))\left(\tau/\alpha(k) - \Delta\right)\right], \quad (16)$$

$$\sigma(k) \triangleq \sum_i \rho_i E[\|\nu_i(k)\|] = \sum_i \rho_i S_i \Theta_i \sqrt{2\left(N_i - n_i(k)\right)/\left(N_i n_i(k)\right)}. \quad (17)$$

The optimality gap in (12) decreases as $T$ increases. More explicitly, as $T, K \to \infty$, $F(\mathbf{w}^K) - F(\mathbf{w}^\star)$ is determined exclusively by $\Psi(\widehat{\alpha})$ terms in (12). Critical to the understanding of the behavior of the optimality gap defined in (12) is $\Psi(\widehat{\alpha})$, which comprises terms $\psi(\alpha(k), k)$ given by (14). $\Psi(\widehat{\alpha})$ ultimately defines the discrepancy between the global model and the theoretical centralized model on one aggregation period.

It is important to note that the last term of (15) (i.e. the term with a negative sign) is decreasing with respect to (w.r.t.) gradient dissimilarity and noise. In most contexts, however, this term is counteracted by the rest of the terms in $\psi(\alpha(k), k)$ that are increasing w.r.t. SGD noise and gradient dissimilarity.

Crucial to the minimization of (12) is the proper choice of $\alpha(k)$. Although the behavior of the expression in (14) is non-trivial to analyze, we experimentally observe in Fig. 2(a) (Sec. V) that $\psi(\alpha(k), k)$ is convex as a function of $\alpha(k) \in (0, 1]$, implying $[F(\mathbf{w}^K) - F(\mathbf{w}^\star) \propto \sqrt{\Psi(\widehat{\alpha})} + \Psi(\widehat{\alpha})$ can be minimized by minimizing each $\psi(\alpha(k), k)$ since $\psi(\alpha(k), k)$'s are independent according to (14). In particular, each $\psi(\alpha(k), k)$ is the solution to the optimization problem

$$\alpha^\star(k) = \underset{\alpha(k) \in (0,1)}{\arg\min} \ \psi(\alpha(k), k), \quad (18)$$

where $\psi(\alpha(k), k)$ is given by (14). Since the closed-form solution of the above problem is non-trivial, this problem can be solved using numerical methods given the bounded range of $\alpha(k)$. Nevertheless, given (14) optimizing over $\psi(\alpha(\infty), \infty)$ would give us the following closed-form solution:

$$\begin{aligned} \alpha^\star(\infty) &= \min\left(1, \sqrt{2\eta\tau(L + \sigma(\infty))[(1 + \eta\beta)^\tau - 1]/A}\right) \\ A &= 2\eta\Delta(L + \sigma(\infty))[(1 + \eta\beta)^\tau - 1] + \eta\Delta L(1 + \eta\beta)^{\tau - \Delta} \\ &\quad - ((\delta + \sigma(\infty))/\beta)(1 + \eta\beta)^{\tau - \Delta}[(1 + \eta\beta)^\Delta - 1] + \eta\delta\Delta. \end{aligned} \quad (19)$$

In practice, to avoid a numerical method, one can use (19) for each $\alpha^\star(k)$ with using $\sigma(k)$ instead of $\sigma(\infty)$ in (19).

## IV. Network Optimization Problem

In this section, we first formulate a problem to jointly minimize energy, delay, and model loss in Sec. IV-A. We then rework the problem and solve it in Sec. IV-B.

### A. Problem Formulation

For period $k$, let $E^{\mathsf{Cmp}}(k)$ be the energy required to compute the gradient over a minibatch of data, $E^{\mathsf{Tx}}(k)$ be the energy required for model transmission, $T^{\mathsf{Cmp}}(k)$ be the computation time, $T^{\mathsf{Tx}}(k)$ be the model transmission time, $Q$ be the number of bits per model, $p_i(k)$ be the transmit power of device $i$, and $R_i(k)$ be the data rate between device $i$ and the BS.

We formulate the following problem to optimize a trade-off between energy consumption, delay, and model performance:

$$\mathcal{P}: \min_{\{\boldsymbol{n}(k)\}_{k=1}^K} \sum_{k=1}^K \Big[ c_1\big[E^{\mathsf{Cmp}}(k)+E^{\mathsf{Tx}}(k)\big] + c_2\big[T^{\mathsf{Cmp}}(k)+T^{\mathsf{Tx}}(k)\big]\Big]$$
$$+ c_3\mathcal{L}(\{n_i(k)\}_{i\in\mathcal{I}, 1\le k\le K}) \qquad (20)$$

s.t.

(C1) $E^{\mathsf{Cmp}}(k) = \sum_{i\in\mathcal{I}} E_i^{\mathsf{Cmp}}(k),$

(C2) $E^{\mathsf{Tx}}(k) = \sum_{i\in\mathcal{I}} E_i^{\mathsf{Tx}}(k),$

(C3) $\sum_{k=1}^K E_i^{\mathsf{Cmp}}(k) + E_i^{\mathsf{Tx}}(k) \le E_i^{\mathsf{Batt}},\ \forall i\in\mathcal{I},$

(C4) $E_i^{\mathsf{Cmp}}(k) = \gamma_i d_i \tau n_i(k)\varrho_i^2/2,\ \forall i\in\mathcal{I},$

(C5) $E_i^{\mathsf{Tx}}(k) = p_i(k)Q/R_i(k),\ \forall i\in\mathcal{I},$

(C6) $T^{\mathsf{Cmp}}(k) = \max_{i\in\mathcal{I}} \tau d_i n_i(k)/\varrho_i,$

(C7) $T^{\mathsf{Tx}}(k) = \max_{i\in\mathcal{I}} Q/R_i(k),$

(C8) $\mathcal{L}(\{\boldsymbol{n}(k)\}_{k=1}^K) = F(\mathbf{w}^K) - F(\mathbf{w}^\star)$ (see (12)),

(C9) $0 \le n_i(k) \le N_i,\ \forall i\in\mathcal{I},$

where $\boldsymbol{n}(k) = \{n_i(k)\}_{i\in\mathcal{I}}$ is the collection of minibatch sizes of the devices over the training interval, and constants $c_1, c_2, c_3 \ge 0$ weigh the importance of the objective terms.

Constraints **C1** and **C2** are, respectively, the total computation and transmission energy consumption during each global aggregation. **C3** limits the amount of energy device $i$ can consume over $K$ according to its battery $E_i^{\mathsf{Batt}}$. **C4** constrains the computation energy of $i$, where $\gamma_i$ is its effective CPU capacitance, $d_i$ is the number of CPU cycles needed to process one datapoint, and $\varrho_i$ is the CPU clocking frequency [5], [8]. **C5** represents the energy needed for transmission, and constraints **C6** and **C7** are the computation and transmission time, respectively, for the network. **C8** constrains the loss gap to its upper bound, and constraint **C9** ensures $\mathcal{P}$'s feasibility.

### B. Geometric Programming-based Optimization

Problem $\mathcal{P}$ is non-convex, particularly due to the behavior of $\mathcal{L}$ in the objective function. However, by fixing the value of $\alpha(k)$, the problem reduces to a signomial programming (SP) problem [16]. While this is still NP-hard in general, the resulting SP can be solved via the method of posynomial condensation and penalty functions [17]. We thus transform $\mathcal{P}$ into an iterative problem in which at each iteration $\ell$, a convex problem is obtained via logarithmic change of optimization variables (c.o.v.), the solution of which is used to determine the value of $\widehat{\boldsymbol{\alpha}}$ using (19). In particular, we write the problem as an optimization problem with a *posynomial* objective function subject to equality on *monomials* and inequality on *posynomials*, which admits the format of geometric programming (GP) [16]. As a result, at each iteration $\ell$, we aim to find the solution to the following optimization problem, which can undergo a logarithmic c.o.v. and be reduced to a convex problem:

$$\widehat{\mathcal{P}}: \min_{\boldsymbol{y}} \sum_{k=1}^K \Big[ c_1\big[E^{\mathsf{Cmp}}(k)+E^{\mathsf{Tx}}(k)\big] + c_2\big[T^{\mathsf{Cmp}}(k)+T^{\mathsf{Tx}}(k)\big]\Big]$$
$$+ c_3\mathcal{L}(\{\boldsymbol{n}(k)\}_{k=1}^K) \qquad (29)$$
$$+ w_1 s_1 + \sum_{k=1}^K \Big[ \sum_{j=2}^4 w_j(k)s_j(k) + \sum_{i\in\mathcal{I}} w_5(k,i)s_5(k,i) \Big]$$

s.t.

($\widehat{\mathbf{C}}$1) $\sum_{i\in\mathcal{I}} E_i^{\mathsf{Cmp}}(k)/E^{\mathsf{Cmp}}(k) \le 1$

($\widehat{\mathbf{C}}$2) $\sum_{i\in\mathcal{I}} E_i^{\mathsf{Tx}}(k)/E^{\mathsf{Tx}}(k) \le 1$

($\widehat{\mathbf{C}}$3) $\sum_{k=1}^K \big(E_i^{\mathsf{Cmp}}(k) + E_i^{\mathsf{Tx}}(k)\big)/E_i^{\mathsf{Batt}} \le 1, \forall i\in\mathcal{I}$

($\widehat{\mathbf{C}}$4) $\gamma_i d_i \tau n_i(k)\varrho_i^2/(2E_i^{\mathsf{Cmp}}(k)) = 1,\ \forall i\in\mathcal{I}$

($\widehat{\mathbf{C}}$5) $p_i Q/(E_i^{\mathsf{Tx}}(k)R_i) = 1,\ \forall i\in\mathcal{I}$

($\widehat{\mathbf{C}}$6) $\tau d_i n_i(k)/(T^{\mathsf{Cmp}}(k)\varrho_i) \le 1, \forall i\in\mathcal{I}$

($\widehat{\mathbf{C}}$7) $Q/(T^{\mathsf{Tx}}(k)R_i) \le 1, \forall i\in\mathcal{I}$

($\widehat{\mathbf{C}}$8.1) $\mathcal{L}^{-1}\big[m_1 + P_1 + L\Psi(\widehat{\boldsymbol{\alpha}})\big] \le 1$

($\widehat{\mathbf{C}}$8.2) $(m_1 2\eta\phi T)^{-1} \le 1$

($\widehat{\mathbf{C}}$8.3) $P_1^{-2}(m_2 + m_3\Psi(\widehat{\boldsymbol{\alpha}})) \le 1$

($\widehat{\mathbf{C}}$8.4) $m_2^{-1}(1/(2\eta\phi T))^2 = 1$

($\widehat{\mathbf{C}}$8.5) $Lm_3^{-1}/(\eta\phi T) = 1$

($\widehat{\mathbf{C}}$8.6) $\Psi^{-1}(\widehat{\boldsymbol{\alpha}}) \sum_{k=1}^K \psi(\alpha(k),k) \le 1$

($\widehat{\mathbf{C}}$8.7) $s_1^{-1}\Psi(\widehat{\boldsymbol{\alpha}})/\widehat{f}_1(\boldsymbol{y},\widehat{\boldsymbol{\alpha}};\ell) \le 1$

($\widehat{\mathbf{C}}$8.8) $\psi^{-1}(k)[\alpha(k)\mathsf{B}_4(k)\epsilon(k)\mathsf{B}_1 + \mathsf{B}_4(k)\mathsf{h}_1(k)$
$\qquad + \alpha(k)\mathsf{h}_2(k) + \alpha(k)\eta\Delta L\mathsf{B}_5 + \eta\sigma(k)\mathsf{B}_6(k)] \le 1$

($\widehat{\mathbf{C}}$8.9) $s_2^{-1}(k)\psi(k)/\widehat{f}_2(\boldsymbol{y},k;\ell) \le 1$

($\widehat{\mathbf{C}}$8.10) $(\mathsf{h}_1^{-1}(k)\mathsf{B}_1\delta\beta^{-1} + \mathsf{h}_1^{-1}(k)\mathsf{B}_1\sigma(k)\beta^{-1})/\widehat{f}_3(\boldsymbol{y},\tau,k,1;\ell) \le 1$

($\widehat{\mathbf{C}}$8.11) $\dfrac{s_3^{-1}(k)\big[1 + \mathsf{h}_1^{-1}(k)\eta\delta\tau + \mathsf{h}_1^{-1}(k)\eta\sigma(k)\tau\big]}{\widehat{f}_4(\boldsymbol{y},\tau,k,1;\ell)} \le 1$

($\widehat{\mathbf{C}}$8.12) $\dfrac{\mathsf{h}_2^{-1}(k)\mathsf{B}_2\delta\beta^{-1} + \mathsf{h}_2^{-1}(k)\mathsf{B}_2\sigma(k)\beta^{-1}}{\widehat{f}_3(\boldsymbol{y},\tau-\Delta,k,2;\ell)} \le 1$

($\widehat{\mathbf{C}}$8.13) $\dfrac{s_4^{-1}(k)\big[1 + \mathsf{h}_2^{-1}(k)\eta\delta\mathsf{B}_7 + \mathsf{h}_2^{-1}(k)\eta\sigma(k)\mathsf{B}_7\big]}{\widehat{f}_4(\boldsymbol{y},\tau-\Delta,k,2;\ell)} \le 1$

($\widehat{\mathbf{C}}$8.14) $\epsilon(k)^{-1}\mathsf{B}_3(k)2\eta(L+\sigma(k))\alpha(k)^{-1}/\widehat{f}_5(\boldsymbol{y},k;\ell) \le 1$

($\widehat{\mathbf{C}}$8.15) $\Big(s_5^{-1}(k)\big[1 + \epsilon(k)^{-1}\mathsf{B}_3(k)2\eta(L+\sigma(k))\Delta\big]\Big)/\widehat{f}_6(\boldsymbol{y},k;\ell) \le 1$

($\widehat{\mathbf{C}}$8.16) $\sum_{i\in\mathcal{I}} \rho_i S_i \Theta_i \sqrt{2}P_i(k)/\sigma(k) \le 1$

($\widehat{\mathbf{C}}$8.17) $s_6^{-1}(k)\sigma(k)/\widehat{f}_7(\boldsymbol{y},k;\ell) \le 1$

($\widehat{\mathbf{C}}$8.18) $P_i^2(k)n_i(k) + n_i(k)N_i^{-1} \le 1,\ \forall i\in\mathcal{I}$

($\widehat{\mathbf{C}}$8.19) $s_7^{-1}(k,i)/\widehat{f}_8(\boldsymbol{y},k;\ell) \le 1,\ \forall i\in\mathcal{I}$

($\widehat{\mathbf{C}}$9) $0 \le n_i(k) \le N_i,\ \forall i\in\mathcal{I}$

($\widehat{\mathbf{C}}$10) $\Big\{s_1, \{\boldsymbol{s}_j(k)\}_{2\le j\le 6, 1\le k\le K}, \{\boldsymbol{s}_7(k,i)\}_{i\in\mathcal{I}, 1\le k\le K}\Big\} \ge 1$

**Variables:** $\boldsymbol{y} \triangleq \Big\{P_1, \Psi(\widehat{\boldsymbol{\alpha}}), \big\{T^{\mathsf{Cmp}}(k), T^{\mathsf{Tx}}(k), E^{\mathsf{Cmp}}(k), E^{\mathsf{Tx}}(k)\big\}_{k=1}^K,$
$\big\{\{\boldsymbol{n}(k)\}, \sigma(k), \mathsf{h}_1(k), \mathsf{h}_2(k), \epsilon(k), \psi(k), \{P_i(k)\}_{i\in\mathcal{I}}\big\}_{k=1}^K,$
$s_1, \{\boldsymbol{s}_j(k)\}_{2\le j\le 6, 1\le k\le K}, \{\boldsymbol{s}_7(k,i)\}_{i\in\mathcal{I}, 1\le k\le K}\Big\},$

where $\mathsf{h}_1(k) = h(\tau,k)$, $\mathsf{h}_2(k) = h(\tau-\Delta,k)$, $\mathsf{B}_1 = (1+\eta\beta)^\tau - 1$, $\mathsf{B}_2 = (1+\eta\beta)^{\tau-\Delta} - 1$, $\mathsf{B}_3(k) = (1-(1-\alpha(k))^k)$, $\mathsf{B}_4(k) = (1-\alpha(k))$, $\mathsf{B}_5 = (1+\eta\beta)^{\tau-\Delta}$, $\mathsf{B}_6(k) = \tau - \alpha(k)\Delta$, and $\mathsf{B}_7 = (\tau - \Delta)$. $\mathsf{B}_j \ge 0, \forall j$. The $\{s_j\}$ terms are added to expand the solution space of each iteration that will be forced to converge to 1 when the problem is solved using the penalty terms (i.e., $w_j \gg 1,\ \forall j$). The terms $\widehat{f}_x(\boldsymbol{y}, ...; \ell)$

$$f_1(\boldsymbol{y}, \widehat{\boldsymbol{\alpha}}) = \sum_{k=1}^{K} \psi(\alpha(k), k) \rightarrow f_1(\boldsymbol{y}, \widehat{\boldsymbol{\alpha}}) \geq \widehat{f}_1(\boldsymbol{y}, \widehat{\boldsymbol{\alpha}}; \ell) \triangleq \prod_{k=1}^{K} \left( \frac{\psi(\alpha(k), k) f_1(\boldsymbol{y}, \widehat{\boldsymbol{\alpha}})^{[\ell-1]}}{\psi(\alpha(k), k)^{[\ell-1]}} \right)^{\frac{\psi(\alpha(k), k)^{[\ell-1]}}{f_1(\boldsymbol{y}, \widehat{\boldsymbol{\alpha}})^{[\ell-1]}}} \tag{21}$$

$$f_2(\boldsymbol{y}, k) = \underbrace{\alpha(k) \mathsf{B}_4(k) \epsilon(k) \mathsf{B}_1}_{q_{2,1}} + \underbrace{\mathsf{B}_4(k) \mathsf{h}_1(k)}_{q_{2,2}} + \underbrace{\alpha(k) \mathsf{h}_2(k)}_{q_{2,3}} + \underbrace{\alpha(k) \eta \Delta L \mathsf{B}_5}_{q_{2,4}} + \underbrace{\eta \sigma(k) \mathsf{B}_6(k)}_{q_{2,5}} \rightarrow$$

$$\tag{22}$$

$$f_2(\boldsymbol{y}, k) \geq \widehat{f}_2(\boldsymbol{y}, k; \ell) \triangleq \prod_{j=1}^{5} \left( \frac{q_{2,j} f_2(\boldsymbol{y}, k)^{[\ell-1]}}{q_{2,j}^{[\ell-1]}} \right)^{\frac{q_{2,j}^{[\ell-1]}}{f_2(\boldsymbol{y}, k)^{[\ell-1]}}}$$

$$f_3(\boldsymbol{y}, x, k, i) = \underbrace{1}_{q_{3,1}} + \underbrace{\mathsf{h}_i^{-1}(k) \eta \delta x}_{q_{3,2}} + \underbrace{\mathsf{h}_i^{-1}(k) \eta \sigma(k) x}_{q_{3,3}} \rightarrow f_3(\boldsymbol{y}, x, k, i) \geq \widehat{f}_3(\boldsymbol{y}, x, k, i; \ell) \triangleq \prod_{j=1}^{3} \left( \frac{q_{3,j} f_3(\boldsymbol{y}, x, k, i)^{[\ell-1]})}{q_{3,j}^{[\ell-1]}} \right)^{\frac{q_{3,j}^{[\ell-1]}}{f_3(\boldsymbol{y}, x, k, i)^{[\ell-1]}}} \tag{23}$$

$$f_4(\boldsymbol{y}, x, k, i) = \underbrace{\mathsf{h}_i^{-1}(k) \mathsf{B}_i \delta \beta^{-1}}_{q_{4,1}} + \underbrace{\mathsf{h}_i^{-1}(k) \mathsf{B}_i \sigma(k) \beta^{-1}}_{q_{4,2}} \rightarrow f_4(\boldsymbol{y}, x, k, i) \geq \widehat{f}_4(\boldsymbol{y}, x, k, i; \ell) \triangleq \prod_{j=1}^{2} \left( \frac{q_{4,j} f_4(\boldsymbol{y}, x, k, i)^{[\ell-1]}}{q_{4,j}^{[\ell-1]}} \right)^{\frac{q_{4,j}^{[\ell-1]}}{f_4(\boldsymbol{y}, x, k, i)^{[\ell-1]}}} \tag{24}$$

$$f_5(\boldsymbol{y}, k) = \underbrace{1}_{q_{5,1}} + \underbrace{\epsilon^{-1}(k) \mathsf{B}_3(k) 2 \eta L \Delta}_{q_{5,2}} + \underbrace{\epsilon^{-1}(k) \mathsf{B}_3(k) 2 \eta \sigma(k) \Delta}_{q_{5,3}} \rightarrow f_5(\boldsymbol{y}, k) \geq \widehat{f}_5(\boldsymbol{y}, k; \ell) \triangleq \prod_{j=1}^{3} \left( \frac{q_{5,j} f_5(\boldsymbol{y}, k)^{[\ell-1]}}{q_{5,j}^{[\ell-1]}} \right)^{\frac{q_{5,j}^{[\ell-1]}}{f_5(\boldsymbol{y}, k)^{[\ell-1]}}} \tag{25}$$

$$f_6(\boldsymbol{y}, k) = \underbrace{\epsilon^{-1}(k) \mathsf{B}_3(k) 2 \eta L \tau \alpha(k)^{-1}}_{q_{6,1}} + \underbrace{\epsilon^{-1}(k) \mathsf{B}_3(k) 2 \eta \sigma(k) \tau \alpha(k)^{-1}}_{q_{6,2}} \rightarrow f_6(\boldsymbol{y}, k) \geq \widehat{f}_6(\boldsymbol{y}, k; \ell) \triangleq \prod_{j=1}^{2} \left( \frac{q_{6,j} f_6(\boldsymbol{y}, k)^{[\ell-1]}}{q_{6,j}^{[\ell-1]}} \right)^{\frac{q_{6,j}^{[\ell-1]}}{f_6(\boldsymbol{y}, k)^{[\ell-1]}}} \tag{26}$$

$$f_7(\boldsymbol{y}, k) = \sum_{j \in \mathcal{I}} \rho_j S_j \Theta_j \sqrt{2} P_j(k) \rightarrow f_7(\boldsymbol{y}, k) \geq \widehat{f}_7(\boldsymbol{y}, k; \ell) \triangleq \prod_{j \in \mathcal{I}} \left( \frac{(\rho_j S_j \Theta_j \sqrt{2} P_j(k)) f_7(\boldsymbol{y}, k)^{[\ell-1]}}{\left\{ \rho_j S_j \Theta_j \sqrt{2} P_j(k) \right\}^{[\ell-1]}} \right)^{\frac{\left\{ \rho_j S_j \Theta_j \sqrt{2} P_j(k) \right\}^{[\ell-1]}}{f_7(\boldsymbol{y}, k)^{[\ell-1]}}} \tag{27}$$

$$f_8(\boldsymbol{y}, k, i) = \underbrace{P_i^2(k) n_i(k)}_{q_{8,1}} + \underbrace{n_i(k) N_i^{-1}}_{q_{8,2}} \rightarrow f_8(\boldsymbol{y}, k, i) \geq \widehat{f}_8(\boldsymbol{y}, k, i; \ell) \triangleq \prod_{j=1}^{2} \left( \frac{q_{8,j} f_8(\boldsymbol{y}, k, i)^{[\ell-1]}}{q_{8,j}^{[\ell-1]}} \right)^{\frac{q_{8,j}^{[\ell-1]}}{f_8(\boldsymbol{y}, k, i)^{[\ell-1]}}} \tag{28}$$

approximate posynomial denominators in $\widehat{\mathcal{P}}$ as monomials to satisfy the requirements of GP, and are outlined in (21)-(28). As the iterations progress, these approximations converge towards the value of the posynomial they represent. After convergence, (19) is applied with $\sigma(k)^{[\ell]}$ to update $\widehat{\boldsymbol{\alpha}}$ and $\mathsf{B}_{\{3,4,5\}}(k)$. A new problem is then solved given the values of these variables, and this *alternative* process is continued upon convergence.

In $\widehat{\mathcal{P}}$, constraints $\widehat{\mathrm{C}}1$-$\widehat{\mathrm{C}}5$ are naturally obtained from problem $\mathcal{P}$'s C1-C5 into ones which fit a geometric programming (GP) paradigm. Constraints $\widehat{\mathrm{C}}6$ and $\widehat{\mathrm{C}}7$ stem from the fact that dividing $\mathcal{P}$'s C6-C7 computation/transmission times by the maximum computation/transmission time across the network will upper-bound the constraint to 1. $\widehat{\mathcal{P}}$'s constraints $\widehat{\mathrm{C}}8.\{1, 2, ..., 19\}$ develop the transformation of the loss gap of (12) into a series of constraints in the form of inequalities on posynomials, which is desired in GP programming to have convergence to a Karush–Kuhn–Tucker condition of $\mathcal{P}$ [16].

## V. EXPERIMENTAL RESULTS

**Experimental Setup:** We consider an edge network of $N = 5$ devices realized according to the parameters described in Table I. Sets of $N$ parameters are uniformly generated then sorted for $\gamma_i$ and $d_i$ (i.e. $\boldsymbol{\gamma} = \{\gamma_1, ..., \gamma_5\}$ and $\boldsymbol{d} = \{d_1, ..., d_5\}$), such that $\gamma_1 = \arg\min\{\gamma_i\}_{i=1}^{5}$, $d_1 = \arg\min\{d_i\}_{i=1}^{5}$ and $\gamma_5 = \arg\max\{\gamma_i\}_{i=1}^{5}$, $d_5 = \arg\max\{d_i\}_{i=1}^{5}$. The first device is modeled using $\gamma_1$ and $d_1$ for its CPU capacitance and number of CPU cycles per datapoint, respectively, making it the most resource-efficient device for data computation; the second device uses $\gamma_2$, $d_2$, and so on. CVX is used to solve

the convex problem at each iteration of $\widehat{\mathcal{P}}$. Each plot in Fig. 2 shows the average of 20 randomized network initializations.

TABLE I: Parameter settings for experiments.

| Parameter(s) | Value / Range |
|---|---|
| $c_1, c_2, c_3$ | $1 \times 10^{-4}, 1 \times 10^{-3}, 2.5 \times 10^6$ |
| $E_i^{\mathrm{Batt}} \forall i$ | $7.5 \times 10^6 (J)$ |
| $n_i(k)$ | $[1, 25]$ |
| $\varrho_i$ | $1 \times 10^6 (Hz)$ |
| $d_i$ | $600 \leq d_i \leq 640$ |
| $\gamma_i$ | $[4 \times 10^{-12}, 6.5 \times 10^{-12}](F)$ |
| $p_i, R_i, Q$ | $0.1 (W) \forall i, 1.0 (Mbps) \forall i, 16 (kbits)$ |
| $\Theta_i, S_i, \delta$ | $2.0 \forall i, 0.2 \forall i, 0.5$ |
| $\eta, \beta, L, \phi$ | $0.02, 1, 25, 0.025$ |
| $\tau, \Delta, K$ | $20, 19, 15$ |
| $w_1, w_{\{2,...,6\}}(k), w_7(k, i)$ | $100000, 100000, 1000000$ |

**Minibatch Optimization:** We first look at minibatch size, which ultimately determines time, energy, and loss across the training interval. Since $\epsilon(k)$ in (16) becomes more dependent on noise as training progresses due to the term $(1 - (1 - \alpha(k))^k$, minibatch size should theoretically increase non-linearly over time. This is corroborated in Fig. 2(b). It can be seen that minibatch size for the devices follows their relative precedence, such that the best edge device, 1, possesses the largest minibatch, 2 the second largest, and so on. Better devices show larger differences between their initial minibatch size and their latest, with device 1 experiencing a nearly 25% increase. This trend reveals that the better edge devices save energy in early training stages for later on when SGD noise is more impactful on the machine learning loss.

**Energy and Minibatch:** In Fig 1(c), we depict average minibatch size across the network while varying the energy
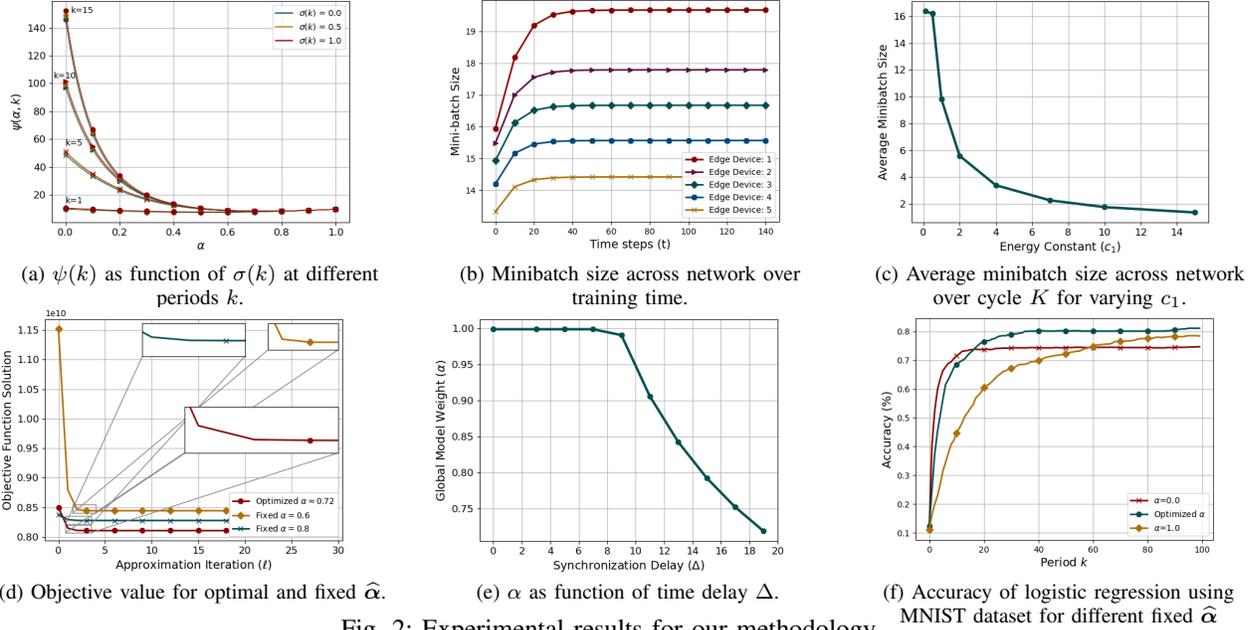
(a) $\psi(k)$ as function of $\sigma(k)$ at different periods $k$.

(b) Minibatch size across network over training time.

(c) Average minibatch size across network over cycle $K$ for varying $c_1$.

(d) Objective value for optimal and fixed $\widehat{\alpha}$.

(e) $\alpha$ as function of time delay $\Delta$.

(f) Accuracy of logistic regression using MNIST dataset for different fixed $\widehat{\alpha}$

Fig. 2: Experimental results for our methodology.

constant, $c_1$ in the objective function of $\mathcal{P}$. The results show that the precedence assigned to energy and the average minibatch size across the network for the complete training cycle exhibit a steep ramp-down from $c_1 \in (0, 1)$.

**Impact and Behavior of $\alpha(k)$:** By allowing the network to choose $\widehat{\alpha}$ per (19), the value for the objective function of the problem $\widehat{\mathcal{P}}$ drops meaningfully, as seen in Fig. 2(d). We determined numerically that the difference between the calculated optimal value of $\alpha(0)$ and that found using (19) was about 0.16. Subsequent values of $\alpha(k > 0)$ were effectively identical to the numerically optimal value, thus proving the efficacy of the proposed method. It is worth noting that this is feasible for the iterative GP approach, as previous values for $\sigma(k)$ can be used, but in real-time this may not be the case.

$\alpha(k)$ is also heavily dependent on delay as shown in Fig. 2(e), where the vertical axis represents the average of elements in $\widehat{\alpha}$ and the horizontal axis $\Delta$. This shows that the proportionality between $\tau$ and $\Delta$ should be carefully considered when choosing $\widehat{\alpha}$. As $\Delta \to 0$, $\alpha(k) \to 1$, as is expected in ideal `FedAvg`.

Figure 2(f) illustrates the impact of $\alpha$ on model accuracy in a real model, trained on the MNIST dataset using $\alpha = 0.0, 1.0$ and $\alpha_{opt}$. It is readily apparent that optimizing $\alpha$ plays a key role in convergence under the `StoFedDelAv` paradigm.

## VI. CONCLUSION AND FUTURE WORK

We proposed a novel methodology for optimizing federated learning implementations over edge networks while explicitly taking into account device-server communication delay and device computation heterogeneity. The loss optimality gap was considered across a training cycle to characterize the performance of the network. We formulated an optimization problem aiming to find the minibatch size of the devices across the training interval to optimize a trade-off between energy consumption, time required to train the model, and ML model performance. This problem was optimized using

an iterative geometric programming-based approach to find the ideal minibatch size for each device across the network. Future works will focus on improving the network and training efficiency, namely distributed device orchestration and delay-aware device sampling. These approaches will enable networks to train models in a more time- and energy-efficient manner.

## REFERENCES

[1] M. Chiang and T. Zhang, "Fog and IoT: An overview of research opportunities," *IEEE Internet Thing J.*, vol. 3, no. 6, pp. 854–864, 2016.

[2] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. Y. Arcas, "Communication-efficient learning of deep networks from decentralized data," vol. 54, pp. 1273–1282, 2017.

[3] J. Konečný *et al.*, "Federated learning: Strategies for improving communication efficiency," in *Proc. NIPS Workshop*, 2016.

[4] H. H. Yang, Z. Liu, T. Q. S. Quek, and H. V. Poor, "Scheduling policies for federated learning in wireless networks," *IEEE Trans. Commun.*, vol. 68, no. 1, pp. 317–333, 2020.

[5] S. Wang, T. Tuor, T. Salonidis, K. K. Leung, C. Makaya, T. He, and K. Chan, "Adaptive federated learning in resource constrained edge computing systems," *IEEE J. Select. Areas Commun.*, vol. 37, no. 6, pp. 1205–1221, 2019.

[6] S. Hosseinalipour, C. G. Brinton, V. Aggarwal, H. Dai, and M. Chiang, "From federated to fog learning: Distributed machine learning over heterogeneous wireless networks," *IEEE Commun. Mag.*, vol. 58, no. 12, pp. 41–47, 2020.

[7] F. Haddadpour and M. Mahdavi, "On the convergence of local descent methods in federated learning," *arXiv preprint arXiv:1910.14425*, 2019.

[8] N. H. Tran, W. Bao, A. Zomaya, N. M. NH, and C. S. Hong, "Federated learning over wireless networks: Optimization model design and analysis," in *IEEE INFOCOM*, 2019, pp. 1387–1395.

[9] F. P.-C. Lin, S. Hosseinalipour, S. S. Azam, C. G. Brinton, and N. Michelusi, "Semi-decentralized federated learning with cooperative D2D local model aggregations," *IEEE J. Sel. Areas Commun.*, 2021.

[10] S. Hosseinalipour, S. S. Azam, C. G. Brinton, N. Michelusi, V. Aggarwal, D. J. Love, and H. Dai, "Multi-stage hybrid federated learning over large-scale D2D-enabled fog networks," *arXiv:2007.09511*, 2020.

[11] M. Chen, H. V. Poor, W. Saad, and S. Cui, "Convergence time minimization of federated learning over wireless networks," in *IEEE ICC*, 2020, pp. 1–6.

[12] Z. Yang, M. Chen, W. Saad, C. S. Hong, and M. Shikh-Bahaei, "Energy efficient federated learning over wireless communication networks," *arXiv preprint arXiv:1911.02417*, 2019.

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

[14] F. P.-C. Lin, C. G. Brinton, and N. Michelusi, "Federated learning with communication delay in edge networks," in *IEEE GLOBECOM*, 2020, pp. 1–6.

[15] "Technical report," https://www.cbrinton.net/icc-2022-tech.pdf.

[16] M. Chiang, *Geometric programming for communication systems.* now Publishers Inc., 2005.

[17] S. Hosseinalipour, A. Rahmati, D. Y. Eun, and H. Dai, "Energy-aware stochastic UAV-assisted surveillance," *IEEE Trans. Wireless Commun.*, vol. 20, no. 5, pp. 2820–2837, 2021.

## VII. APPENDIX A

### A. Proof of Lemma 1

**Lemma 1.** *For ease of manipulation in bounding equations using the triangular inequality, the noise can be defined as*

$$\mathbb{E}\left[\|\nu_i(k)\|\right] \leq S_i \Theta_i \sqrt{2} \sqrt{\frac{N_i - n_i(k)}{N_i n_i(k)}} \tag{30}$$

*Proof.* We begin by defining the variance of the gradients, $\widehat{\mathcal{S}_i^2}$. With $\boldsymbol{\lambda}_i$ and $S_i$ denoting the mean and sample variance of the device's datapoints, respectively, and using Definition 1, we say:

$$\widehat{\mathcal{S}_i^2} = \frac{\sum_{x_1 \in \mathcal{D}_i} \|\nabla f_i(x_1, y_1; \mathbf{w}) - \sum_{x_2 \in \mathcal{D}_i} \frac{\nabla f_i(x_2, y_2; \mathbf{w})}{N_i})\|^2}{N_i - 1}$$

$$= \frac{\sum_{x_1 \in \mathcal{D}_i} \frac{1}{N_i^2} \|N_i \nabla f_i(x_1, y_1; \mathbf{w}) - \sum_{x_2 \in \mathcal{D}_i} \nabla f_i(x_2, y_2; \mathbf{w})\|^2}{N_i - 1}$$

$$\leq \frac{\sum_{x_1 \in \mathcal{D}_i(k)} \frac{N_i - 1}{N_i^2} \sum_{x_2 \in \mathcal{D}_i} \|\nabla f_i(x_1, y_1; \mathbf{w}) - \nabla f_i(x_2, y_2; \mathbf{w})\|^2}{Z_1 - 1}$$

$$\leq \frac{\sum_{\mathbf{x}_1 \in \mathcal{D}_i} \frac{(N_i - 1)\Theta_i^2}{N_i^2} \sum_{\mathbf{x}_2 \in \mathcal{D}_i} \|\mathbf{x}_1 - \mathbf{x}_2\|^2}{N_i - 1}$$

$$\leq \frac{(N_i - 1)\Theta_i^2}{N_i^2} \frac{\sum_{x_1 \in \mathcal{D}_i} \sum_{x_2 \in \mathcal{D}_i} \|\mathbf{x}_1 - \mathbf{x}_2 + \boldsymbol{\lambda}_i - \boldsymbol{\lambda}_i\|^2}{N_i - 1}$$

$$= \frac{(N_i - 1)\Theta_i^2}{N_i^2} \times$$

$$\left[\frac{\sum_{x_1 \in \mathcal{D}_i} \sum_{x_2 \in \mathcal{D}_i} \left[\|\mathbf{x}_1 - \boldsymbol{\lambda}_i\|^2 + \|\mathbf{x}_2 - \boldsymbol{\lambda}_i\|^2 - 2(\mathbf{x}_1 - \boldsymbol{\lambda}_i)^\mathsf{T}(\mathbf{x}_2 - \boldsymbol{\lambda}_i)\right]}{N_i - 1}\right]$$

$$= \frac{(N_i - 1)\Theta_i^2}{N_i^2} \frac{N_i \sum_{x_1 \in \mathcal{D}_i} \|\mathbf{x}_1 - \boldsymbol{\lambda}_i\|^2 + N_i \sum_{x_2 \in \mathcal{D}_i} \|\mathbf{x}_2 - \boldsymbol{\lambda}_i\|^2}{N_i - 1}$$

$$= \frac{2(N_i - 1)\Theta_i^2 S_i^2}{N_i} \leq 2(\Theta_i S_i)^2, \tag{31}$$

where the first inequality is found using the Cauchy-Schwarz inequality, and the second to last line stems from the fact that $\sum_{x_1 \in \mathcal{D}_i}(\mathbf{x}_1 - \boldsymbol{\lambda}_i) = \mathbf{0}$.

We now look to the variance of the SGD noise itself. As defined in (11), the variance of the noise for any iteration $\ell$ is

$$\mathbb{E}[\|\nu_i(t)\|^2] = \left(1 - \frac{n_i(t)}{N_i}\right) \frac{\widehat{\mathcal{S}_i^2}}{n_i(t)}. \tag{32}$$

Using the above derivation of $\widehat{\mathcal{S}_i^2}$, we can upper-bound this as

$$\mathbb{E}[\|\nu_i(t)\|^2] \leq \left(1 - \frac{|\mathcal{D}(t)|}{N_i}\right) \frac{2(\Theta_i S_i)^2}{|\mathcal{D}(t)|}. \tag{33}$$

Since the minibatch size, $|\mathcal{D}_i(t)|, \forall i$, is fixed during each local training period (i.e. it only varies across global aggregations), with some abuse of notation we replace $t$ with $k$ in the above definition and express the SGD variance during period $k$ as

$$\mathbb{E}[\|\nu_i(k)\|^2] \leq \left(1 - \frac{n_i(k)}{N_i}\right) \frac{2(\Theta_i S_i)^2}{n_i(k)}, \tag{34}$$

For use in future derivations, we then take the square root of both sides of the equation:

$$\sqrt{\mathbb{E}[\|\nu_i(k)\|^2]} \leq \Theta_i S_i \sqrt{2} \sqrt{\left(1 - \frac{n_i(k)}{N_i}\right) \frac{1}{n_i(k)}}. \tag{35}$$

Additionally, by the concavity of a square root function and Jensen's inequality, which states that $\mathbb{E}[f(X)] \leq f(\mathbb{E}[X])$ for some differentiable, concave function $f$,

$$\mathbb{E}[\sqrt{\|\nu_i(k)\|^2}] \leq \sqrt{\mathbb{E}[\|\nu_i(k)\|^2]} \leq \Theta_i S_i \sqrt{2} \sqrt{\frac{N_i - n_i(k)}{N_i \times n_i(k)}}. \tag{36}$$

Thus the lemma is proven. $\square$

### B. Proof of Lemma 2

**Lemma 2.** *Taking the weighted average of Lemma 1 yields a form useful to manipulations necessary in later lemmas, i.e.*

$$\sigma(k) \triangleq \sum_i \rho_i \mathbb{E}[\|\nu_i(k)\|] = \sum_i \rho_i \Theta_i S_i \sqrt{2} \sqrt{\frac{N_i - n_i(k)}{N_i \times n_i(k)}} \tag{37}$$

### C. Proof of Lemma 3

**Lemma 3.**

$$\|\nabla F_i(\mathbf{w})\| \leq L, \forall i, \forall \mathbf{w} \tag{38}$$

*Proof.* From the convexity and $L$-Lipschitz conditions, for $\forall \mathbf{w}', \mathbf{w}$,

$$\langle \mathbf{w}' - \mathbf{w}, \nabla F_i(\mathbf{w}) \rangle \leq F_i(\mathbf{w}') - F_i(\mathbf{w}) \tag{39}$$

$$F_i(\mathbf{w}') - F_i(\mathbf{w}) \leq L\|\mathbf{w}' - \mathbf{w}\| \tag{40}$$

Letting $\mathbf{w}' = \mathbf{w} - \nabla F_i(\mathbf{w})$,

$$\|\nabla F_i(w) \leq L\| \tag{41}$$

### D. Proof of Lemma 4

**Lemma 4.** *With $\eta < \frac{2}{\beta}$, under Assumption 1,*

$$\begin{aligned} \epsilon_i(k) &\triangleq \|\mathbf{w}_i(k\tau - \Delta) - \mathbf{w}(k\tau - \Delta)\| \\ &\leq (1 - (1 - \alpha(k))^k)[2\eta L\left(\frac{\tau}{\alpha(k)} - \Delta\right) \\ &\quad + \eta\left(\frac{\tau}{\alpha(k)} - \Delta\right) \sum_j \rho_j \|\nu_j\| + \|\nu_i(k)\| \end{aligned} \tag{42}$$

*Proof.* Using the SGD approximation $g_i$ (which for brevity will have the $\mathcal{D}$ term not included), and letting $\ell = k\tau - r$ and $m = (k+1)\tau - r - \Delta$, we can say

$$
\begin{aligned}
&\mathbf{w}_i((k+1)\tau - \Delta) - \mathbf{w}((k+1)\tau - \Delta) \\
&= \mathbf{w}_i((k+1)\tau - \Delta) - \sum_j \rho_j \mathbf{w}_j((k+1)\tau - \Delta) \\
&= (1 - \alpha(k))[\mathbf{w}_i(k\tau - \Delta) - \mathbf{w}(k\tau - \Delta)] \\
&\quad - (1 - \alpha(k))\eta \sum_{r=1}^{\Delta} g_i(\mathbf{w}_i(\ell)) \\
&\quad + (1 - \alpha(k))\eta \sum_j \rho_j \sum_{r=1}^{\Delta} g_j(\mathbf{w}_j(\ell)) \\
&\quad - (1 - \alpha(k))\eta \sum_{r=1}^{\tau - \Delta} g_i(\mathbf{w}_i(m)) \\
&\quad + (1 - \alpha(k))\eta \sum_j \rho_j \sum_{r=1}^{\tau - \Delta} g_j(\mathbf{w}_j(m)) \\
&= (1 - \alpha(k))[\mathbf{w}_i(k\tau - \Delta) - \mathbf{w}(k\tau - \Delta)] \\
&\quad + (1 - \alpha(k))\eta \sum_{r=1}^{\Delta} \left[ \sum_j \rho_j g_j(\ell) - g_i(\mathbf{w}_i(\ell)) \right] \\
&\quad + \eta \sum_{r=1}^{\tau - \Delta} \left[ \sum_j \rho_j g_j(\mathbf{w}_j(m)) - g_i(\mathbf{w}_i(m)) \right]
\end{aligned}
\tag{43}
$$

Now expanding the SGD approximations into their gradients and noises,

$$
\begin{aligned}
&\mathbf{w}_i((k+1)\tau - \Delta) - \mathbf{w}((k+1)\tau - \Delta) \\
&= (1 - \alpha(k))[\mathbf{w}_i(k\tau - \Delta) - \mathbf{w}(k\tau - \Delta)] \\
&\quad + (1 - \alpha(k))\eta \sum_{r=1}^{\Delta} [\sum_{j \neq i} \rho_j \nabla F_j(\mathbf{w}_j(\ell)) \\
&\quad + \rho_i \nabla F_i(\mathbf{w}_i(\ell)) - \nabla F_i(\mathbf{w}_i(\ell)) \\
&\quad + \sum_j \rho_j \nu_j(k) - \nu_i(k)] \\
&\quad + \eta \sum_{r=1}^{\tau - \Delta} [\sum_{j \neq i} \rho_j \nabla F_j(\mathbf{w}_j(m)) \\
&\quad + \rho_i \nabla F_i(\mathbf{w}_i(m)) - \nabla F_i(\mathbf{w}_i(m)) \\
&\quad + \sum_j \rho_j \nu_j(k) - \nu_i(k)]
\end{aligned}
\tag{44}
$$

Using the triangle inequality and rearranging terms,

$$
\begin{aligned}
&\|\mathbf{w}_i((k+1)\tau - \Delta) - \mathbf{w}((k+1)\tau - \Delta)\| \\
&\leq (1 - \alpha(k))\|[\mathbf{w}_i(k\tau - \Delta) - \mathbf{w}(k\tau - \Delta)]\| \\
&\quad + (1 - \alpha(k))\eta(1 - \rho_i) \sum_{r=1}^{\Delta} \|\nabla F_i(\mathbf{w}_i(\ell))\| \\
&\quad + (1 - \alpha(k))\eta \sum_{r=1}^{\Delta} \sum_{j \neq i} \rho_j \|\nabla F_j(\mathbf{w}_j(\ell))\| \\
&\quad + (1 - \alpha(k))\eta \sum_{r=1}^{\Delta} \left[ \sum_j \rho_j \|\nu_j(k)\| + \|\nu_i(k)\| \right] \\
&\quad + \eta(1 - \rho_i) \sum_{r=1}^{\tau - \Delta} \|\nabla F_i(\mathbf{w}_i(m))\| \\
&\quad + \eta \sum_{r=1}^{\tau - \Delta} \sum_{j \neq i} \rho_j \|\nabla F_j(\mathbf{w}_j(m))\| \\
&\quad + \eta \sum_{r=1}^{\tau - \Delta} \left[ \sum_j \rho_j \|\nu_j(k)\| + \|\nu_i(k)\| \right]
\end{aligned}
\tag{45}
$$

Applying Lemma 3 and Assumption 1,

$$
\begin{aligned}
&\|\mathbf{w}_i((k+1)\tau - \Delta) - \mathbf{w}((k+1)\tau - \Delta)\| \\
&\leq (1 - \alpha(k))\|[\mathbf{w}_i(k\tau - \Delta) - \mathbf{w}(k\tau - \Delta)]\| \\
&\quad + 2\eta L(1 - \rho_i)(\tau - \alpha(k)\Delta) \\
&\quad + \eta(\tau - \alpha(k)\Delta) \left[ \sum_j \rho_j \|\nu_j(k)\| + \|\nu_i(k)\| \right]
\end{aligned}
\tag{46}
$$

Recursively unpacking the term until $t = -\Delta$, since $\mathbf{w}_i(-\Delta) = \mathbf{w}(-\Delta)$,

$$
\begin{aligned}
&\|\mathbf{w}_i(k\tau - \Delta) - \mathbf{w}(k\tau - \Delta)\| \\
&\leq (1 - \alpha(k))\|\mathbf{w}_i(-\Delta) - \mathbf{w}(-\Delta)\| \\
&\quad + (1 - (1 - \alpha(k))^k) \left[ 2\eta L \left( \frac{\tau}{\alpha(k)} - \Delta \right) \right] \\
&\quad + (1 - (1 - \alpha(k))^k) \left[ \eta \left( \frac{\tau}{\alpha(k)} - \Delta \right) \left[ \sum_j \rho_j \|\nu_j(k)\| + \|\nu_i(k)\| \right] \right] \\
&\triangleq \epsilon_i(k)
\end{aligned}
\tag{47}
$$

*E. Proof of Lemma 5*

**Lemma 5.** *Taking weighted average of Lemma 4 and applying Lemma 2,*

$$
\begin{aligned}
\epsilon(k) &\triangleq \mathbb{E} \left[ \sum_i \rho_i \epsilon_i(k) \right] \\
&= (1 - (1 - \alpha(k))^k) \left[ 2\eta(L + \sigma(k)) \left( \frac{\tau}{\alpha(k)} - \Delta \right) \right]
\end{aligned}
\tag{48}
$$

*Proof.* First taking the weighted average of all $\epsilon_i(k)$ terms,

$$\sum_i \rho_i \epsilon_i(k)$$
$$= (1 - (1 - \alpha(k))^k)\{2\eta L(\tau/\alpha(k) - \Delta) \qquad (49)$$
$$+ \eta(\tau/\alpha(k) - \Delta)[\sum_j \rho_j \|\nu_j(k)\| + \sum_i \rho_i \|\nu_i(k)\|]\}$$

Now taking the expectation,

$$\mathbb{E}\left[\sum_i \rho_i \epsilon_i(k)\right]$$
$$= (1 - (1 - \alpha(k))^k)\{2\eta L(\tau/\alpha(k) - \Delta)$$
$$+ \eta(\tau/\alpha(k) - \Delta)[\sum_j \rho_j \mathbb{E}[\|\nu_j(k)\|]] \qquad (50)$$
$$+ \sum_i \rho_i \mathbb{E}[\|\nu_i(k)\|]]\}$$
$$= (1 - (1 - \alpha(k))^k)\{2\eta L(\tau/\alpha(k) - \Delta)$$
$$+ \eta(\tau/\alpha(k) - \Delta)(2\sigma(k))\}$$

Thus proving the lemma after algebraic manipulations.

### F. Proof of Lemma 6.

**Lemma 6.** *Under Assumption 1, we have*

$$\|[\mathbf{w}_1 - \eta\nabla F(\mathbf{w}_1)] - [\mathbf{w}_2 - \eta\nabla F(\mathbf{w}_2)]\| \leq (1+\eta\beta)\|\mathbf{w}_1 - \mathbf{w}_2\| \qquad (51)$$

*Proof.* From the convexit of $F$,

$$F(\mathbf{w}_2) \leq F(\mathbf{w}_1) + (\mathbf{w}_2 - \mathbf{w}_1)^T \nabla F(\mathbf{w}_1) \qquad (52)$$
$$F(\mathbf{w}_1) \leq F(\mathbf{w}_2) + (\mathbf{w}_1 - \mathbf{w}_2)^T \nabla F(\mathbf{w}_2). \qquad (53)$$

Now summing the inequalities,

$$(\mathbf{w}_2 - \mathbf{w}_1)^T(\nabla F(\mathbf{w}_2) - \nabla F(\mathbf{w}_1)) \geq 0. \qquad (54)$$

By using the $\beta$-smoothness outlined in Assumption 1,

$$\|[\mathbf{w}_1\eta\nabla F(\mathbf{w}_1] - [\mathbf{w}_2 - \eta\nabla F(\mathbf{w}_2)]\|^2 \qquad (55)$$
$$= \|\mathbf{w}_1 - \mathbf{w}_2\|^2 + \eta^2\|\nabla F(\mathbf{w}_1) - \nabla F(\mathbf{w}_2)\|^2 \qquad (56)$$
$$- 2[\mathbf{w}_2 - \mathbf{w}_1][\nabla F(\mathbf{w}_2) - \eta\nabla F(\mathbf{w}_1)] \qquad (57)$$
$$\leq (1 + (\eta\beta)^2)\|\mathbf{w}_1 - \mathbf{w}_2\|^2. \qquad (58)$$

The result of the lemma follow accordingly.

### G. Proof of Lemma 7

**Lemma 7.** *Using Assumption 1, with learning rate $\eta < \frac{2}{\beta}$, for $t \in (k\tau - \Delta, (k+1)\tau - \Delta), t \neq k\tau$,*

$$\mathbb{E}[\|\mathbf{w}_i(t) - \mathbf{c}_k(t)\|] \leq \mathbb{E}[(1+\eta\beta)\|\mathbf{w}_i(t-1)$$
$$- \mathbf{c}_k(t-1)\|]$$
$$+ \eta\delta_i \qquad (59)$$
$$+ \eta\Theta_i S_i\sqrt{2}\sqrt{\frac{N_i - n_i(k)}{N_i n_i(k)}}$$

*Proof.* For $t \in (k\tau - \Delta, (k+1)\tau - \Delta), t \neq k\tau$,

$$\mathbf{w}_i(t) - \mathbf{c}_k(t)$$
$$= (\mathbf{w}_i(t-1) - \eta g_i(\mathbf{w}_i(t-1); \xi_i(t-1)))$$
$$- (\mathbf{c}_k(t-1) - \eta\nabla F(\mathbf{c}_k(t-1)))$$
$$= \mathbf{w}_i(t-1) - \mathbf{c}_k(t-1) \qquad (60)$$
$$- \eta[\nabla F_i(\mathbf{w}_i(t-1)) - \nabla F_i(\mathbf{c}_k(t-1))]$$
$$- \eta[\nabla F(\mathbf{c}_k(t-1) - \nabla F_i(\mathbf{c}_k(t-1))]$$
$$- \eta\nu_i(k)$$

Taking the norm and applying the triangle inequality,

$$\|\mathbf{w}_i(t) - \mathbf{c}_k\|$$
$$\leq \eta\|\nabla F_i(\mathbf{w}_i(t-1)) - \nabla F_i(\mathbf{c}_k(t-1))\|$$
$$+ \eta\|\nabla F_i(\mathbf{c}_k(t-1) - \nabla F(\mathbf{c}_k(t-1))\| \qquad (61)$$
$$+ \eta\|\nu_i(k)\|$$

Using Lemma 6 and Assumption 2,

$$\|\mathbf{w}_i(t) - \mathbf{c}_k(t)\| \leq (1+\eta\beta)\|\mathbf{w}_i(t-1) - \mathbf{c}_k(t-1)\|$$
$$+ \eta\delta_i \qquad (62)$$
$$+ \eta\|\nu_i(k)\|$$

Lastly taking the expectation and applying Lemma 1,

$$\mathbb{E}[\|\mathbf{w}_i(t) - \mathbf{c}_k(t)\|]$$
$$\leq \mathbb{E}[(1+\eta\beta)\|\mathbf{w}_i(t-1) - \mathbf{c}_k(t-1)\|]$$
$$+ \eta\delta_i \qquad (63)$$
$$+ \eta S_i\sqrt{\frac{N_i - n_i(k)}{N_i n_i(k)}}$$

### H. Proof of Lemma 8

**Lemma 8.** *Under Assumption 1 with $\eta < \frac{2}{\beta}$,*

$$\mathbb{E}[\|\mathbf{w}(k\tau) - \mathbf{c}_k(k\tau)\|]$$
$$\leq \alpha(k)\Delta L\eta$$
$$+ (1 - \alpha(k))\left[((1+\eta\beta)^\Delta - 1)\epsilon(k) + h(\Delta, k) + \eta\Delta\sigma(k)\right] \qquad (64)$$

*Where $h(x, k) = \frac{\delta + \sigma(k)}{\beta}[(1+\eta\beta)^x - 1] - \eta(\delta + \sigma(k))x$*

*Proof.* By the definitions of $\mathbf{w}(t)$ and $\mathbf{c}_k(t)$, after some algebraic manipulations,

$$\mathbf{w}(k\tau) = \sum_i \rho_i \mathbf{w}_i(k\tau)$$
$$= \mathbf{w}(k\tau - \Delta)$$
$$- (1 - \alpha(k))\sum_{r=1}^{\Delta}\sum_i \rho_i g_i(\mathbf{w}_i(k\tau - r); \xi_i(k\tau - r))$$
$$= \sum_i \rho_i \mathbf{w}_i(k\tau) \qquad (65)$$
$$- (1 - \alpha(k))\eta\sum_{r=1}^{\Delta}\sum_i \rho_i \nabla F_i(\mathbf{w}_i(k\tau - r))$$
$$- (1 - \alpha(k))\eta\sum_{r=1}^{\Delta}\sum_i \rho_i \nu_i(k)$$

and

$$\mathbf{c}_k(k\tau) = \mathbf{c}_k(k\tau - \Delta) - \eta \sum_{r=1}^{\Delta} \sum_i \rho_i \nabla F_i(\mathbf{c}_k(k\tau - r)) \quad (66)$$

Now we take the difference between the two previously defined terms,

$$\mathbf{w}(k\tau) - \mathbf{c}_k(k\tau) =$$
$$\eta\alpha(k) \sum_{r=1}^{\Delta} \sum_i \rho_i \nabla F_i(\mathbf{c}_k(k\tau - r))$$
$$- (1 - \alpha(k))\eta \sum_{r=1}^{\Delta} \sum_i \rho_i [\nabla F_i(\mathbf{w}_i(k\tau - r)) - \nabla F_i(\mathbf{c}_k(k\tau - r))]$$
$$- (1 - \alpha(k))\eta \sum_{r=1}^{\Delta} \sum_i \rho_i \nu_i(k), \quad (67)$$

and by taking the norm and applying the triangle inequality, we obtain

$$\|\mathbf{w}(k\tau) - \mathbf{c}_k(k\tau)\| \le \eta\alpha(k) \sum_{r=1}^{\Delta} \sum_i \rho_i \|\nabla F_i(\mathbf{c}_k(k\tau - r))\|$$
$$+ (1 - \alpha(k))\eta \sum_{r=1}^{\Delta} \sum_i \rho_i \|\nabla F_i(\mathbf{w}_i(k\tau - r)) - \nabla F_i(\mathbf{c}_k(k\tau - r))\|$$
$$+ (1 - \alpha(k))\eta \sum_{r=1}^{\Delta} \sum_i \rho_i \|\nu_i(k)\|. \quad (68)$$

Recursively unpacking terms ending at $\mathbf{w}(k\tau - \Delta) = \mathbf{c}_k(k\tau - \Delta)$, taking the expectation, applying Assumption 1, and using Lemma 7,

$$\mathbb{E}[\|\mathbf{w}(k\tau) - \mathbf{c}_k(k\tau)\|] \le \eta\alpha(k)\Delta L$$
$$+ (1 - \alpha(k))\eta\beta[$$
$$\sum_{r=1}^{\Delta} (1 + \eta\beta)^{\Delta - r} \sum_i \rho_i \mathbb{E}[\|\mathbf{w}_i(k\tau - \Delta) - \mathbf{w}(k\tau - \Delta)\|]]$$
$$+ (1 - \alpha(k))\eta\beta \sum_{r=1}^{\Delta} \sum_{j=0}^{\Delta - r - 1} (1 + \eta\beta)^j \sum_i \rho_i \delta_i$$
$$+ (1 - \alpha(k))\eta\beta \sum_{r=1}^{\Delta} \sum_{j=0}^{\Delta - r - 1} (1 + \eta\beta)^j \sum_i \rho_i \mathbb{E}[\|\nu_i(k)\|]$$
$$+ (1 - \alpha(k))\eta \sum_{r=1}^{\Delta} \sum_i \rho_i \mathbb{E}[\|\nu_i(k)\|] \quad (69)$$

Lastly, we apply Lemmas 2 and 5 and use Assumption 2 to

conclude that

$$\mathbb{E}[\|\mathbf{w}(k\tau) - \mathbf{c}_k(k\tau)\|]$$
$$\le \eta\alpha(k)\Delta L$$
$$+ (1 - \alpha(k))\eta\beta\epsilon(k) \sum_{r=1}^{\Delta} (1 + \eta\beta)^{\Delta - r}$$
$$+ \frac{\delta + \sigma(k)}{\beta}(1 - \alpha(k))\eta\beta \sum_{r=1}^{\Delta} [(1 + \eta\beta)^{\Delta - r} - 1]$$
$$+ (1 - \alpha(k))\eta\Delta\sigma(k), \quad (70)$$

with algebraic simplifications leading to the result of the lemma described above.

*I. Proof of Proposition 1*

**Proposition 1.** *Under Assumption 1 with $\eta < \frac{2}{\beta}$,*

$$\mathbb{E}[\|\mathbf{w}((k+1)\tau - \Delta) - \mathbf{c}_k((k+1)\tau - \Delta)\|$$
$$\le (1 - \alpha(k))\epsilon(k)([1 + \eta\beta]^\tau - 1)$$
$$+ (1 - \alpha(k))h(\tau, k) + \alpha(k)h(\tau - \Delta, k) \quad (71)$$
$$+ \alpha(k)\eta\Delta L[1 + \eta\beta]^{\tau - \Delta}$$
$$+ \eta\sigma(k)[\tau - \alpha(k)\Delta] \triangleq \psi(\alpha(k), k)$$

*Proof.* Let $t \in (k\tau - \Delta, (k+1)\tau - \Delta]$. Using (7),

$$\mathbf{w}_i = \alpha_t(k)\mathbf{w}(k\tau - \Delta) + (1 - \alpha_t(k))[\mathbf{w}_i(t - 1)$$
$$- \eta g_i(\mathbf{w}_i(t - 1); \xi_i(t - 1))] \quad (72)$$
$$\mathbf{c}_k(t) = \mathbf{c}_k(t - 1) - \eta \nabla F(\mathbf{c}_k(t - 1)) \quad (73)$$

Since

$$\mathbf{c}_k(k\tau - 1) = \mathbf{w}(k\tau - \Delta) - \eta \sum_{r=0}^{\Delta - 2} \nabla F(\mathbf{c}_k(k\tau - \Delta + r)) \quad (74)$$

it follows that (by taking $\sum_i \rho_i \mathbf{w}_i$) and expanding $g_i$ into its gradient and noise,

$$\mathbf{w}(t) - \mathbf{c}_k(t)$$
$$= (1 - \alpha_t(k)[\mathbf{w}(t - 1) - \mathbf{c}_k(t - 1)]$$
$$- (1 - \alpha_t(k))\eta \sum_i \rho_i [\nabla F_i(\mathbf{w}_i(t - 1)) - \nabla F_i(\mathbf{c}_k(t - 1))]$$
$$- (1 - \alpha_t(k))\eta \sum_i \rho_i \nu_i(k)$$
$$+ \eta\alpha_t(k) \sum_{r=0}^{\Delta - 1} \nabla F(\mathbf{c}_k(k\tau - \Delta + r)) \quad (75)$$

Applying the triangle inequality to the norm and applying Assumption 1 and Lemma 51,

$$\|\mathbf{w}(t) - \mathbf{c}_k(t)\|$$
$$\le (1 - \alpha_t(k))\|\mathbf{w}(t - 1) - \mathbf{c}_k(t - 1)\|$$
$$(1 - \alpha_t(k))\eta\beta \sum_i \rho_i \|\mathbf{w}_i(t - 1) - \mathbf{c}_k(t - 1)\|$$
$$+ \alpha_t(k)\eta\Delta L \quad (76)$$
$$+ (1 - \alpha_t(k))\eta \sum_i \rho_i \|\nu_i(k)\|$$

For $t \in [k\tau - \Delta, k\tau - 1]$, where $\alpha_t(k) = 0$, and using $\mathbf{c}_k(k\tau - \Delta) = \mathbf{w}(k\tau - \Delta)$

$$\|\mathbf{w}(t) - \mathbf{c}_k(t)\| \leq \eta\beta \sum_{\ell=k\tau-\Delta}^{t-1} \sum_i \rho_i \|\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)\|$$
$$+ \eta \sum_{\ell=k\tau-\Delta}^{t-1} \sum_i \rho_i \|\nu_i(k)\| \tag{77}$$

And for $t \in [k\tau, (k+t)\tau - \Delta]$, with $\alpha_{k\tau}(k) = \alpha(k)$, $\alpha_t(k) = 0, \forall t > k\tau$

$$\|\mathbf{w}(t) - \mathbf{c}_k(t)\| \leq (1 - \alpha(k))\eta\beta \sum_{\ell=k\tau-\Delta}^{k\tau-1} \sum_i \rho_i \|\nu_i(k)\|$$
$$+ \eta\beta \sum_{\ell=k\tau}^{t-1} \sum_i \rho_i \|\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)\|$$
$$+ \alpha(k)\eta\Delta L$$
$$+ (1 - \alpha(k)) + \eta \sum_{\ell=k\tau-\Delta}^{k\tau-1} \sum_i \rho_i \|\nu_i(k)\|$$
$$+ \eta \sum_{\ell=k\tau}^{t-1} \sum_i \rho_i \|\nu_i(k)\| \tag{78}$$

Which that implies that for $t = (k+1)\tau - \Delta$, by taking the expectation and applying Lemma 2 and Assumption 1,

$$\mathbb{E}[\|\mathbf{w}((k+1)\tau - \Delta) - \mathbf{c}_k((k+1)\tau - \Delta)\|]$$
$$\leq (1 - \alpha(k))\eta\beta \sum_{\ell=k\tau-\Delta}^{k\tau-1} \sum_i \rho_i \mathbb{E}[\|\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)\|]$$
$$+ \eta\beta \sum_{\ell=k\tau}^{(k+1)\tau-\Delta-1} \sum_i \rho_i \mathbb{E}[\|\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)\|]$$
$$+ \alpha(k)\eta\Delta L \tag{79}$$
$$+ (1 - \alpha(k))\eta \sum_{\ell=k\tau-\Delta}^{k\tau-1} \sigma(k)$$
$$+ \eta \sum_{\ell=k\tau}^{(k+1)\tau-\Delta-1} \sigma(k)$$

With everything else solved for, the term $\sum_i \rho_i[\|\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)\|]$, will now be derived, beginning with

$$\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)$$
$$= (1 - \alpha_\ell(k))[\mathbf{w}_i(\ell-1) - \mathbf{c}_k(\ell-1)]$$
$$- \eta(1 - \alpha_\ell(k))[\nabla F_i(\mathbf{w}_i(\ell-1)) - \nabla F_i(\mathbf{c}_k(\ell-1))]$$
$$- (1 - \alpha_\ell(k))\eta[\nabla F_i(\mathbf{c}_k(\ell-1)) - \nabla F(\mathbf{c}_k(\ell-1))] \tag{80}$$
$$+ \alpha_\ell(k)\eta \sum_{r=0}^{\Delta-1} \nabla F(\mathbf{c}_k(k\tau - \Delta + r))$$
$$- (1 - \alpha_\ell(k))\eta\nu_i(k)$$

Applying $\sum_i \rho_i$ and taking the norm,

$$\sum_i \rho_i \|\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)\|$$
$$\leq (1 - \alpha_\ell(k))[(1 + \eta\beta) \sum_i \rho_i \|\mathbf{w}_i(\ell-1) - \mathbf{c}_k(\ell-1)\|]$$
$$+ (1 - \alpha_\ell(k))\eta\delta$$
$$+ \alpha_\ell(k)\eta\Delta L$$
$$+ (1 - \alpha_\ell(k))\eta \sum_i \rho_i \|\nu_i(k)\|$$

$$= (1 - \alpha_\ell(k))[(1 + \eta\beta) \sum_i \rho_i \|\mathbf{w}_i(\ell-1) - \mathbf{c}_k(\ell-1)\|]$$
$$+ (1 - \alpha_\ell(k))\eta(\delta + \sum_i \rho_i \|\nu_i(k)\|)$$
$$+ \alpha_\ell(k)\eta\Delta L \tag{81}$$

Following a similar approach to dividing the time interval into separate parts, we first begin with the period $\ell \in [k\tau - \Delta, k\tau - 1], \alpha_\ell(k) = 0$,

$$\sum_i \rho_i \|\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)\|$$
$$\leq (1 + \eta\beta) \sum_i \rho_i \|\mathbf{w}_i(\ell-1) - \mathbf{c}_k(\ell-1)\| \tag{82}$$
$$+ \eta(\delta + \sum_i \rho_i \|\nu_i(k)\|)$$

Recursively unpacking the first term and using the fact that $\mathbf{w}(k\tau - \Delta) = \mathbf{c}_k(k\tau - \Delta)$,

$$\sum_i \rho_i \|\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)\|$$
$$\leq [1 + \eta\beta]^{\ell-(k\tau-\Delta)} \sum_i \rho_i \|\mathbf{w}_i(k\tau - \Delta) - \mathbf{w}(k\tau - \Delta)\|$$
$$+ (\delta + \sum_i \rho_i \|\nu_i(k)\|) \frac{[1 + \eta\beta]^{\ell-(k\tau-\Delta)} - 1}{\beta} \tag{83}$$

Taking the expectation and using Lemmas 2 and 5,

$$\sum_i \rho_i \mathbb{E}[\|\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)\|$$
$$\leq \epsilon(k)[1 + \eta\beta]^{\ell-k\tau+\Delta} \tag{84}$$
$$+ (\delta + \sigma(k)) \frac{[1 + \eta\beta]^{\ell-k\tau+\Delta} - 1}{\beta}$$

Similarly for $\ell \in [k\tau, (k+1)\tau - \Delta]$,

$$\sum_i \rho_i \mathbb{E}[\|\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)\|]$$
$$\leq (1 - \alpha)[1 + \eta\beta]^{\ell-(k\tau-\Delta)} \epsilon_k$$
$$+ (1 - \alpha)(\delta + \sigma(k))[1 + \eta\beta]^{\ell-k\tau} \frac{[1 + \eta\beta]^\Delta - 1}{\beta} \tag{85}$$
$$+ (\delta + \sigma(k)) \frac{[1 + \eta\beta]^{\ell-k\tau} - 1}{\beta}$$
$$+ \alpha\eta\Delta L[1 + \eta\beta]^{\ell-k\tau}$$

Which leads to

$$\sum_{\ell=k\tau-\Delta}^{k\tau-1} \sum_i \rho_i \, \mathbb{E}[\|\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)\|] \qquad (86)$$
$$\leq \epsilon(k)\frac{[1+\eta\beta]^\Delta - 1}{\eta\beta} + \frac{h(\Delta,k)}{\eta\beta}$$

and

$$\sum_{\ell=k\tau}^{(k+1)\tau-\Delta-1} \sum_i \rho_i \, \mathbb{E}[\|\mathbf{w}_i(\ell) - \mathbf{c}_k(\ell)\|]$$
$$\leq (1-\alpha)[1+\eta\beta]^\Delta \epsilon(k)\frac{[1+\eta\beta]^{\tau-\Delta} - 1}{\eta\beta} \qquad (87)$$
$$+ (1-\alpha)\frac{h(\tau,k) - h(\Delta,k)}{\eta\beta} + \alpha\frac{h(\tau-\Delta,k)}{\eta\beta}$$
$$+ \alpha\Delta L\frac{[1+\eta\beta]^{\tau-\Delta} - 1}{\beta}$$

The result of the lemma is thus yielded by plugging in the above into (79):

$$\mathbb{E}[\|\mathbf{w}((k+1)\tau-\Delta) - \mathbf{c}_k((k+1)\tau-\Delta)\|]$$
$$\leq (1-\alpha)\epsilon(k)([1+\eta\beta]^\tau - 1)$$
$$+ (1-\alpha)h(\tau,k) + \alpha h(\tau-\Delta,k) \qquad (88)$$
$$+ \alpha\eta\Delta L[1+\eta\beta]^{\tau-\Delta}$$
$$+ \eta\sigma(k)[\tau - \alpha\Delta] \triangleq \psi(\alpha,k)$$

$\square$

### J. Proof of Proposition 2

**Proposition 2.** *Let*

$$\omega = \frac{1}{\max_{k\in\{0,\dots,K-1\}} \|\mathbf{c}_k(k\tau-\Delta) - \mathbf{w}^\star\|^2}. \qquad (89)$$

*Under Assumption 1, and if the following conditions are met,*

1) $\eta < \frac{2}{\beta}$
2) $T\eta\phi - \frac{L\Psi(\widehat{\alpha})}{\Xi^2} > 0$
3) $F(\mathbf{c}_k((k+1)\tau-\Delta)) - F(\mathbf{w}^\star) \geq \Xi, \forall k$
4) $F(\mathbf{w}((K+1)\tau-\Delta)) - F(\mathbf{w}^\star) \geq \Xi,$

*for some $\Xi > 0$, we can upper-bound the convergence of* `StoFedDelAv` *as*

$$F(\mathbf{w}((K+1)\tau-\Delta) - F(\mathbf{w}^\star) \leq \frac{1}{T\eta\phi - \frac{L\Psi(\alpha)}{\Xi^2}}, \qquad (90)$$

*where $\Psi(\widehat{\alpha}) \triangleq \sum_{k=1}^K \psi(\alpha(k),k)$.*

*Proof.* We consider the case $\omega < \infty$ since $\omega = \infty$ is trivially tied to $\mathbf{w}((K+1)\tau-\Delta) = \mathbf{c}((K+1)\tau-\Delta) = \mathbf{w}^\star \Rightarrow F(\mathbf{w}((K+1)\tau-\Delta) = F(\mathbf{w}^\star)$. Then for every $k$ and $t \in [k\tau-\Delta, (k+1)\tau-\Delta]$, we define the sub-optimality gap of the centralized GD scheme as

$$\Gamma_{[k]}(t) = F(\mathbf{c}_k(t)) - F(\mathbf{w}^\star), \qquad (91)$$

noting that $\Gamma_{[k]}(t) \geq 0, \forall k$. Since $\mathbf{w}((K+1)\tau - \delta)) = \mathbf{c}_{[K+1]}((K+1)\tau - \Delta)$, we wish to prove that

$$\Gamma_{[K+1]}((K+1)\tau - \Delta))^{-1} \geq T\eta\phi - \frac{L\Psi(\widehat{\alpha})}{\Xi^2}. \qquad (92)$$

From the results of [5]'s Lemma 6, we know that

$$\Gamma_{[k]}^{-1}(t+1) - \Gamma_{[k]}^{-1}(t) \geq \frac{\eta(1 - (\eta\beta)/2)}{\|\mathbf{c}_k(t) - \mathbf{w}^\star\|^2}$$
$$\geq \frac{\eta(1 - (\eta\beta)/2)}{\max_k \|\mathbf{c}_k(t) - \mathbf{w}^\star\|^2} = \eta\omega\left(1 - \frac{\eta\beta}{2}\right) = \eta\phi. \qquad (93)$$

We therefore conclude that

$$\Gamma_{[k]}^{-1}((k+1)\tau - \Delta) - \Gamma_{[k]}^{-1}(k\tau - \Delta) \qquad (94)$$
$$= \sum_{t=k\tau-\Delta}^{(k+1)\tau-\Delta-1} \left[\Gamma_{[k]}^{-1}(t+1) - \Gamma_{[k]}^{-1}(t)\right] \geq \tau\eta\phi. \qquad (95)$$

With this in mind, we can conclude the following:

$$\sum_{k=1}^K \left[\Gamma_{[k]}^{-1}((k+1)\tau - \Delta) - \Gamma_{[k]}^{-1}(k\tau - \Delta)\right] \qquad (96)$$
$$= \Gamma_{[K+1]}^{-1}((K+1)\tau - \Delta)) - \Gamma_{[1]}^{-1}(\tau - \Delta) \qquad (97)$$
$$- \sum_{k=1}^K \left[\Gamma_{[k+1]}^{-1}((k+1)\tau - \Delta) - \Gamma_{[k]}^{-1}((k+1)\tau - \Delta)\right]$$
$$\geq T\eta\phi$$

To prove (92), we need to show that

$$\sum_{k=1}^K \left[\Gamma_{[k]}^{-1}((k+1)\tau - \Delta) - \Gamma_{[k+1]}^{-1}((k+1)\tau - \Delta)\right] \leq \frac{L\Psi(\widehat{\alpha})}{\Xi^2}. \qquad (98)$$

Since $\Psi(\widehat{\alpha}) = \sum_{k=1}^K \psi(\alpha(k),k)$, (98) is implied by

$$\Gamma_{[k+1]}((k+1)\tau - \Delta) - \Gamma_{[k]}((k+1)\tau - \Delta) \qquad (99)$$
$$\leq \frac{L\psi(\alpha(k),k)}{\Xi^2}\Gamma_{[k]}((k+1)\tau - \Delta)\Gamma_{[k+1]}((k+1)\tau - \Delta). \qquad (100)$$

Conditions (3) and (4) from the proposition statement imply that

$$\Gamma_{[k]}((k+1)\tau - \Delta) \geq \Xi, \forall k, \qquad (101)$$
$$\Gamma_{[K+1]}((K+1)\tau - \Delta) \geq \Xi. \qquad (102)$$

Using (94), with $k < K - 1$,

$$\Gamma_{[k+1]}((k+1)\tau - \Delta) \geq \frac{\Gamma_{[k+1]}((k+2)\tau - \Delta)}{1 - \tau\eta\phi\Gamma_{[k+1]}((k+2)\tau - \Delta)} \qquad (103)$$
$$\geq \Gamma_{[k+1]}((k+2)\tau - \Delta) \geq \Xi. \qquad (104)$$

The above statements show that (99) can be proven by showing that

$$\Gamma_{[k+1]}((k+1)\tau - \Delta) - \Gamma_{[k]}((k+1)\tau - \Delta) \leq L\psi(\alpha(k), k). \tag{105}$$

This is in fact the case. By combining with Proposition 1, we obtain:

$$\Gamma_{[k+1]}((k+1)\tau - \Delta) - \Gamma_{[k]}((k+1)\tau - \Delta) \tag{106}$$
$$= F(\mathbf{w}((k+1)\tau - \Delta)) - F(\mathbf{c}_k((k+1)\tau - \Delta)) \tag{107}$$
$$\leq L\mathbb{E}[\|\mathbf{w}((k+1)\tau - \delta) - \mathbf{c}_k((k+1)\tau - \Delta)]\|]. \tag{108}$$

The result of Proposition 2 directly follows. $\qquad\square$

*K. Proof of Theorem 1*

**Theorem 1.***With $\eta < \frac{2}{\beta}$ and under Assumption 1.*

$$F(\mathbf{w}^K) - F(\mathbf{w}^\star)$$
$$\leq \frac{1}{2\eta\phi T} + \sqrt{\frac{1}{(2\eta\phi T)^2} + \frac{L\Psi(\alpha)}{\eta\phi T}} + L\Psi(\widehat{\alpha}) \tag{109}$$

*where $\Psi(\widehat{\alpha}) = \sum_{k=1}^{K}\psi(\alpha, k)$.*

*Proof.* To prove Theorem 1, we first begin by defining an auxiliary variable $\Xi^* > 0$ given $\eta \leq \frac{1}{\beta}$ such that $T\eta\phi - \frac{L*\Psi(\widehat{\alpha})}{\Xi^{*2}} > 0$ and $\Xi^* = \frac{1}{T\eta\phi - \frac{L*\Psi\widehat{\alpha}}{\Xi^{*2}}}$. Solving these equations for $\Xi^*$ yields

$$\Xi^* = \frac{1}{2\eta\phi T} + \sqrt{\left(\frac{1}{2\eta\phi T}\right)^2 + \frac{L\Psi(\widehat{\alpha})}{\eta\phi T}} \tag{110}$$

Letting $\Xi > \Xi^*$, and assuming that the conditions of Proposition 2 are satisfied, it follows that

$$F(\mathbf{w}((K+1)\tau - \Delta)) - F(\mathbf{w}^\star) < \frac{1}{T\eta\phi - \frac{L\Psi(\widehat{\alpha})}{\Xi^2}}$$
$$\leq \frac{1}{T\eta\phi - \frac{L\Psi(\widehat{\alpha})}{\Xi^{*2}}} \tag{111}$$
$$\Rightarrow \Xi^* < \Xi$$

This presents a contradiction with the fourth condition of Prop. 2, meaning that at least one of those conditions cannot be satisfied given $\Xi > \Xi^*$. The first two conditions are readily satisfied and

$$\Xi > \Xi^* = \frac{1}{T\eta\phi - \frac{L\Psi(\widehat{\alpha})}{\Xi^{*2}}} \tag{112}$$

With either the third or fourth conditions not met, we therefore conclude that

$$\min\left\{F(\mathbf{w}((K+1)\tau - \Delta)), \min F(\mathbf{c}_k((k+1)\tau - \Delta))\right\}$$
$$- F(\mathbf{w}^\star) \leq \Xi^\star \tag{113}$$

Therefore, using Prop. 1 and noting that $\psi(\alpha(k), k)$ is increasing as a function of $k$,

$$F(\mathbf{w}((k+1)\tau - \Delta)) \leq F(\mathbf{c}_k((k+1)\tau - \Delta)) \tag{114}$$
$$+ \left|F(\mathbf{w}((k+1)\tau - \Delta)) - F(\mathbf{c}_k((k+1)\tau - \Delta))\right|$$

Taking the norm and expectation,

$$\leq F(\mathbf{c}_k((k+1)\tau - \Delta)) \tag{115}$$
$$+ L\mathbb{E}\left[\|\mathbf{w}((k+1)\tau - \Delta) - \mathbf{c}_k((k+1)\tau - \Delta)\|\right]$$
$$\leq F(\mathbf{c}_k((k+1)\tau - \Delta)) + L\psi(\alpha(k), k) \tag{116}$$
$$\leq F(\mathbf{c}_k((k+1)\tau - \Delta)) + \Psi(\widehat{\alpha}) \tag{117}$$

Implying

$$\min_k\{F(\mathbf{c}_k((k+1)\tau - \Delta))\} \tag{118}$$
$$\geq \min_k F(\mathbf{w}((k+1)\tau - \Delta)) - L\Psi(\widehat{\alpha})$$

Using the result of (113),

$$\min_{k \leq K} F(\mathbf{w}((k+1)\tau - \Delta)) - L\Psi(\widehat{\alpha}) - F(\mathbf{w}^*) \leq \Xi^* \tag{119}$$

with the theorem following as a direct consequence. $\qquad\square$