

Supplementary Material: Federated Learning under Heterogeneous and Correlated Client Availability

APPENDIX A PROOF OF THEOREM 1

Theorem 1 (Decomposing the total error). *Let $\kappa := L/\mu$. Under Assumptions 2–4, the optimization error of the target global objective $\epsilon = F(\mathbf{w}) - F^*$ can be bounded as follows:*

$$\epsilon \leq 2\kappa^2 \underbrace{(F_B(\mathbf{w}) - F_B^*)}_{:=\epsilon_{\text{opt}}} + \underbrace{F(\mathbf{w}_B^*) - F^*}_{:=\epsilon_{\text{bias}}}. \quad (10)$$

Moreover, let $\chi_{\alpha\|\mathbf{p}}^2 := \sum_{k=1}^N (\alpha_k - p_k)^2 / p_k$. Then:

$$\epsilon_{\text{bias}} \leq \kappa^2 \cdot \underbrace{\chi_{\alpha\|\mathbf{p}}^2}_{:=\epsilon_{\text{bias}}} \cdot \Gamma. \quad (11)$$

The proof of Theorem 1 employs well-established techniques from convex optimization. It is based on the proof presented in [33, Theorem 2].

Proof of Theorem 1. By leveraging the L -smoothness and μ -strong convexity properties of F , we obtain:

$$F(\mathbf{w}) - F^* \leq \frac{1}{2\mu} \|\nabla F(\mathbf{w})\|^2 \quad (21)$$

$$\leq \frac{L^2}{2\mu} \|\mathbf{w} - \mathbf{w}^*\|^2 \quad (22)$$

$$\leq \frac{L^2}{\mu} (\|\mathbf{w} - \mathbf{w}_B^*\|^2 + \|\mathbf{w}_B^* - \mathbf{w}^*\|^2) \quad (23)$$

$$\leq \frac{2L^2}{\mu^2} \left(\underbrace{F_B(\mathbf{w}) - F_B^*}_{:=\epsilon_{\text{opt}}} + \underbrace{F(\mathbf{w}_B^*) - F^*}_{:=\epsilon_{\text{bias}}} \right), \quad (24)$$

where the inequality in (21) follows from Assumption 4 and is commonly referred to as the *Polyak-Lojasiewicz inequality*; the inequality in (22) is derived using the fact that $\nabla F(\mathbf{w}^*) = 0$ (Assumption 2) and the definition of L -Lipschitz continuous gradient for F (Assumption 3); the inequality in (23) is based on $(a + b)^2 \leq 2(a^2 + b^2)$; lastly, the inequality in (24) follows from the μ -strong convexity of both F_B and F (Assumptions 4), and uses $\nabla F_B(\mathbf{w}_B^*) = 0$ and $\nabla F(\mathbf{w}^*) = 0$ (Assumption 2). The obtained results complete the first part of the proof, establishing the bound in (10).

Next, to prove the relation in (11), we proceed by bounding the term ϵ_{bias} as follows:

$$\epsilon_{\text{bias}} := (F(\mathbf{w}_B^*) - F^*) \leq \frac{1}{2\mu} \|\nabla F(\mathbf{w}_B^*)\|^2, \quad (25)$$

where the inequality in (25) directly follows from the Polyak-Lojasiewicz inequality (Assumption 4).

Furthermore, we bound the term $\|\nabla F(\mathbf{w}_B^*)\|$ as follows:

$$\|\nabla F(\mathbf{w}_B^*)\| = \left\| \sum_{k=1}^N (\alpha_k - p_k) \nabla F_k(\mathbf{w}_B^*) \right\| \quad (26)$$

$$\leq \sum_{k=1}^N |\alpha_k - p_k| \|\nabla F_k(\mathbf{w}_B^*)\| \quad (27)$$

$$\leq L \sum_{k=1}^N |\alpha_k - p_k| \|\mathbf{w}_B^* - \mathbf{w}_k^*\| \quad (28)$$

$$\leq L \sqrt{\frac{2}{\mu}} \sum_{k=1}^N |\alpha_k - p_k| \sqrt{(F_k(\mathbf{w}_B^*) - F_k^*)}, \quad (29)$$

where, in (26), we use $\nabla F_B(\mathbf{w}_B^*) = 0$ (Assumption 2) and apply the definitions of F and F_B given in (1) and (4), respectively. The bound in (27) follows from the triangle inequality. Next, the inequality in (28) uses $\nabla F_k(\mathbf{w}_k^*) = 0$ (Assumption 2) and

the L -smoothness of F_k (Assumption 3). Finally, the inequality in (29) leverages the μ -strong convexity of F_k (Assumption 4) and $\nabla F_k(\mathbf{w}_k^*) = 0$ (Assumption 2), and follows multiplying and dividing by $\sqrt{p_k}$.

By squaring both sides of Equation (29), we obtain:

$$\|\nabla F(\mathbf{w}_B^*)\|^2 \leq \frac{2L^2}{\mu} \left(\sum_{k=1}^N \frac{|\alpha_k - p_k|}{\sqrt{p_k}} \sqrt{p_k(F_k(\mathbf{w}_B^*) - F_k^*)} \right)^2 \quad (30)$$

$$\leq \frac{2L^2}{\mu} \left(\sum_{k=1}^N \frac{(\alpha_k - p_k)^2}{p_k} \right) \left(\sum_{k=1}^N p_k(F_k(\mathbf{w}_B^*) - F_k^*) \right) \quad (31)$$

$$\leq \frac{2L^2}{\mu} \cdot \chi_{\alpha\|p}^2 \cdot \Gamma, \quad (32)$$

where the inequality in (31) follows from the Cauchy-Schwarz inequality. Furthermore, the inequality in (32) holds because:

$$\sum_{k=1}^N p_k(F_k(\mathbf{w}_B^*) - F_k^*) = F_B^* - \sum_{k=1}^N p_k F_k^* \quad (33)$$

$$\leq F_B(\mathbf{w}^*) - \sum_{k=1}^N p_k F_k^* \quad (34)$$

$$= \sum_{k=1}^N p_k(F_k(\mathbf{w}^*) - F_k^*) \quad (35)$$

$$\leq \max_{k \in \mathcal{K}} \{F_k(\mathbf{w}^*) - F_k^*\} := \Gamma. \quad (36)$$

We remark that the inequality in (34) only holds if \mathbf{w}_B^* is the global minimizer of F_B , as guaranteed by Assumption 2. By replacing (32) into (25), we have:

$$\epsilon_{\text{bias}} \leq \frac{1}{2\mu} \|\nabla F(\mathbf{w}_B^*)\|^2 \leq \frac{L^2}{\mu^2} \cdot \chi_{\alpha\|p}^2 \cdot \Gamma, \quad (37)$$

which concludes the proof of Equation (11), and therefore, of Theorem 1. \square

APPENDIX B PROOF OF THEOREM 2

B1. Algorithm Overview and Supplementary Notation

Let $\mathbf{w}_{t,j}^k$ represent the model parameter maintained by the k -th client during the t -th global communication round and the j -th local step. The t -th global communication round can be described as follows: 1) The server broadcasts the model parameter $\mathbf{w}_{t,0}$ to the active clients, which adopt it as their local model, i.e., $\mathbf{w}_{t,0}^k = \mathbf{w}_{t,0}$ for $k \in \mathcal{A}_t$; 2) Each active client $k \in \mathcal{A}_t$ generates a sequence of local models $\{\mathbf{w}_{t,j}^k\}_{j=1}^E$ using the local-SGD update rule defined in (2); 3) The active clients send their model updates $\Delta_t^k := \mathbf{w}_{t,E}^k - \mathbf{w}_{t,0}$ back to the server; 4) The server aggregates the model updates using the aggregation rule specified in (3), resulting in the new global model parameter $\mathbf{w}_{t+1,0}$.

$$\begin{cases} \mathbf{w}_{t,j+1}^k = \mathbf{w}_{t,j}^k - \eta_t \nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) & \text{for } j = 0, \dots, E-1; \\ \mathbf{w}_{t+1,0} = \mathbf{Proj}_W(\mathbf{w}_{t,0} + \sum_{k \in \mathcal{A}_t} q_k (\mathbf{w}_{t,E}^k - \mathbf{w}_{t,0})) & \text{for } j = E. \end{cases} \quad (2)$$

$$\quad (3)$$

The projection operator in (3) ensures that the current iterate $\mathbf{w}_{t+1,0}$ in the optimization algorithm defined by (2) and (3) remains within the feasible region W .

Sources of randomness: In the system, we model two sources of randomness. The first arises from the availability of random clients, which follows a Markov process as stated in Assumption 1. The second source of randomness originates from the random sampling of batches for computing stochastic gradients. Remember that \mathcal{A}_t denotes the random set of clients available at the t -th communication round and that $\mathcal{B}_{t,j}^k$ denotes the random batch independently sampled from client- k 's local dataset at round t , local iteration j . For the analysis, we introduce the following additional notation:

- $\mathcal{A}_{i:j} := \{\mathcal{A}_i, \dots, \mathcal{A}_j\}$: the family of random sets of clients available from the i -th to the j -th communication rounds, $i < j$;
- $\mathcal{B}_t^k := \{\mathcal{B}_{t,j}^k\}_{j=0}^{E-1}$: the set of random batches sampled by the k -th client at the t -th communication round;
- $\mathcal{B}_t := \{\mathcal{B}_t^k\}_{k \in \mathcal{A}_t}$: the set of random batches sampled by the available clients (\mathcal{A}_t) in the t -th communication round;
- $\mathcal{B}_{t,i:j}^k := \{\mathcal{B}_{t,i}^k, \dots, \mathcal{B}_{t,j}^k\}$: the set of random batches sampled by the k -th client at the t -th communication round between the i -th and the j -th local iterations, $i < j$;
- $\mathcal{B}_{i:j} := \{\mathcal{B}_i, \dots, \mathcal{B}_j\}$: the set of random batches sampled by the available clients ($\mathcal{A}_{i:j}$) between the i -th and j -th communication rounds, $i < j$.

With this notation established, the randomness in the t -th communication round, which starts with the initial model $\mathbf{w}_{t,0}$ and yields the updated model $\mathbf{w}_{t+1,0}$, is fully determined by the sets \mathcal{A}_t and \mathcal{B}_t . This implies that the evolution of the algorithm, governed by the update rules in (2) and (3), from round 0 to round t can be completely described by the tuple:

$$\mathcal{H}_t := (\mathcal{A}_0, \dots, \mathcal{A}_{t-1}; \mathcal{B}_0, \dots, \mathcal{B}_{t-1}), \quad (38)$$

which represents the historical information up to the t -th communication round.

We introduce the following additional quantities for our analysis:

$$\mathbf{g}_t(\mathcal{A}_t, \mathcal{B}_t) := \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k), \quad (39)$$

and

$$\bar{\mathbf{g}}_t(\mathcal{A}_t, \mathcal{B}_t) := \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \nabla F_k(\mathbf{w}_{t,j}^k), \quad (40)$$

where $\mathbf{g}_t(\mathcal{A}_t, \mathcal{B}_t)$ denotes the global pseudo-gradient computed at communication round t , aggregated from the active clients in \mathcal{A}_t , and $\bar{\mathbf{g}}_t(\mathcal{A}_t, \mathcal{B}_t)$ denotes its expected value with respect to the choices of the random batches $\mathcal{B}_{t,j}^k$, for all $j = 0, \dots, E-1$ and $k \in \mathcal{A}_t$. With this notation established, the global update rule for the t -th communication round can be expressed as:

$$\mathbf{w}_{t+1,0} = \mathbf{Proj}_W(\mathbf{w}_{t,0} - \eta_t \mathbf{g}_t(\mathcal{A}_t, \mathcal{B}_t)). \quad (41)$$

B2. Supporting Lemmas

In this section, we introduce several lemmas that are instrumental in proving Theorem 2. Firstly, we prove Lemma 1, introduced in Section III-A. Its proof relies on the convexity and compactness of the hypothesis class W (Assumption 2), on the L -smoothness of the functions $\{F_k\}_{k \in \mathcal{K}}$ (Assumption 3), and on the bounded variance of the stochastic gradients (Assumption 5).

Lemma 1. *Under Assumptions 2, 3, and 5, there exist constants D , G , and $H > 0$, such that, for $\mathbf{w} \in W$ and $k \in \mathcal{K}$, we have:*

$$\|\nabla F_k(\mathbf{w})\| \leq D, \quad (6)$$

$$\mathbb{E} \|\nabla F_k(\mathbf{w}, \xi)\|^2 \leq G^2, \quad (7)$$

$$|F_k(\mathbf{w}) - F_k(\mathbf{w}_B^*)| \leq H. \quad (8)$$

Proof of Lemma 1. The boundedness of the hypothesis class W (Assumption 2) provides a bound on the sequence $(\mathbf{w}_{t,0})_{t \geq 0}$ generated by the scheme defined in Equations (2) and (3). Moreover, since \mathbf{w}_k^* minimizes $\nabla F_k(\mathbf{w})$, we have $\nabla F_k(\mathbf{w}_k^*) = 0$. Furthermore, the L -smoothness of $\{F_k\}_{k \in \mathcal{K}}$ (Assumption 3) leads to the following inequality:

$$\|\nabla F_k(\mathbf{w})\| = \|\nabla F_k(\mathbf{w}) - \nabla F_k(\mathbf{w}_k^*)\| \leq L \|\mathbf{w} - \mathbf{w}_k^*\| := D < +\infty. \quad (42)$$

The bound in (6) is directly derived from (42), while the bound in (8) follows from the continuity of $\{F_k\}_{k \in \mathcal{K}}$ over the compact set W (Assumption 2). Finally, the inequality in (7) requires a bound on the variance of the stochastic gradients (Assumption 5). In particular, it holds that:

$$\mathbb{E} \|\nabla F_k(\mathbf{w}, \xi)\|^2 \leq D^2 + \max_{k \in \mathcal{K}} \{\sigma_k^2\} := G^2. \quad (43)$$

□

The following lemma proves that the global pseudo-gradient $\mathbf{g}_t(\mathcal{A}_t, \mathcal{B}_t)$ is an unbiased estimator of $\bar{\mathbf{g}}_t(\mathcal{A}_t, \mathcal{B}_t)$. A similar result has been used in previous works, specifically in [33, Appendix C1]. Here, we provide a comprehensive proof for this result.

Lemma 2. *Let $\mathbf{g}_t(\mathcal{A}_t, \mathcal{B}_t)$ and $\bar{\mathbf{g}}_t(\mathcal{A}_t, \mathcal{B}_t)$ be defined as in (39) and (40), respectively. The following equality holds:*

$$\mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} [\mathbf{g}_t(\mathcal{A}_t, \mathcal{B}_t)] = \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} [\bar{\mathbf{g}}_t(\mathcal{A}_t, \mathcal{B}_t)]. \quad (44)$$

Proof of Lemma 2.

$$\mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} [\mathbf{g}_t(\mathcal{A}_t, \mathcal{B}_t)] = \quad (45)$$

$$= \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) \right] \quad (46)$$

$$= \sum_{k \in \mathcal{A}_t} q_k \mathbb{E}_{\mathcal{B}_t^k} \left[\sum_{j=0}^{E-1} \nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) \right] \quad (47)$$

$$= \sum_{k \in \mathcal{A}_t} q_k \left[\mathbb{E}_{\mathcal{B}_{t,0}^k} [\nabla F_k(\mathbf{w}_{t,0}, \mathcal{B}_{t,0}^k)] + \mathbb{E}_{\mathcal{B}_{t,0}^k, \mathcal{B}_{t,1}^k} [\nabla F_k(\mathbf{w}_{t,1}^k, \mathcal{B}_{t,1}^k)] + \cdots + \mathbb{E}_{\mathcal{B}_{t,0:E-1}^k} [\nabla F_k(\mathbf{w}_{t,E-1}^k, \mathcal{B}_{t,E-1}^k)] \right] \quad (48)$$

$$= \sum_{k \in \mathcal{A}_t} q_k \left[\nabla F_k(\mathbf{w}_{t,0}) + \mathbb{E}_{\mathcal{B}_{t,0}^k} \left[\mathbb{E}_{\mathcal{B}_{t,1}^k | \mathcal{B}_{t,0}^k} [\nabla F_k(\mathbf{w}_{t,1}^k, \mathcal{B}_{t,1}^k)] \right] + \cdots + \mathbb{E}_{\mathcal{B}_{t,0:E-2}^k} \left[\mathbb{E}_{\mathcal{B}_{t,E-1}^k | \mathcal{B}_{t,0:E-2}^k} [\nabla F_k(\mathbf{w}_{t,E-1}^k, \mathcal{B}_{t,E-1}^k)] \right] \right] \quad (49)$$

$$= \sum_{k \in \mathcal{A}_t} q_k \left[\nabla F_k(\mathbf{w}_{t,0}) + \mathbb{E}_{\mathcal{B}_{t,0}^k} [\nabla F_k(\mathbf{w}_{t,1}^k)] + \cdots + \mathbb{E}_{\mathcal{B}_{t,0:E-2}^k} [\nabla F_k(\mathbf{w}_{t,E-1}^k)] \right] \quad (50)$$

$$= \sum_{k \in \mathcal{A}_t} q_k \mathbb{E}_{\mathcal{B}_{t,0:E-2}^k} \left[\sum_{j=0}^{E-1} \nabla F_k(\mathbf{w}_{t,j}^k) \right] \quad (51)$$

$$= \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \nabla F_k(\mathbf{w}_{t,j}^k) \right] = \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} [\bar{\mathbf{g}}_t(\mathcal{A}_t, \mathcal{B}_t)], \quad (52)$$

where, in (47), we considered that both the evolution of the local models $\{\mathbf{w}_{t,j}^k\}_{j=0}^{E-1}$ and the choices of the random batches $\{\mathcal{B}_{t,j}^k\}_{j=0}^{E-1}$ are independent among different clients $k \in \mathcal{A}_t$ within the same communication round $t \in \mathcal{T}$. □

For the sake of simplicity, we will henceforth denote $\mathbf{g}_t(\mathcal{A}_t, \mathcal{B}_t)$ and $\bar{\mathbf{g}}_t(\mathcal{A}_t, \mathcal{B}_t)$ as \mathbf{g}_t and $\bar{\mathbf{g}}_t$, respectively. The following lemma decomposes the optimization error into multiple components, which we will bound separately in subsequent lemmas.

Lemma 3 (Decomposition of the error in a global communication round). *Let Assumption 2 hold. We have:*

$$\begin{aligned} \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \|\mathbf{w}_{t+1,0} - \mathbf{w}_B^*\|^2 &\leq \|\mathbf{w}_{t,0} - \mathbf{w}_B^*\|^2 - \underbrace{2\eta_t \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \langle \mathbf{w}_{t,0} - \mathbf{w}_B^*, \bar{\mathbf{g}}_t \rangle}_{\text{bounded in Lemma 4}} + \underbrace{\eta_t^2 \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \|\bar{\mathbf{g}}_t\|^2}_{\text{bounded in Lemma 5}} \\ &\quad + \underbrace{2\eta_t \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \langle \mathbf{w}_{t,0} - \mathbf{w}_B^* - \eta_t \bar{\mathbf{g}}_t, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle}_{\text{bounded in Lemma 6}} + \underbrace{\eta_t^2 \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \|\mathbf{g}_t - \bar{\mathbf{g}}_t\|^2}_{\text{bounded in Lemma 7}}. \end{aligned} \quad (53)$$

Proof of Lemma 3.

$$\|\mathbf{w}_{t+1,0} - \mathbf{w}_B^*\|^2 = \|\mathbf{Proj}_W(\mathbf{w}_{t,0} - \eta_t \mathbf{g}_t) - \mathbf{Proj}_W(\mathbf{w}_B^*)\|^2 \quad (54)$$

$$\leq \|\mathbf{w}_{t,0} - \eta_t \mathbf{g}_t - \mathbf{w}_B^* + \eta_t \bar{\mathbf{g}}_t - \eta_t \bar{\mathbf{g}}_t\|^2 \quad (55)$$

$$= \|\mathbf{w}_{t,0} - \mathbf{w}_B^* - \eta_t \bar{\mathbf{g}}_t\|^2 + 2\eta_t \langle \mathbf{w}_{t,0} - \mathbf{w}_B^* - \eta_t \bar{\mathbf{g}}_t, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle + \eta_t^2 \|\mathbf{g}_t - \bar{\mathbf{g}}_t\|^2 \quad (56)$$

$$= \|\mathbf{w}_{t,0} - \mathbf{w}_B^*\|^2 - 2\eta_t \langle \mathbf{w}_{t,0} - \mathbf{w}_B^*, \bar{\mathbf{g}}_t \rangle + \eta_t^2 \|\bar{\mathbf{g}}_t\|^2 + 2\eta_t \langle \mathbf{w}_{t,0} - \mathbf{w}_B^* - \eta_t \bar{\mathbf{g}}_t, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle + \eta_t^2 \|\mathbf{g}_t - \bar{\mathbf{g}}_t\|^2, \quad (57)$$

where, in (54), we used Assumption 2; whereas, the inequality in (55) is due to the contracting property of projection. We observe that (55) does not hold in general if $\mathbf{w}_B^* \notin W$. \square

In what follows, we present a series of lemmas to establish bounds for the error in (53).

Lemma 4. *Let Assumption 3 hold and the local functions $\{F_k\}_{k=1}^N$ be convex. We have:*

$$\begin{aligned} -2\eta_t \langle \mathbf{w}_{t,0} - \mathbf{w}_B^*, \bar{\mathbf{g}}_t \rangle &\leq -2\eta_t(1 - \eta_t L) \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_B^*)) \\ &\quad + \underbrace{\sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\mathbf{w}_{t,j}^k - \mathbf{w}_{t,0}\|^2}_{\text{bounded in Lemma 9}} + 2\eta_t^2 L E \underbrace{\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*)}_{\text{bounded in Lemma 10}}. \end{aligned} \quad (58)$$

Proof of Lemma 4. We decompose the term $-2\eta_t \langle \mathbf{w}_{t,0} - \mathbf{w}_B^*, \bar{\mathbf{g}}_t \rangle$, by adding and subtracting $\mathbf{w}_{t,j}^k$:

$$-2\eta_t \langle \mathbf{w}_{t,0} - \mathbf{w}_B^*, \bar{\mathbf{g}}_t \rangle = \underbrace{-2\eta_t \langle \mathbf{w}_{t,0} - \mathbf{w}_{t,j}^k, \bar{\mathbf{g}}_t \rangle}_{\text{developed in Eq. (60)}} - \underbrace{2\eta_t \langle \mathbf{w}_{t,j}^k - \mathbf{w}_B^*, \bar{\mathbf{g}}_t \rangle}_{\text{developed in Eq. (64)}}. \quad (59)$$

We bound the two terms separately. We bound the first term in (59) as:

$$-2\eta_t \langle \mathbf{w}_{t,0} - \mathbf{w}_{t,j}^k, \bar{\mathbf{g}}_t \rangle = -2\eta_t \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \langle \nabla F_k(\mathbf{w}_{t,j}^k), \mathbf{w}_{t,0} - \mathbf{w}_{t,j}^k \rangle \quad (60)$$

$$\leq \eta_t^2 \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\nabla F_k(\mathbf{w}_{t,j}^k)\|^2 + \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\mathbf{w}_{t,j}^k - \mathbf{w}_{t,0}\|^2 \quad (61)$$

$$\leq 2\eta_t^2 L \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (F_k(\mathbf{w}_{t,j}^k) - F_k^*) + \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\mathbf{w}_{t,j}^k - \mathbf{w}_{t,0}\|^2 \quad (62)$$

$$= 2\eta_t^2 L \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_B^*)) + \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\mathbf{w}_{t,j}^k - \mathbf{w}_{t,0}\|^2 + 2\eta_t^2 L E \sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*), \quad (63)$$

where, in (61), we used $|\langle \mathbf{a}, \mathbf{b} \rangle| \leq \frac{1}{2} \|\mathbf{a}\|^2 + \frac{1}{2} \|\mathbf{b}\|^2$; in (62), we applied the L -smoothness of $\{F_k(\mathbf{w})\}_{k \in \mathcal{K}}$ (Assumption 3); in (63), we added and subtracted $F_k(\mathbf{w}_B^*)$.

We bound the second term in (59) as:

$$-2\eta_t \langle \mathbf{w}_{t,j}^k - \mathbf{w}_B^*, \bar{\mathbf{g}}_t \rangle = -2\eta_t \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \langle \mathbf{w}_{t,j}^k - \mathbf{w}_B^*, \nabla F_k(\mathbf{w}_{t,j}^k) \rangle \quad (64)$$

$$\leq -2\eta_t \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_B^*)), \quad (65)$$

where, in (65), we use the convexity of $\{F_k(\mathbf{w})\}_{k \in \mathcal{K}}$.

By summing the bounds provided in (63) and (65), we conclude the proof. \square

Lemma 5 (Bound on the squared norm of a global gradient step). *Let Assumption 3 hold. We have:*

$$\eta_t^2 \|\bar{\mathbf{g}}_t\|^2 \leq 2\eta_t^2 L E Q \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_B^*)) + 2\eta_t^2 L E^2 Q \underbrace{\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*)}_{\text{bounded in Lemma 10}}. \quad (66)$$

Proof of Lemma 5.

$$\eta_t^2 \|\bar{\mathbf{g}}_t\|^2 = \eta_t^2 \left\| \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \nabla F_k(\mathbf{w}_{t,j}^k) \right\|^2 \quad (67)$$

$$\leq \eta_t^2 \sum_{k' \in \mathcal{A}_t} q_{k'} \sum_{k \in \mathcal{A}_t} q_k \left\| \sum_{j=0}^{E-1} \nabla F_k(\mathbf{w}_{t,j}^k) \right\|^2 \quad (68)$$

$$\leq \eta_t^2 Q E \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\nabla F_k(\mathbf{w}_{t,j}^k)\|^2 \quad (69)$$

$$\leq 2\eta_t^2 Q L E \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (F_k(\mathbf{w}_{t,j}^k) - F_k^*) \quad (70)$$

$$= 2\eta_t^2 L E Q \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_B^*)) + 2\eta_t^2 L E^2 Q \sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*), \quad (71)$$

where, in (68) and in (69), we applied the Jensen's inequality; in (69), we also observed that $\sum_{k \in \mathcal{A}_t} q_k \leq \sum_{k \in \mathcal{K}} q_k := Q$; in (70), we used the L -smoothness of $\{F_k(\mathbf{w})\}_{k \in \mathcal{K}}$ (Assumption 3); in (71), we added and subtracted $F_k(\mathbf{w}_B^*)$ to the sum. \square

Lemma 6. *Let Assumption 5 hold. We have:*

$$\begin{aligned} 2\eta_t \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} [\langle \mathbf{w}_{t,0} - \mathbf{w}_B^* - \eta_t \bar{\mathbf{g}}_t, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle] &\leq 2\eta_t^2 L E Q \sum_{k \in \mathcal{A}_t} q_k \sum_{j=1}^{E-1} \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} [F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_B^*)] \\ &\quad + \frac{1}{2} \eta_t^2 E (E-1) \sum_{k \in \mathcal{A}_t} q_k^2 \sigma_k^2 \\ &\quad + 2\eta_t^2 L E^2 Q \underbrace{\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*)}_{\text{bounded in Lemma 10}}. \end{aligned} \quad (72)$$

Proof of Lemma 6. We decompose the term $\langle \mathbf{w}_{t,0} - \mathbf{w}_B^* - \eta_t \bar{\mathbf{g}}_t, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle$ in two parts:

$$2\eta_t \langle \mathbf{w}_{t,0} - \mathbf{w}_B^* - \eta_t \bar{\mathbf{g}}_t, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle = 2\eta_t \langle \mathbf{w}_{t,0} - \mathbf{w}_B^*, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle - 2\eta_t^2 \langle \bar{\mathbf{g}}_t, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle. \quad (73)$$

From Lemma 2, we conclude that $\mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \langle \mathbf{w}_{t,0} - \mathbf{w}_B^*, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle = 0$.

We now focus on:

$$-2\eta_t^2 \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} [\langle \bar{\mathbf{g}}_t, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle] = \quad (74)$$

$$= -2\eta_t^2 \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} \sum_{k' \in \mathcal{A}_t} q_k q_{k'} \sum_{j=0}^{E-1} \sum_{j'=0}^{E-1} \langle \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_{k'}(\mathbf{w}_{t,j'}^{k'}) - \nabla F_{k'}(\mathbf{w}_{t,j'}^{k'}, \mathcal{B}_{t,j'}^{k'}) \rangle \right] \quad (75)$$

$$= -2\eta_t^2 \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=0}^{E-1} \sum_{j'=0}^{E-1} \langle \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}^k, \mathcal{B}_{t,j'}^k) \rangle \right]$$

$$- 2\eta_t^2 \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} \sum_{\substack{k' \in \mathcal{A}_t \\ k' \neq k}} q_k q_{k'} \sum_{j=0}^{E-1} \sum_{j'=0}^{E-1} \langle \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_{k'}(\mathbf{w}_{t,j'}^{k'}) - \nabla F_{k'}(\mathbf{w}_{t,j'}^{k'}, \mathcal{B}_{t,j'}^{k'}) \rangle \right] \quad (76)$$

$$= -2\eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{j=0}^{E-1} \sum_{j'=0}^{E-1} \langle \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}^k, \mathcal{B}_{t,j'}^k) \rangle \right]$$

$$-2\eta_t^2 \sum_{k \in \mathcal{A}_t} \sum_{\substack{k' \in \mathcal{A}_t \\ k' \neq k}} q_k q_{k'} \sum_{j=0}^{E-1} \sum_{j'=0}^{E-1} \left\langle \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} [\nabla F_k(\mathbf{w}_{t,j}^k)], \underbrace{\mathbb{E}_{\mathcal{B}_{t,0:j'-1}^{k'} | \mathcal{A}_t, \mathcal{H}_t} \left[\mathbb{E}_{\mathcal{B}_{t,j'}^{k'} | \mathcal{B}_{t,0:j'-1}^{k'}, \mathcal{A}_t, \mathcal{H}_t} [\nabla F_{k'}(\mathbf{w}_{t,j'}^{k'}) - \nabla F_{k'}(\mathbf{w}_{t,j'}, \mathcal{B}_{t,j'}^{k'})] \right]}_{=0} \right\rangle, \quad (77)$$

where, in (75), we replaced the definitions of g_t and \bar{g}_t given in (39) and in (40), respectively; in (76), we consider the cases $k = k'$ and $k \neq k'$ separately; (77) follows from the consideration that local models of different clients evolve independently and then all the terms with $k' \neq k$ equal zero because $\nabla F_k(\mathbf{w}, \mathcal{B})$ is an unbiased estimator of $\nabla F_k(\mathbf{w})$. It follows that:

$$-2\eta_t^2 \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} [\langle \bar{\mathbf{g}}_t, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle] = \quad (78)$$

$$= -2\eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{j=0}^{E-1} \sum_{j'=0}^{E-1} \langle \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}, \mathcal{B}_{t,j'}^k) \rangle \right] \quad (79)$$

$$= -2\eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{j=0}^{E-1} \sum_{\substack{j'=0 \\ j' < j}}^{E-1} \langle \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}, \mathcal{B}_{t,j'}^k) \rangle \right]$$

$$- 2\eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{j=0}^{E-1} \sum_{\substack{j'=0 \\ j' \geq j}}^{E-1} \langle \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}, \mathcal{B}_{t,j'}^k) \rangle \right] \quad (80)$$

$$= -2\eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=0}^{E-1} \sum_{\substack{j'=0 \\ j' < j}}^{E-1} \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} [\langle \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}, \mathcal{B}_{t,j'}^k) \rangle]$$

$$- 2\eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=0}^{E-1} \sum_{\substack{j'=0 \\ j' \geq j}}^{E-1} \mathbb{E}_{\mathcal{B}_{t,0:j'-1}^k | \mathcal{A}_t, \mathcal{H}_t} \left[\mathbb{E}_{\mathcal{B}_{t,j'}^k | \mathcal{B}_{t,0:j'-1}^k, \mathcal{A}_t, \mathcal{H}_t} [\langle \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}, \mathcal{B}_{t,j'}^k) \rangle] \right] \quad (81)$$

$$= -2\eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=0}^{E-1} \sum_{\substack{j'=0 \\ j' < j}}^{E-1} \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} [\langle \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}, \mathcal{B}_{t,j'}^k) \rangle]$$

$$- 2\eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=0}^{E-1} \sum_{\substack{j'=0 \\ j' \geq j}}^{E-1} \mathbb{E}_{\mathcal{B}_{t,0:j'-1}^k | \mathcal{A}_t, \mathcal{H}_t} \left[\underbrace{\mathbb{E}_{\mathcal{B}_{t,j'}^k | \mathcal{B}_{t,0:j'-1}^k, \mathcal{A}_t, \mathcal{H}_t} [\nabla F_k(\mathbf{w}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}, \mathcal{B}_{t,j'}^k)]}_{=0} \right], \quad (82)$$

where, in (80), we consider the cases $j' < j$ and $j' \geq j$ separately; then, in (81) and in (82), we use the law of total expectation.

Finally, we bound the remaining term in the right-hand side of (82) as follows:

$$-2\eta_t^2 \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} [\langle \bar{\mathbf{g}}_t, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle] = \quad (83)$$

$$= -2\eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=1}^{E-1} \sum_{j' < j} \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} \langle \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}, \mathcal{B}_{t,j'}^k) \rangle \quad (84)$$

$$= \eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=1}^{E-1} \sum_{j' < j} \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} \left[\|\nabla F_k(\mathbf{w}_{t,j}^k)\|^2 + \|\nabla F_k(\mathbf{w}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}, \mathcal{B}_{t,j'}^k)\|^2 \right] \quad (85)$$

$$= \eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=1}^{E-1} \sum_{j' < j} \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} \left[\|\nabla F_k(\mathbf{w}_{t,j}^k)\| \right] +$$

$$+ \eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=1}^{E-1} \sum_{j' < j} \underbrace{\mathbb{E}_{\mathcal{B}_{t,0:j'-1}^k | \mathcal{A}_t, \mathcal{H}_t} \left[\mathbb{E}_{\mathcal{B}_{t,j'}^k | \mathcal{B}_{t,0:j'-1}^k, \mathcal{A}_t, \mathcal{H}_t} \|\nabla F_k(\mathbf{w}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}, \mathcal{B}_{t,j'}^k)\|^2 \right]}_{\text{bounded with Assumption 5}} \quad (86)$$

$$\leq \eta_t^2 \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=1}^{E-1} \sum_{j' < j} \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} \|\nabla F_k(\mathbf{w}_{t,j}^k)\|^2 + \frac{1}{2} \eta_t^2 E(E-1) \sum_{k \in \mathcal{A}_t} q_k^2 \sigma_k^2 \quad (87)$$

$$\leq \eta_t^2 L(E-1) \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=1}^{E-1} \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} [(F_k(\mathbf{w}_{t,j}^k) - F_k^*)] + \frac{1}{2} \eta_t^2 E(E-1) \sum_{k \in \mathcal{A}_t} q_k^2 \sigma_k^2 \quad (88)$$

$$\begin{aligned} &= \eta_t^2 L(E-1) \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=1}^{E-1} \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} [(F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_B^*))] \\ &\quad + \eta_t^2 L E(E-1) \sum_{k \in \mathcal{A}_t} q_k^2 (F_k(\mathbf{w}_B^*) - F_k^*) + \frac{1}{2} \eta_t^2 E(E-1) \sum_{k \in \mathcal{A}_t} q_k^2 \sigma_k^2 \end{aligned} \quad (89)$$

$$\begin{aligned} &\leq \eta_t^2 L(E-1) Q \sum_{k \in \mathcal{A}_t} q_k \sum_{j=1}^{E-1} \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} [(F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_B^*))] \\ &\quad + \underbrace{\eta_t^2 L E(E-1) Q \sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*)}_{\text{bounded in Lemma 10}} + \frac{1}{2} \eta_t^2 E(E-1) \sum_{k \in \mathcal{A}_t} q_k^2 \sigma_k^2, \end{aligned} \quad (90)$$

where, in (85), we used $|\langle \mathbf{a}, \mathbf{b} \rangle| \leq \frac{1}{2} \|\mathbf{a}\|^2 + \frac{1}{2} \|\mathbf{b}\|^2$; in (87), we applied Assumption 5; in (88), we used the L -smoothness of $\{F_k(\mathbf{w})\}_{k \in \mathcal{K}}$; in (89), we added and subtracted $F_k(\mathbf{w}_B^*)$ from the sum; finally, in (90), we used $\sum_{k \in \mathcal{A}_t} q_k^2 f(k) \leq (\sum_{k \in \mathcal{A}_t} q_k)(\sum_{k \in \mathcal{A}_t} q_k f(k))$ and $\sum_{k \in \mathcal{A}_t} q_k \leq \sum_{k=1}^N q_k := Q$. Noting that $E-1 < 2E$ concludes the proof of Lemma 6. \square

Lemma 7 (Bound on the variance of the stochastic gradients). *Let Assumption 5 hold. Similarly to [23, Lemma 2], we have:*

$$\eta_t^2 \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \|\mathbf{g}_t - \bar{\mathbf{g}}_t\|^2 \leq \eta_t^2 E \sum_{k \in \mathcal{A}_t} q_k^2 \sigma_k^2. \quad (91)$$

Proof of Lemma 7.

$$\mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \|\mathbf{g}_t - \bar{\mathbf{g}}_t\|^2 = \quad (92)$$

$$= \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \left\| \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (\nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) - \nabla F_k(\mathbf{w}_{t,j}^k)) \right\|^2 \quad (93)$$

$$\begin{aligned} &= \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=0}^{E-1} \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} \|\nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) - \nabla F_k(\mathbf{w}_{t,j}^k)\|^2 \\ &\quad + \sum_{k \in \mathcal{A}_t} q_k^2 \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{j=0}^{E-1} \sum_{\substack{j'=0 \\ j' \neq j}}^{E-1} \langle \nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) - \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k, \mathcal{B}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}^k) \rangle \right] \\ &\quad + \sum_{k \in \mathcal{A}_t} \sum_{\substack{k' \in \mathcal{A}_t \\ k' \neq k}} q_k q_{k'} \sum_{j=0}^{E-1} \underbrace{\mathbb{E}_{\mathcal{B}_{t,0;j-1}^k | \mathcal{A}_t, \mathcal{H}_t} \left[\mathbb{E}_{\mathcal{B}_{t,j}^k | \mathcal{B}_{t,0;j-1}^k, \mathcal{A}_t, \mathcal{H}_t} [\nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) - \nabla F_k(\mathbf{w}_{t,j}^k)] \right]}_{=0}, \\ &\quad \quad \quad \mathbb{E}_{\mathcal{B}_{t,0;j-1}^{k'} | \mathcal{A}_t, \mathcal{H}_t} \left[\mathbb{E}_{\mathcal{B}_{t,j}^{k'} | \mathcal{B}_{t,0;j-1}^{k'}, \mathcal{A}_t, \mathcal{H}_t} [\nabla F_{k'}(\mathbf{w}_{t,j}^{k'}, \mathcal{B}_{t,j}^{k'}) - \nabla F_{k'}(\mathbf{w}_{t,j}^{k'})] \right] \rangle \\ &\quad + \sum_{k \in \mathcal{A}_t} \sum_{\substack{k' \in \mathcal{A}_t \\ k' \neq k}} q_k q_{k'} \sum_{j=0}^{E-1} \sum_{\substack{j'=0 \\ j' \neq j}}^{E-1} \underbrace{\mathbb{E}_{\mathcal{B}_{t,0;j-1}^k | \mathcal{A}_t, \mathcal{H}_t} \left[\mathbb{E}_{\mathcal{B}_{t,j}^k | \mathcal{B}_{t,0;j-1}^k, \mathcal{A}_t, \mathcal{H}_t} [\nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) - \nabla F_k(\mathbf{w}_{t,j}^k)] \right]}_{=0}, \\ &\quad \quad \quad \mathbb{E}_{\mathcal{B}_{t,0;j-1}^{k'} | \mathcal{A}_t, \mathcal{H}_t} \left[\mathbb{E}_{\mathcal{B}_{t,j}^{k'} | \mathcal{B}_{t,0;j-1}^{k'}, \mathcal{A}_t, \mathcal{H}_t} [\nabla F_{k'}(\mathbf{w}_{t,j'}^{k'}, \mathcal{B}_{t,j'}^{k'}) - \nabla F_{k'}(\mathbf{w}_{t,j'}^{k'})] \right] \rangle \end{aligned} \quad (94)$$

$$\begin{aligned}
&= \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=0}^{E-1} \underbrace{\mathbb{E}_{\mathcal{B}_{t,j}^k | \mathcal{A}_t, \mathcal{H}_t} \|\nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) - \nabla F_k(\mathbf{w}_{t,j}^k)\|^2}_{\text{bounded with Assumption 5}} \\
&+ \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=0}^{E-1} \sum_{\substack{j'=0 \\ j' < j}}^{E-1} \mathbb{E}_{\mathcal{B}_{t,0:j-1}^k | \mathcal{A}_t, \mathcal{H}_t} \left[\mathbb{E}_{\mathcal{B}_{t,j}^k | \mathcal{B}_{t,0:j-1}^k, \mathcal{A}_t, \mathcal{H}_t} [\langle \nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) - \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k, \mathcal{B}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}^k) \rangle] \right] \\
&+ \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=0}^{E-1} \sum_{\substack{j'=0 \\ j' > j}}^{E-1} \mathbb{E}_{\mathcal{B}_{t,0:j'-1}^k | \mathcal{A}_t, \mathcal{H}_t} \left[\mathbb{E}_{\mathcal{B}_{t,j'}^k | \mathcal{B}_{t,0:j'-1}^k, \mathcal{A}_t, \mathcal{H}_t} [\langle \nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) - \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k, \mathcal{B}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}^k) \rangle] \right] \\
\end{aligned} \tag{95}$$

$$\begin{aligned}
&= \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=0}^{E-1} \underbrace{\mathbb{E}_{\mathcal{B}_{t,j}^k | \mathcal{A}_t, \mathcal{H}_t} \|\nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) - \nabla F_k(\mathbf{w}_{t,j}^k)\|^2}_{\text{bounded with Assumption 5}} \\
&+ \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=0}^{E-1} \sum_{\substack{j'=0 \\ j' < j}}^{E-1} \mathbb{E}_{\mathcal{B}_{t,0:j-1}^k | \mathcal{A}_t, \mathcal{H}_t} \left[\underbrace{\mathbb{E}_{\mathcal{B}_{t,j}^k | \mathcal{B}_{t,0:j-1}^k, \mathcal{A}_t, \mathcal{H}_t} [\langle \nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) - \nabla F_k(\mathbf{w}_{t,j}^k), \nabla F_k(\mathbf{w}_{t,j'}^k, \mathcal{B}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}^k) \rangle]}_{=0} \right] \\
&+ \sum_{k \in \mathcal{A}_t} q_k^2 \sum_{j=0}^{E-1} \sum_{\substack{j'=0 \\ j' > j}}^{E-1} \mathbb{E}_{\mathcal{B}_{t,0:j'-1}^k | \mathcal{A}_t, \mathcal{H}_t} \left[\underbrace{\langle \nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) - \nabla F_k(\mathbf{w}_{t,j}^k), \mathbb{E}_{\mathcal{B}_{t,j'}^k | \mathcal{B}_{t,0:j'-1}^k, \mathcal{A}_t, \mathcal{H}_t} [\nabla F_k(\mathbf{w}_{t,j'}^k, \mathcal{B}_{t,j'}^k) - \nabla F_k(\mathbf{w}_{t,j'}^k)] \rangle}_{=0} \right] \\
\end{aligned} \tag{96}$$

$$\leq E \sum_{k \in \mathcal{A}_t} q_k^2 \sigma_k^2, \tag{97}$$

where, in (94), (95), and (96), we used the law of total expectation; in (97), we applied Assumption 5. Multiplying both sides of (97) by η_t^2 completes the proof of Lemma 7. \square

Lemma 8. *Let Assumption 3 hold and let the local functions $\{F_k\}_{k=1}^N$ be convex. Define $\gamma_t := 2\eta_t(1 - \eta_t L(1 + 2EQ))$. For a diminishing step-size $0 < \eta_t \leq \frac{1}{2L(1+2EQ)}$, satisfying $\gamma_t > 0$, we have:*

$$\begin{aligned}
-\gamma_t \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_B^*)) &\leq -\frac{1}{2}\eta_t E \sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) \\
&+ \underbrace{\sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\mathbf{w}_{t,j}^k - \mathbf{w}_{t,0}\|^2}_{\text{bounded in Lemma 9}} + 2\eta_t^2 L E \underbrace{\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*)}_{\text{bounded in Lemma 10}}, \tag{98}
\end{aligned}$$

Proof of Lemma 8. In the following, we require $\gamma_t > 0$.

$$-\gamma_t \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_B^*)) \tag{99}$$

$$= -\gamma_t \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_{t,0})) - \gamma_t \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) \tag{100}$$

$$\leq -\gamma_t \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \langle \nabla F_k(\mathbf{w}_{t,0}), \mathbf{w}_{t,j}^k - \mathbf{w}_{t,0} \rangle - \gamma_t E \sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) \tag{101}$$

$$\leq \gamma_t \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \frac{1}{2} \left[\eta_t \|\nabla F_k(\mathbf{w}_{t,0})\|^2 + \frac{1}{\eta_t} \|\mathbf{w}_{t,j}^k - \mathbf{w}_{t,0}\|^2 \right] - \gamma_t E \sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) \tag{102}$$

$$\leq \gamma_t \eta_t L E \sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k^*) + \frac{\gamma_t}{2\eta_t} \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\mathbf{w}_{t,j}^k - \mathbf{w}_{t,0}\|^2 - \gamma_t E \sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) \quad (103)$$

$$\leq -\gamma_t E (1 - \eta_t L) \sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) + \frac{\gamma_t}{2\eta_t} \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\mathbf{w}_{t,j}^k - \mathbf{w}_{t,0}\|^2 + \gamma_t \eta_t L E \sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*) \quad (104)$$

where, in (100), we added and subtracted $F_k(\mathbf{w}_{t,0})$ to the sum; in (101), we used the convexity of $\{F_k(\mathbf{w})\}_{k \in \mathcal{K}}$; note that (101) also requires $\gamma_t > 0$; in (102), we used the inequality $|\langle \mathbf{a}, \mathbf{b} \rangle| \leq \frac{1}{2} \|\mathbf{a}\|^2 + \frac{1}{2} \|\mathbf{b}\|^2$; in (103), we applied the L -smoothness of $\{F_k(\mathbf{w})\}_{k \in \mathcal{K}}$ (Assumption 3); finally, in (104), we added and subtracted $F_k(\mathbf{w}_B^*)$ to the sum.

In particular, for $\gamma_t := 2\eta_t(1 - \eta_t L(1 + 2EQ)) > 0$, since $0 < \eta_t \leq \frac{1}{2L(1+2EQ)}$, we further obtain:

$$\begin{aligned} & -\gamma_t \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} (F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_B^*)) \\ & \leq -\frac{1}{2} \eta_t E \sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) + \underbrace{\sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\mathbf{w}_{t,j}^k - \mathbf{w}_{t,0}\|^2}_{\text{bounded in Lemma 9}} + 2\eta_t^2 L E \underbrace{\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*)}_{\text{bounded in Lemma 10}}, \end{aligned} \quad (105)$$

where, in (105), we used $0 < \eta_t \leq \frac{1}{2L(1+2EQ)}$, which gives $-\gamma_t E(1 - \eta_t L) = -2\eta_t E(1 - \eta_t L(1 + 2EQ))(1 - \eta_t L) \leq -\frac{1}{2} \eta_t E$. Moreover, since $\gamma_t \leq 2\eta_t$, we also used $\gamma_t \eta_t \leq 2\eta_t^2$, and $\frac{\gamma_t}{2\eta_t} \leq 1$. \square

Lemma 9 (Bound on the divergence of local models). *Let Assumption 2, 3, and 5 hold, the local functions $\{F_k\}_{k=1}^N$ be convex and G be defined as in Lemma 1, Equation (7). Similarly to [23, Lemma 3], we obtain the following inequality:*

$$\mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\mathbf{w}_{t,j}^k - \mathbf{w}_{t,0}\|^2 \right] \leq \frac{1}{2} \eta_t^2 E^3 G^2 \left(\sum_{k \in \mathcal{A}_t} q_k \right). \quad (106)$$

Proof of Lemma 9.

$$\mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\mathbf{w}_{t,j}^k - \mathbf{w}_{t,0}\|^2 \right] = \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k \sum_{j=1}^{E-1} \eta_t^2 \left\| \sum_{j'=0}^{j-1} \nabla F_k(\mathbf{w}_{t,j'}^k, \mathcal{B}_{t,j'}^k) \right\|^2 \right] \quad (107)$$

$$\leq \eta_t^2 \sum_{k \in \mathcal{A}_t} q_k \sum_{j=1}^{E-1} j \sum_{j'=0}^{j-1} \mathbb{E}_{\mathcal{B}_t^k | \mathcal{A}_t, \mathcal{H}_t} \left[\|\nabla F_k(\mathbf{w}_{t,j'}^k, \mathcal{B}_{t,j'}^k)\|^2 \right] \quad (108)$$

$$\leq \eta_t^2 G^2 \left(\sum_{j=1}^{E-1} j^2 \right) \left(\sum_{k \in \mathcal{A}_t} q_k \right) \quad (109)$$

$$= \frac{1}{6} \eta_t^2 E(E-1)(2E-1) G^2 \left(\sum_{k \in \mathcal{A}_t} q_k \right), \quad (110)$$

where, in (108), we used the triangle and the Jensen's inequalities; in (109), we applied the bound in Lemma 1, Equation (7); finally, in (110), we developed the sum of sequence of squares $\sum_{j=1}^{E-1} j^2 = \frac{1}{6} E(E-1)(2E-1) \leq \frac{1}{2} E^3$ since $E \geq 1$. \square

Lemma 10 (Bound on the dissimilarity of local functions). *Let Assumption 1 hold and $(\mathcal{A}_t)_{t \geq 0}$ defined therein. We have:*

$$\mathbb{E} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*) \right] \leq \left(\sum_{k=1}^N \pi_k q_k \right) \Gamma, \quad (111)$$

where Γ is defined in (9).

Proof of Lemma 10.

$$\mathbb{E} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*) \right] = \sum_{k=1}^N \pi_k q_k (F_k(\mathbf{w}_B^*) - F_k^*) \quad (112)$$

$$= \left(\sum_{k'=1}^N \pi_{k'} q_{k'} \right) \sum_{k=1}^N p_k (F_k(\mathbf{w}_B^*) - F_k^*) \quad (113)$$

$$\leq \left(\sum_{k'=1}^N \pi_{k'} q_{k'} \right) \sum_{k=1}^N p_k (F_k(\mathbf{w}^*) - F_k^*) \quad (114)$$

$$\leq \left(\sum_{k'=1}^N \pi_{k'} q_{k'} \right) \underbrace{\max_{k \in \mathcal{K}} \{ (F_k(\mathbf{w}^*) - F_k^*) \}}_{:= \Gamma} = \left(\sum_{k=1}^N \pi_k q_k \right) \Gamma, \quad (115)$$

where, in (112), we solved the total expectation, observing that $\mathbb{E} [\sum_{k \in \mathcal{A}_t} q_k f(k)] = \sum_{k=1}^N \pi_k q_k f(k)$ (Assumption 1); in (113), we applied $p_k := \frac{\pi_k q_k}{\sum_{k'=1}^N \pi_{k'} q_{k'}}$; in (114), we used $F_B(\mathbf{w}) := \sum_{k=1}^N p_k F_k(\mathbf{w})$ and we observed $F_B(\mathbf{w}_B^*) \leq F_B(\mathbf{w}^*)$; finally, in (115), we used $\sum_{k=1}^N p_k = 1$ and $\Gamma := \max_{k \in \mathcal{K}} \{ (F_k(\mathbf{w}^*) - F_k^*) \}$. \square

Lemma 11 (Convergence results under heterogeneous client availability). *Let Assumptions 1–3 and 5 hold and the functions $\{F_k\}_{k=1}^N$ be convex. For a diminishing step-size $0 < \eta_t \leq \frac{1}{2L(1+2EQ)}$ satisfying $\sum_{t=1}^{+\infty} \eta_t^2 < +\infty$, for any $t_0 \leq T$, we have:*

$$\begin{aligned} \sum_{t=t_0}^T \eta_t \mathbb{E} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) \right] &\leq \frac{2}{E} \text{diam}(W)^2 + (E+1) \left(\sum_{k=1}^N \pi_k q_k^2 \sigma_k^2 \right) \left(\sum_{t=1}^{+\infty} \eta_t^2 \right) \\ &\quad + 2E^2 G^2 \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=1}^{+\infty} \eta_t^2 \right) \\ &\quad + 4L(1+EQ) \Gamma \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=1}^{+\infty} \eta_t^2 \right) \\ &:= C_0 < +\infty. \end{aligned} \quad (116)$$

Proof of Lemma 11. We take expectation over $\mathcal{B}_t \mid \mathcal{A}_t, \mathcal{H}_t$ on Lemma 3:

$$\begin{aligned} \mathbb{E}_{\mathcal{B}_t \mid \mathcal{A}_t, \mathcal{H}_t} \|\mathbf{w}_{t+1,0} - \mathbf{w}_B^*\|^2 &\leq \underbrace{\|\mathbf{w}_{t,0} - \mathbf{w}_B^*\|^2 - 2\eta_t \mathbb{E}_{\mathcal{B}_t \mid \mathcal{A}_t, \mathcal{H}_t} \langle \mathbf{w}_{t,0} - \mathbf{w}_B^*, \bar{\mathbf{g}}_t \rangle}_{\text{bounded in Lemma 4}} + \underbrace{\eta_t^2 \mathbb{E}_{\mathcal{B}_t \mid \mathcal{A}_t, \mathcal{H}_t} \|\bar{\mathbf{g}}_t\|^2}_{\text{bounded in Lemma 5}} \\ &\quad + \underbrace{2\eta_t \mathbb{E}_{\mathcal{B}_t \mid \mathcal{A}_t, \mathcal{H}_t} \langle \mathbf{w}_{t,0} - \mathbf{w}_B^* - \eta_t \bar{\mathbf{g}}_t, \bar{\mathbf{g}}_t - \mathbf{g}_t \rangle}_{\text{bounded in Lemma 6}} + \underbrace{\eta_t^2 \mathbb{E}_{\mathcal{B}_t \mid \mathcal{A}_t, \mathcal{H}_t} \|\mathbf{g}_t - \bar{\mathbf{g}}_t\|^2}_{\text{bounded in Lemma 7}}. \end{aligned} \quad (117)$$

Replacing Lemmas 4–7 in (117), we obtain:

$$\begin{aligned} \mathbb{E}_{\mathcal{B}_t \mid \mathcal{A}_t, \mathcal{H}_t} \|\mathbf{w}_{t+1,0} - \mathbf{w}_B^*\|^2 &\leq \|\mathbf{w}_{t,0} - \mathbf{w}_B^*\|^2 + 2\eta_t^2 L E (1 + 2EQ) \mathbb{E}_{\mathcal{B}_t \mid \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*) \right] \\ &\quad - \underbrace{2\eta_t (1 - \eta_t L (1 + 2EQ)) \mathbb{E}_{\mathcal{B}_t \mid \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k \sum_{j=1}^{E-1} (F_k(\mathbf{w}_{t,j}^k) - F_k(\mathbf{w}_B^*)) \right]}_{\gamma_t} \\ &\quad + \frac{1}{2} \eta_t^2 E (E + 1) \sum_{k \in \mathcal{A}_t} q_k^2 \sigma_k^2 + \underbrace{\mathbb{E}_{\mathcal{B}_t \mid \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \|\mathbf{w}_{t,j}^k - \mathbf{w}_{t,0}\|^2 \right]}_{\text{bounded in Lemma 9}} \end{aligned} \quad (118)$$

We apply Lemmas 8 and 9 to (118) with $\gamma_t := 2\eta_t(1 - \eta_t L(1 + 2EQ))$. We observe that $\gamma_t > 0$ because:

$$0 \leq \eta_t \leq \frac{1}{2L(1 + 2EQ)}. \quad (119)$$

We obtain:

$$\begin{aligned} \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \|\mathbf{w}_{t+1,0} - \mathbf{w}_B^*\|^2 &\leq \|\mathbf{w}_{t,0} - \mathbf{w}_B^*\|^2 - \frac{1}{2}\eta_t E \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) \right] \\ &\quad + \frac{1}{2}\eta_t^2 E(E+1) \sum_{k \in \mathcal{A}_t} q_k^2 \sigma_k^2 + \eta_t^2 E^3 G^2 \sum_{k \in \mathcal{A}_t} q_k \\ &\quad + 4\eta_t^2 LE(1 + EQ) \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*) \right]. \end{aligned} \quad (120)$$

Computing the total expectation on (120), we have:

$$\begin{aligned} \mathbb{E}_{\mathcal{A}_t, \mathcal{B}_t, \mathcal{H}_t} \|\mathbf{w}_{t+1,0} - \mathbf{w}_B^*\|^2 &\leq \mathbb{E}_{\mathcal{H}_t} \|\mathbf{w}_{t,0} - \mathbf{w}_B^*\|^2 - \frac{1}{2}\eta_t E \mathbb{E}_{\mathcal{A}_t, \mathcal{B}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) \right] \\ &\quad + \frac{1}{2}\eta_t^2 E(E+1) \mathbb{E}_{\mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k^2 \sigma_k^2 \right] + \eta_t^2 E^3 G^2 \mathbb{E}_{\mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k \right] \\ &\quad + 4\eta_t^2 LE(1 + EQ) \underbrace{\mathbb{E}_{\mathcal{A}_t, \mathcal{H}_t} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_B^*) - F_k^*) \right]}_{\text{bounded in Lemma 10}} \end{aligned} \quad (121)$$

Applying Lemma 10 to (121) and considering $\mathbb{E} [\sum_{k \in \mathcal{A}_t} a_k] = \sum_{k=1}^N \pi_k a_k$ (Assumption 1), the following inequality holds:

$$\begin{aligned} \mathbb{E} \|\mathbf{w}_{t+1,0} - \mathbf{w}_B^*\|^2 &\leq \mathbb{E} \|\mathbf{w}_{t,0} - \mathbf{w}_B^*\|^2 - \frac{1}{2}\eta_t E \mathbb{E} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) \right] \\ &\quad + \frac{1}{2}\eta_t^2 E(E+1) \left(\sum_{k=1}^N \pi_k q_k^2 \sigma_k^2 \right) + \eta_t^2 E^3 G^2 \left(\sum_{k=1}^N \pi_k q_k \right) + 4\eta_t^2 LE(1 + EQ) \Gamma \left(\sum_{k=1}^N \pi_k q_k \right). \end{aligned} \quad (122)$$

Rearranging and summing over $t = t_0, \dots, T$, we obtain the following inequality:

$$\begin{aligned} \sum_{t=t_0}^T \eta_t \mathbb{E} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) \right] &\leq \frac{2}{E} \sum_{t=t_0}^T \mathbb{E} \left[\left(\|\mathbf{w}_{t,0} - \mathbf{w}_B^*\|^2 - \|\mathbf{w}_{t+1,0} - \mathbf{w}_B^*\|^2 \right) \right] \\ &\quad + (E+1) \left(\sum_{k=1}^N \pi_k q_k^2 \sigma_k^2 \right) \left(\sum_{t=t_0}^T \eta_t^2 \right) \\ &\quad + 2E^2 G^2 \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=t_0}^T \eta_t^2 \right) \\ &\quad + 4L(1 + EQ) \Gamma \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=t_0}^T \eta_t^2 \right). \end{aligned} \quad (123)$$

The first term in the right-hand side of (123) is a telescoping sum and we remove the negative term $-\mathbb{E} \|\mathbf{w}_{T+1,0} - \mathbf{w}_B^*\|^2$:

$$\begin{aligned} \sum_{t=t_0}^T \eta_t \mathbb{E} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_B^*)) \right] &\leq \frac{2}{E} \mathbb{E} \|\mathbf{w}_{t_0,0} - \mathbf{w}_B^*\|^2 + (E+1) \left(\sum_{k=1}^N \pi_k q_k^2 \sigma_k^2 \right) \left(\sum_{t=t_0}^T \eta_t^2 \right) \\ &\quad + 2E^2 G^2 \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=t_0}^T \eta_t^2 \right) \end{aligned}$$

$$+ 4L(1 + EQ)\Gamma \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=t_0}^T \eta_t^2 \right). \quad (124)$$

Finally, by noting that $\|\mathbf{w}_{t_0,0} - \mathbf{w}_B^*\| \leq \text{diam}(W)$ and $\sum_{t=t_0}^T \eta_t^2 \leq \sum_{t=1}^{+\infty} \eta_t^2 < +\infty$, we complete the proof of Lemma 11. \square

Lemma 12. *Let Assumptions 2 and 3 hold, and the local functions $\{F_k\}_{k=1}^N$ be convex. We have:*

$$|F_k(\mathbf{v}) - F_k(\mathbf{w})| \leq D \cdot \|\mathbf{v} - \mathbf{w}\|, \quad \forall \mathbf{v}, \mathbf{w} \in W \quad (125)$$

Proof of Lemma 12. In Lemma 1, under Assumptions 2 and 3, we have already proved that:

$$\|\nabla F_k(\mathbf{w})\| \leq D. \quad (6)$$

Moreover, from the convexity of $\{F_k\}_{k \in \mathcal{K}}$, it follows that:

$$\langle \nabla F_k(\mathbf{v}), \mathbf{v} - \mathbf{w} \rangle \leq F_k(\mathbf{v}) - F_k(\mathbf{w}) \leq \langle \nabla F_k(\mathbf{w}), \mathbf{v} - \mathbf{w} \rangle. \quad (126)$$

The Cauchy–Schwarz inequality completes the proof of Lemma 12:

$$|F_k(\mathbf{v}) - F_k(\mathbf{w})| \leq \max\{\|\nabla F_k(\mathbf{v})\|, \|\nabla F_k(\mathbf{w})\|\} \cdot \|\mathbf{v} - \mathbf{w}\| \leq D \cdot \|\mathbf{v} - \mathbf{w}\|. \quad (127)$$

\square

Lemma 13. *Let Assumptions 2, 3, and 5 hold. We have:*

$$\mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \|\mathbf{w}_{t+1,0} - \mathbf{w}_{t,0}\| \leq \eta_t EG \left(\sum_{k \in \mathcal{A}_t} q_k \right). \quad (128)$$

Proof of Lemma 13. The proof is based on [15, Proposition 1.4].

$$\mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \|\mathbf{w}_{t+1,0} - \mathbf{w}_{t,0}\| = \mathbb{E}_{\mathcal{B}_t | \mathcal{A}_t, \mathcal{H}_t} \left\| -\eta_t \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k) \right\| \quad (129)$$

$$\leq \eta_t \sum_{k \in \mathcal{A}_t} q_k \sum_{j=0}^{E-1} \mathbb{E}_{\mathcal{B}_{t,0:j-1}^k | \mathcal{A}_t, \mathcal{H}_t} \left[\mathbb{E}_{\mathcal{B}_{t,j}^k | \mathcal{B}_{t,0:j-1}^k, \mathcal{A}_t, \mathcal{H}_t} [\nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k)] \right] \quad (130)$$

$$\leq \eta_t EG \left(\sum_{k \in \mathcal{A}_t} q_k \right), \quad (131)$$

where, in (130), we used the triangle inequality and the law of total expectation; in (131), we applied Lemma 1, Equation (7). \square

Similarly to [15, Theorem 1], we provide the following definition.

Definition 1. For communication round $t \geq 1$, denote the positive integer \mathcal{J}_t as follows:

$$\mathcal{J}_t := \min \left\{ \max \left\{ \left\lceil \frac{\ln(2C_P H t)}{\ln(1/\lambda(\mathbf{P}))} \right\rceil, T_P \right\}, t \right\}. \quad (132)$$

The parameter \mathcal{J}_t is crucial in our analysis: it represents the communication rounds needed to bound the stationary distribution convergence of the Markov process $(\mathcal{A}_t)_{t>0}$. It will play a key role in Lemmas 14–18 and in the proof of Theorem 2. We remark that, by definition: $T_P \leq \mathcal{J}_t \leq t$.

Our definition of \mathcal{J}_t corrects a typo in [15, (6.27)], which considered $\ln(t/(2C_P H))$ rather than $\ln(2C_P H t)$. In fact, we observe that [15, (6.28)] and consequently [15, (6.35)] do not hold when \mathcal{J}_t is defined as in [15, (6.27)].

Lemma 14 (Convergence results under heterogeneous and correlated client availability after \mathcal{J}_t communication rounds). *Let Assumptions 1–3, and 5 hold, the local functions $\{F_k\}_{k=1}^N$ be convex, and the parameter $\mathcal{J}_t \leq t$ be as in Definition 1. For a diminishing step-size $\{\eta_t\}_{t \geq 1}$ satisfying $\sum_{t=1}^{+\infty} \ln(t) \cdot \eta_t^2 < +\infty$, for any $t_0 \leq T$, we have:*

$$\sum_{t=t_0}^T \eta_t \mathbb{E} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_{t,0})) \right] \leq \frac{C_1}{\ln(1/\lambda(\mathbf{P}))} < +\infty, \quad (133)$$

where:

$$C_1 := EDGQ \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=1}^{+\infty} \ln(2C_P H t) \eta_{t-\mathcal{J}_t}^2 \right). \quad (134)$$

Proof of Lemma 14. This proof is based on [15, Equation (6.31)].

$$\sum_{t=t_0}^T \eta_t \mathbb{E} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_{t,0})) \right] \leq Q \sum_{t=t_0}^T \eta_t \mathbb{E} \left[\max_{k \in \mathcal{K}} \{F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_{t,0})\} \right] \quad (135)$$

$$\leq DQ \sum_{t=t_0}^T \eta_t \mathbb{E} \|\mathbf{w}_{t-\mathcal{J}_t,0} - \mathbf{w}_{t,0}\| \quad (136)$$

$$\leq DQ \sum_{t=t_0}^T \eta_t \sum_{d=t-\mathcal{J}_t}^{t-1} \mathbb{E}_{\mathcal{A}_d, \mathcal{H}_d} \left[\mathbb{E}_{\mathcal{B}_d | \mathcal{A}_d, \mathcal{H}_d} \|\mathbf{w}_{d,0} - \mathbf{w}_{d+1,0}\| \right] \quad (137)$$

$$\leq EDGQ \sum_{t=t_0}^T \sum_{d=t-\mathcal{J}_t}^{t-1} \eta_t \eta_d \mathbb{E} \left[\sum_{k \in \mathcal{A}_d} q_k \right] \quad (138)$$

$$\leq EDGQ \left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=t_0}^T \sum_{d=t-\mathcal{J}_t}^{t-1} \eta_t \eta_d \quad (139)$$

$$\leq \frac{EDGQ}{2} \left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=t_0}^T \sum_{d=t-\mathcal{J}_t}^{t-1} (\eta_t^2 + \eta_d^2) \quad (140)$$

$$\leq EDGQ \left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=t_0}^T \mathcal{J}_t \eta_{t-\mathcal{J}_t}^2, \quad (141)$$

where, in (135), we used $\sum_{k \in \mathcal{A}_t} q_k a_k \leq \sum_{k=1}^N q_k a_k \leq (\sum_{k=1}^N q_k) \cdot \max_{k \in \mathcal{K}} \{a_k\} = Q \cdot \max_{k \in \mathcal{K}} \{a_k\}$; in (136), we applied Lemma 12; in (137), we used the triangle inequality and the law of total expectation; in (138), we applied Lemma 13 and again the law of total expectation; in (139), we observed that $\mathbb{E} \left[\sum_{k \in \mathcal{A}_d} q_k \right] = \sum_{k=1}^N \pi_k q_k$ (Assumption 1); in (140), we used $2ab \leq a^2 + b^2$; finally, in (141), we applied $\eta_t < \eta_d \leq \eta_{t-\mathcal{J}_t}$ due to the diminishing learning rate.

We apply then the definition of \mathcal{J}_t in (132) and we observe that $\sum_{t=t_0}^T \ln(t) \eta_{t-\mathcal{J}_t}^2 \leq \sum_{t=1}^{+\infty} \ln(t) \eta_{t-\mathcal{J}_t}^2$:

$$\sum_{t=t_0}^T \eta_t \mathbb{E} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_{t,0})) \right] \leq EDGQ \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=t_0}^T \frac{\ln(2C_P H t)}{\ln(1/\lambda(\mathbf{P}))} \eta_{t-\mathcal{J}_t}^2 \right) \quad (142)$$

$$\leq EDGQ \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=1}^{+\infty} \frac{\ln(2C_P H t)}{\ln(1/\lambda(\mathbf{P}))} \eta_{t-\mathcal{J}_t}^2 \right) = \frac{C_1}{\ln(1/\lambda(\mathbf{P}))}. \quad (143)$$

Finally, we conclude that C_1 is finite. To this purpose, we observe that $\mathcal{J}_t \leq a \ln(t) + b$, for opportune positive values a and b . Let t' be a positive integer such that $t \geq a \ln(t) + b$ for any $t \geq t'$. Then:

$$\sum_{t=t'}^T \ln(t) \cdot \eta_{t-\mathcal{J}_t}^2 = \sum_{t=t'-\mathcal{J}_t}^{T-\mathcal{J}_t} \ln(t + \mathcal{J}_t) \cdot \eta_t^2 \quad (144)$$

$$\leq \sum_{t=1}^{+\infty} \ln(t + a \ln t + b) \cdot \eta_t^2 \quad (145)$$

$$\leq \sum_{t=1}^{+\infty} \ln((1 + a + b)t) \cdot \eta_t^2 < +\infty. \quad (146)$$

□

Lemma 15. Let Assumptions 2, 3 and 5 hold, the local functions $\{F_k\}_{k=1}^N$ be convex, and $\mathcal{J}_t \leq t$ be as in Definition 1. Let the step-size be decreasing and satisfy: $\sum_{t=1}^{+\infty} \ln(t) \cdot \eta_t^2 < +\infty$. For any $t_0 \leq T$, we have:

$$\left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=t_0}^T \eta_t \mathbb{E} [F_B(\mathbf{w}_{t,0}) - F_B(\mathbf{w}_{t-\mathcal{J}_t,0})] \leq \frac{C_1}{\ln(1/\lambda(\mathbf{P}))} < +\infty, \quad (147)$$

where:

$$C_1 := EDGQ \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=1}^{+\infty} \ln(2C_P H t) \eta_{t-\mathcal{J}_t}^2 \right). \quad (148)$$

Proof of Lemma 15. This proof is based on [15, Equation (6.38)].

$$\left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=t_0}^T \eta_t \mathbb{E} [F_B(\mathbf{w}_{t,0}) - F_B(\mathbf{w}_{t-\mathcal{J}_t,0})] = \sum_{t=t_0}^T \eta_t \sum_{k=1}^N \pi_k q_k \mathbb{E} [F_k(\mathbf{w}_{t,0}) - F_k(\mathbf{w}_{t-\mathcal{J}_t,0})] \quad (149)$$

$$\leq D \left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=t_0}^T \eta_t \mathbb{E} \|\mathbf{w}_{t-\mathcal{J}_t,0} - \mathbf{w}_{t,0}\| \quad (150)$$

$$\leq D \left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=t_0}^T \eta_t \sum_{d=t-\mathcal{J}_t}^{t-1} \mathbb{E}_{\mathcal{A}_d, \mathcal{H}_d} \left[\mathbb{E}_{\mathcal{B}_d | \mathcal{A}_d, \mathcal{H}_d} \|\mathbf{w}_{d,0} - \mathbf{w}_{d+1,0}\| \right] \quad (151)$$

$$\leq DEGQ \left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=t_0}^T \sum_{d=t-\mathcal{J}_t}^{t-1} \eta_t \eta_d \quad (152)$$

$$\leq \frac{DEGQ}{2} \left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=t_0}^T \sum_{d=t-\mathcal{J}_t}^{t-1} (\eta_t^2 + \eta_d^2) \quad (153)$$

$$\leq DEGQ \left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=t_0}^T \mathcal{J}_t \cdot \eta_{t-\mathcal{J}_t}^2 \quad (154)$$

$$\leq EDGQ \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=1}^{+\infty} \frac{\ln(2C_P H t)}{\ln(1/\lambda(\mathbf{P}))} \eta_{t-\mathcal{J}_t}^2 \right) = \frac{C_1}{\ln(1/\lambda(\mathbf{P}))}, \quad (155)$$

where, in (149), we applied $F_B(\mathbf{w}) = \sum_{k=1}^N p_k F_k(\mathbf{w})$, where $p_k = \frac{\pi_k q_k}{\sum_{h=1}^N \pi_h q_h}$; in (150), we applied Lemma 12; in (151), we applied the triangle inequality and the law of total expectation; in (152), we applied Lemma 13; in (153), we used $2ab \leq a^2 + b^2$; in (154), we observed that $\eta_t^2 + \eta_d^2 \leq 2\eta_{t-\mathcal{J}_t}^2$ due to the diminishing learning rate; finally, in (155), we applied the definition of \mathcal{J}_t given in (132) and we observed that $\sum_{t=t_0}^{+\infty} \ln(t) \eta_{t-\mathcal{J}_t}^2 \leq \sum_{t=1}^{+\infty} \ln(t) \eta_{t-\mathcal{J}_t}^2 < +\infty$ and then $C_1 < +\infty$. \square

Lemma 16 (Bound on the distance dynamics between the current and the stationary distributions of the Markov process). Let Assumption 1 hold, and \mathbf{P} , ρ defined therein. The following inequality holds:

$$\max_{i,j \in [M]} |[\mathbf{P}^t]_{i,j} - \rho_j| \leq C_P \cdot \lambda(\mathbf{P})^t, \quad \text{for } t \geq T_P, \quad (5)$$

where C_P and T_P are positive constants defined as:

$$C_P := \left(\sum_{i=2}^d n_i^2 \right)^{\frac{1}{2}} \cdot \|\mathbf{U}\|_F \|\mathbf{U}^{-1}\|_F, \quad (156)$$

$$T_P := \max \left\{ \max_{1 \leq i \leq d} \left\{ \left\lceil \frac{2n_i(n_i-1)(\ln(\frac{2n_i}{\ln \lambda(\mathbf{P})/|\lambda_2(\mathbf{P})|}) - 1)}{(n_i+1) \ln(\lambda(\mathbf{P})/|\lambda_2(\mathbf{P})|)} \right\rceil \right\}, 0 \right\}. \quad (157)$$

Here, d , n_i , and \mathbf{U} are quantities related to the Jordan canonical form of \mathbf{P} . Specifically, $\mathbf{P} = \mathbf{U} \mathbf{J} \mathbf{U}^{-1}$, where \mathbf{J} denotes the Jordan $M \times M$ matrix with d blocks \mathbf{J}_i , $i = 2, \dots, d$. Each block \mathbf{J}_i , $i = 2, 3, \dots, d$, has a dimension $n_i \geq 1$, and $\sum_{i=1}^d n_i = M$. Moreover, $\|\mathbf{U}\|_F$ denotes the Frobenius norm of the matrix \mathbf{U} .

Furthermore, let Assumptions 2 and 3 hold, H be defined as in Lemma 1, Equation (8), and $T_P \leq \mathcal{J}_t \leq t$ be defined in (132). We obtain the additional inequality:

$$|[\mathbf{P}^{\mathcal{J}_t}]_{i,j} - \rho_j| \leq C_P \cdot \lambda(\mathbf{P})^t \leq C_P \lambda(\mathbf{P})^{\mathcal{J}_t} = \frac{1}{2Ht}, \quad \forall i, j \in [M] \text{ and } \forall t \geq T_P. \quad (158)$$

Proof of Lemma 16. The inequality in (5) is proven in [15, Lemma 1] and holds for any $t \geq T_P$. Here, T_P is a constant dependent on the transition matrix \mathbf{P} of the Markov chain $(\mathcal{A}_t)_{t \geq 0}$ defined in Assumption 1. To prove (158), we further observe that $0 < \lambda(\mathbf{P}) \leq 1$ and $T_P \leq \mathcal{J}_t \leq t$. The last inequality in (158) follows from the definition of \mathcal{J}_t in (132). \square

We remark that the bounds in [15, Lemma 1], and consequently our (158), require $t \geq T_P$. Therefore, the derivations in [15, (6.28)] and [15, (6.35)–(6.37)] are not accurate, since they hold for $t \geq T_P$. We address this problem with Lemmas 17 and 18.

Lemma 17. *Let Assumptions 1–3 hold, and T_P be defined as in (157). The following inequality holds:*

$$\left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=1}^{T_P-1} \eta_t \mathbb{E} [F_B(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_B^*] \leq C_2 < +\infty, \quad (159)$$

where:

$$C_2 := H \left(\sum_{t=1}^{T_P-1} \eta_t \right) \left(\sum_{k=1}^N \pi_k q_k \right) < +\infty. \quad (160)$$

Proof of Lemma 17.

$$\left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=1}^{T_P-1} \eta_t \mathbb{E} [F_B(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_B^*] = \sum_{t=1}^{T_P-1} \eta_t \sum_{k=1}^N \pi_k q_k \mathbb{E} [F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_B^*)] \quad (161)$$

$$\leq H \left(\sum_{t=1}^{T_P-1} \eta_t \right) \left(\sum_{k=1}^N \pi_k q_k \right) := C_2 < +\infty, \quad (162)$$

where, in (161), we used the definition of F_B from (4), and in (162), we applied Lemma 1, Equation (8), which holds for any $\mathbf{w} \in W$. Lastly, it is worth noting that C_2 is a sum of finite elements, and is therefore finite. \square

Lemma 18. *Let Assumptions 1–3 and 5 hold, and $\{F_k\}_{k=1}^N$ be convex. Recall the definitions of \mathcal{J}_t and T_P in (132) and in (157), respectively. Let the step-size $(\eta_t)_{t \geq 1}$ decrease and satisfy $\eta_1 \leq \frac{1}{2L(1+2EQ)}$, $\sum_{t=1}^{+\infty} \eta_t^2 < +\infty$, and $\sum_{t=1}^{+\infty} \ln(t) \cdot \eta_t^2 < +\infty$. For $t \geq T_P$, we have:*

$$\left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=T_P}^T \eta_t \mathbb{E} [F_B(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_B^*] \leq \frac{C_1}{\ln(1/\lambda(\mathbf{P}))} + C_3 < +\infty, \quad (163)$$

where:

$$C_1 := EDGQ \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=1}^{+\infty} \ln(2C_P H t) \cdot \eta_{t-\mathcal{J}_t}^2 \right) < +\infty. \quad (164)$$

$$C_3 := C_0 + \frac{MQ}{4} \sum_{t=1}^{+\infty} \left(\eta_t^2 + \frac{1}{t^2} \right) < +\infty; \quad (165)$$

Proof of Lemma 18. Assume $t \geq T_P$. With a similar proof technique to [15, (6.35)], we derive the following lower bound:

$$\begin{aligned} \mathbb{E}_{\mathcal{A}_t | \mathcal{A}_{t-\mathcal{J}_t}, \mathcal{H}_{t-\mathcal{J}_t}} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_B^*)) \right] &= \\ &= \sum_{a \in \mathcal{M}} \mathbb{P}(\mathcal{A}_t = a | \mathcal{A}_{t-\mathcal{J}_t}, \mathcal{H}_{t-\mathcal{J}_t}) \sum_{k \in a} q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_B^*)) \end{aligned} \quad (166)$$

$$= \sum_{a \in \mathcal{M}} [\mathbf{P}^{\mathcal{J}_t}]_{\mathcal{A}_{t-\mathcal{J}_t}, a} \sum_{k \in a} q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_B^*)) \quad (167)$$

$$\geq \sum_{a \in \mathcal{M}} \left(\rho_a - \frac{1}{2Ht} \right) \sum_{k \in a} q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_B^*)) \quad (168)$$

$$= \sum_{k=1}^N \mathbb{E} [\mathbb{1}_{k \in \mathcal{A}_t}] q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_B^*)) - \frac{1}{2Ht} \sum_{a \in \mathcal{M}} \sum_{k \in a} q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_B^*)) \quad (169)$$

$$\geq \sum_{k=1}^N \pi_k q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_B^*)) - \frac{MQ}{2Ht} \max_{k \in \mathcal{K}} \{F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_B^*)\} \quad (170)$$

$$\geq \left(\sum_{k=1}^N \pi_k q_k \right) \cdot (F_B(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_B^*) - \frac{MQ}{2t}, \quad (171)$$

where, in (166), we applied the definition of expected value to the random variable \mathcal{A}_t , with a representing a realization of \mathcal{A}_t , that is a state in the state space \mathcal{M} , and $\mathbb{P}(\mathcal{A}_t = a \mid \mathcal{A}_{t-\mathcal{J}_t}, \mathcal{H}_{t-\mathcal{J}_t})$ denoting the conditional probability of the event $\mathcal{A}_t = a$ given $(\mathcal{A}_{t-\mathcal{J}_t}, \mathcal{H}_{t-\mathcal{J}_t})$; in (167), we applied the Markov property (Assumption 1), observing that $\mathbb{P}(\mathcal{A}_t = a \mid \mathcal{A}_{t-\mathcal{J}_t}) = [\mathbf{P}^{\mathcal{J}_t}]_{\mathcal{A}_{t-\mathcal{J}_t}, a}$, where $[\mathbf{P}^k]_{i,j}$ denotes the (i, j) -th element of the k -th power of the transition matrix \mathbf{P} ; in (168), we applied Lemma 16, Equation (158); for the first term in (169), we used $\sum_{a \in \mathcal{M}} \rho_a \sum_{k \in a} f(k) = \sum_{a \in \mathcal{M}} \rho_a \sum_{k=1}^N \mathbb{1}_{\{k \in a\}} f(k) = \sum_{k=1}^N f(k) \sum_{a \in \mathcal{M}} \rho_a \mathbb{1}_{k \in a} = \sum_{k=1}^N f(k) \mathbb{E}[\mathbb{1}_{k \in \mathcal{A}_t}]$, where $\mathbb{1}_{k \in \mathcal{A}_t}$ is the indicator function that equals 1 if and only if $k \in \mathcal{A}_t$; in (170), we used $\mathbb{E}[\mathbb{1}_{k \in \mathcal{A}_t}] = \mathbb{P}(k \in \mathcal{A}_t) := \pi_k$ for the first term, and $\sum_{k \in a} q_k f(k) \leq \sum_{k=1}^N q_k f(k) \leq (\sum_{k=1}^N q_k) (\max_{k \in \mathcal{K}} f(k)) = Q \max_{k \in \mathcal{K}} f(k)$ and $\sum_{a \in \mathcal{M}} 1 = M$ for the second term; finally, in (171), we used the definition of F_B in (4) for the first term, and we used Lemma 1, Equation (8) for the second term.

Our derivations in (170) and (171) correct a typo in [15, (6.35)], which considered $Q/(2t)$ instead of $(MQ)/(2t)$. In (171), the dimension (M) of the state space (\mathcal{M}) of the Markov chain $(\mathcal{A}_t)_{t \geq 0}$ appears in the numerator of the second term.

Note that the steps in (168)–(171) require $t \geq T_P$. Multiplying by η_t and summing for $t = T_P, \dots, T$, rearranging, and computing the total expectation, we obtain the following inequality:

$$\left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=T_P}^T \eta_t \mathbb{E} [F_B(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_B^*] \leq \sum_{t=T_P}^T \eta_t \mathbb{E} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_B^*)) \right] + \frac{MQ}{2} \sum_{t=T_P}^T \frac{\eta_t}{t} \quad (172)$$

$$\leq \underbrace{\sum_{t=T_P}^T \eta_t \mathbb{E} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_B^*)) \right]}_{\text{bounded with Lemma 11 + Lemma 14}} + \frac{MQ}{4} \sum_{t=1}^T \left(\eta_t^2 + \frac{1}{t^2} \right), \quad (173)$$

where, in (173), we used $2ab \leq a^2 + b^2$ and we observed that $\sum_{t=T_P}^T (\eta_t^2 + \frac{1}{t^2}) \leq \sum_{t=1}^T (\eta_t^2 + \frac{1}{t^2})$ since $t > 0$ and $\eta_t > 0$.

Moreover, if the step-size $(\eta_t)_{t \geq 1}$ decreases and satisfies $\eta_1 \leq \frac{1}{2L(1+2EQ)}$, $\sum_{t=1}^{+\infty} \eta_t^2 < +\infty$, and $\sum_{t=1}^{+\infty} \ln(t) \cdot \eta_t^2 < +\infty$, we can further bound the first term in (173) by combining Lemma 11 and Lemma 14 for $t_0 = T_P$, and we obtain:

$$\sum_{t=T_P}^T \eta_t \mathbb{E} \left[\sum_{k \in \mathcal{A}_t} q_k (F_k(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_k(\mathbf{w}_B^*)) \right] \leq C_0 + \frac{C_1}{\ln(1/\lambda(\mathbf{P}))} < +\infty, \quad (174)$$

where:

$$\begin{aligned} C_0 &:= \frac{2}{E} \text{diam}(W)^2 + (E+1) \left(\sum_{k=1}^N \pi_k q_k^2 \sigma_k^2 \right) \left(\sum_{t=1}^{+\infty} \eta_t^2 \right) \\ &\quad + 2E^2 G^2 \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=1}^{+\infty} \eta_t^2 \right) \\ &\quad + 4L(1+EQ) \Gamma \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=1}^{+\infty} \eta_t^2 \right). \end{aligned} \quad (175)$$

Finally, plugging (174) into (173), observing that $\sum_{t=1}^T (\eta_t^2 + \frac{1}{t^2}) \leq \sum_{t=1}^{+\infty} (\eta_t^2 + \frac{1}{t^2}) < +\infty$ because $\sum_{t=1}^{+\infty} \eta_t^2 < +\infty$ and $\sum_{t=1}^{+\infty} \frac{1}{t^2} = \frac{\pi}{6} < +\infty$, and denoting $C_3 := C_0 + \frac{MQ}{4} \sum_{t=1}^{+\infty} (\eta_t^2 + \frac{1}{t^2}) < +\infty$, we conclude the proof of Lemma 18. \square

B3. Proof of Theorem 2

Theorem 2 (Convergence of the optimization error ϵ_{opt}). *Let Assumptions 1–3 and 5 hold and the functions $\{F_k\}_{k=1}^N$ be convex. Recall the constants $M, L, D, G, H, \Gamma, \sigma_k, C_P, T_P, \mathcal{J}_t$, and $\lambda(\mathbf{P})$ defined above. Let $Q = \sum_{k \in \mathcal{K}} q_k$. Let the step-size $\eta_t > 0$ decrease and satisfy:*

$$\eta_1 \leq \frac{1}{2L(1+2EQ)}, \quad \sum_{t=1}^{+\infty} \eta_t = +\infty, \quad \sum_{t=1}^{+\infty} \ln(t) \cdot \eta_t^2 < +\infty. \quad (12)$$

Let T denote the total communication rounds.

For $T \geq T_P$, the expected optimization error $\mathbb{E}[F_B(\bar{\mathbf{w}}_{T,0}) - F_B^*]$ can be bounded as follows:

$$\mathbb{E}[F_B(\bar{\mathbf{w}}_{T,0}) - F_B^*] \leq \frac{\frac{1}{2} \mathbf{q}^\top \boldsymbol{\Sigma} \mathbf{q} + v + \psi + \frac{\phi}{\ln(1/\lambda(\mathbf{P}))}}{\left(\sum_{t=1}^T \eta_t\right)}, \quad (13)$$

where $\bar{\mathbf{w}}_{T,0} = \frac{\sum_{t=1}^T \eta_t \mathbf{w}_{t,0}}{\sum_{t=1}^T \eta_t}$, and:

$$\boldsymbol{\Sigma} := \text{diag} \left(2(E+1) \pi_k \sigma_k^2 \sum_{t=1}^{+\infty} \eta_t^2 \right); \quad (176)$$

$$v := \frac{2}{E} \text{diam}(W)^2 + \frac{MQ}{4} \sum_{t=1}^{+\infty} \left(\eta_t^2 + \frac{1}{t^2} \right); \quad (177)$$

$$\psi := 4L(1+EQ)\Gamma \left(\sum_{t=1}^{+\infty} \eta_t^2 \right) + 2E^2G^2 \left(\sum_{t=1}^{+\infty} \eta_t^2 \right) + H \left(\sum_{t=1}^{T_P-1} \eta_t \right); \quad (178)$$

$$\phi := 2EDGQ \left(\sum_{t=1}^{+\infty} \ln(2C_P H t) \cdot \eta_{t-\mathcal{J}_t}^2 \right). \quad (179)$$

Proof of Theorem 2. The proof involves three main steps.

Step 1: From Lemma 15, observe that:

$$\left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=1}^T \eta_t \mathbb{E}[F_B(\mathbf{w}_{t,0}) - F_B(\mathbf{w}_{t-\mathcal{J}_t,0})] \leq \frac{C_1}{\ln(1/\lambda(\mathbf{P}))} < +\infty, \quad (180)$$

where:

$$C_1 := EDGQ \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=1}^{+\infty} \ln(2C_P H t) \cdot \eta_{t-\mathcal{J}_t}^2 \right) < +\infty. \quad (181)$$

Step 2: By combining Lemma 17 and Lemma 18, we obtain:

$$\left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=1}^T \eta_t \mathbb{E}[F_B(\mathbf{w}_{t-\mathcal{J}_t,0}) - F_B^*] \leq \frac{C_1}{\ln(1/\lambda(\mathbf{P}))} + C_2 + C_3 < +\infty, \quad (182)$$

where C_1 is defined in (181), and:

$$C_2 := H \left(\sum_{t=1}^{T_P-1} \eta_t \right) \left(\sum_{k=1}^N \pi_k q_k \right) < +\infty; \quad (183)$$

$$\begin{aligned} C_3 &:= \frac{2}{E} \text{diam}(W)^2 + (E+1) \left(\sum_{k=1}^N \pi_k q_k^2 \sigma_k^2 \right) \left(\sum_{t=1}^{+\infty} \eta_t^2 \right) \\ &\quad + 2E^2G^2 \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=1}^{+\infty} \eta_t^2 \right) \\ &\quad + 4L(1+EQ)\Gamma \left(\sum_{k=1}^N \pi_k q_k \right) \left(\sum_{t=1}^{+\infty} \eta_t^2 \right) + \frac{MQ}{4} \sum_{t=1}^{+\infty} \left(\eta_t^2 + \frac{1}{t^2} \right) < +\infty. \end{aligned} \quad (184)$$

Step 3: By summing the results from Steps 1 and 2, given in (180) and (182), respectively, we have:

$$\left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=1}^T \eta_t \mathbb{E}[F_B(\mathbf{w}_{t,0}) - F_B^*] \leq \frac{2C_1}{\ln(1/\lambda(\mathbf{P}))} + C_2 + C_3 < +\infty. \quad (185)$$

With the convexity of $F_B(\cdot)$, applying the Jensen's inequality, we complete Step 3:

$$\left(\sum_{t=1}^T \eta_t \right) \left(\sum_{k=1}^N \pi_k q_k \right) \mathbb{E}[F_B(\bar{\mathbf{w}}_{T,0}) - F_B^*] \leq \left(\sum_{k=1}^N \pi_k q_k \right) \sum_{t=1}^T \eta_t \mathbb{E}[F_B(\mathbf{w}_{t,0}) - F_B^*] \quad (186)$$

$$\leq \frac{2C_1}{\ln(1/\lambda(\mathbf{P}))} + C_2 + C_3 < +\infty, \quad (187)$$

where $\bar{\mathbf{w}}_{T,0} := \frac{\sum_{t=1}^T \eta_t \mathbf{w}_{t,0}}{\sum_{t=1}^T \eta_t}$, and the constants C_1 , C_2 , and C_3 are defined in (181), (183), and (184), respectively.

By dividing (186) and (187) by $\left(\sum_{t=1}^T \eta_t \right) \cdot \left(\sum_{k=1}^N \pi_k q_k \right)$, we obtain the expression for Theorem 2 given in (13). \square

APPENDIX C PROOF OF THEOREM 3

Theorem 3 (An alternative bound on the bias error ϵ_{bias}). *Under the same assumptions of Theorem 1, define $\Gamma' := \max_k \{F_k(\mathbf{w}_B^*) - F_k^*\}$. The following result holds:*

$$\epsilon_{\text{bias}} \leq 4\kappa^2 \cdot \underbrace{d_{TV}^2(\boldsymbol{\alpha}, \mathbf{p})}_{:= \bar{\epsilon}'_{\text{bias}}} \cdot \Gamma', \quad (15)$$

where $d_{TV}(\boldsymbol{\alpha}, \mathbf{p}) := \frac{1}{2} \sum_{k=1}^N |\alpha_k - p_k|$ denotes the total variation distance between the probability distributions $\boldsymbol{\alpha}$ and \mathbf{p} .

Proof of Theorem 3. The proof follows the same steps as in Theorem 1, proceeding from (29) as follows:

$$\|\nabla F(\mathbf{w}_B^*)\| \leq L \sqrt{\frac{2}{\mu}} \sum_{k=1}^N |\alpha_k - p_k| \sqrt{(F_k(\mathbf{w}_B^*) - F_k^*)} \quad (29)$$

$$\leq 2L \sqrt{\frac{2}{\mu}} d_{TV}(\boldsymbol{\alpha}, \mathbf{p}) \sqrt{\Gamma'}, \quad (188)$$

where, in (188), we applied the definitions of $d_{TV}(\boldsymbol{\alpha}, \mathbf{p}) := \frac{1}{2} \sum_{k=1}^N |\alpha_k - p_k|$ and $\Gamma' := \max_k \{F_k(\mathbf{w}_B^*) - F_k^*\}$.

Squaring (188), we obtain the following expression:

$$\|\nabla F(\mathbf{w}_B^*)\|^2 \leq \frac{8L^2}{\mu} d_{TV}^2(\boldsymbol{\alpha}, \mathbf{p}) \Gamma'. \quad (189)$$

Then, replacing (189) in (25), we obtain:

$$\epsilon_{\text{bias}} := (F(\mathbf{w}_B^*) - F^*) \leq \frac{1}{2\mu} \|\nabla F(\mathbf{w}_B^*)\|^2 \leq 4 \frac{L^2}{\mu^2} \underbrace{d_{TV}^2(\boldsymbol{\alpha}, \mathbf{p}) \Gamma'}_{:= \bar{\epsilon}'_{\text{bias}}}, \quad (190)$$

which concludes the proof of Theorem 3. \square

APPENDIX D CONVEXITY OF $\bar{\epsilon}_{\text{OPT}} + \bar{\epsilon}_{\text{BIAS}}$

For the proof of the convexity of $\bar{\epsilon}_{\text{opt}}(\mathbf{q})$, please refer to Appendix E1. To prove that $\bar{\epsilon}_{\text{bias}}(\mathbf{q})$ is also convex, we need to study the convexity of $\chi_{\boldsymbol{\alpha} \parallel \mathbf{p}}^2 := \sum_{k=1}^N (\alpha_k - p_k)^2 / p_k$ in $\mathbf{q} \in \{q_k > 0 \forall k, \|\mathbf{q}\|_1 = Q > 0\}$. To this purpose, we define the following functions:

$$h_k : \mathbb{R}_{\geq 0}^N \setminus \{\mathbf{0}\} \rightarrow \mathbb{R}_{\geq 0}, \quad h_k(\mathbf{q}) := \frac{\pi_k q_k}{\sum_{k'=1}^N \pi_{k'} q_{k'}}; \quad (191)$$

$$g_k : \mathbb{R}_{> 0} \rightarrow \mathbb{R}_{\geq 0}, \quad g_k(p_k) := \frac{(p_k - \alpha_k)^2}{p_k}. \quad (192)$$

Finally, we write the chi-square divergence $\chi_{\alpha\|\mathbf{p}}^2$ between the target and biased probability distributions α and \mathbf{p} as:

$$\chi_{\alpha\|\mathbf{p}}^2(\mathbf{q}) = \sum_{k=1}^N (g_k \circ h_k)(\mathbf{q}) = \sum_{k=1}^N g_k(h_k(\mathbf{q})). \quad (193)$$

We observe that:

- $h_k(\mathbf{q})$ is a particular case of linear-fractional functions [40, Example 3.32, p. 97];
- $g_k(\cdot)$ is a convex in p_k over $\mathbb{R}_{>0}$ because sum of convex functions;
- each $g_k \circ h_k$ is quasi-convex in $\mathbf{q} \in \mathbb{R}_{>0}^N$ because composition of a convex function (g_k) and a linear-fractional function (h_k) [40, p. 102].

However, note that the sum of quasi-convex functions is not necessarily quasi-convex.

Proposition 1. *The function $\chi_{\alpha\|\mathbf{p}}^2(\mathbf{q})$ is not convex over $\mathbb{R}_{>0}^N$.*

Proof of Proposition 1. To analyze the convexity of $\chi_{\alpha\|\mathbf{p}}^2(\mathbf{q}) = \sum_{k=1}^N (g_k \circ h_k)(\mathbf{q})$ over $\mathbb{R}_{>0}^N$, a possible approach is to check whether each function $(g_k \circ h_k)(\mathbf{q})$ is convex over $\mathbb{R}_{>0}^N$. In what follows, we show that $(g_k \circ h_k)$ is not convex over $\mathbb{R}_{>0}^N$.

Consider the case when $\pi_k = 1 \forall k \in \mathcal{K}$. We can rewrite $(g_k \circ h_k)(\mathbf{q})$ as follows:

$$(g_k \circ h_k)(\mathbf{q}) = \frac{\left(\frac{q_k}{\|\mathbf{q}\|_1} - \alpha_k\right)^2}{\frac{q_k}{\|\mathbf{q}\|_1}}. \quad (194)$$

We show that this function fails to satisfy the definition of convexity, i.e., $\exists \mathbf{q}, \mathbf{q}' \in \mathbb{R}_{>0}^N, \zeta \in [0, 1]$ such that:

$$(g_k \circ h_k)(\zeta \mathbf{q} + (1 - \zeta)\mathbf{q}') > \zeta (g_k \circ h_k)(\mathbf{q}) + (1 - \zeta) (g_k \circ h_k)(\mathbf{q}'). \quad (195)$$

The left-hand side (LHS) of (195) is:

$$(g_k \circ h_k)(\zeta \mathbf{q} + (1 - \zeta)\mathbf{q}') = \frac{\left(\frac{\zeta q_k + (1 - \zeta)q'_k}{\zeta \|\mathbf{q}\|_1 + (1 - \zeta)\|\mathbf{q}'\|_1} - \alpha_k\right)^2}{\frac{\zeta q_k + (1 - \zeta)q'_k}{\zeta \|\mathbf{q}\|_1 + (1 - \zeta)\|\mathbf{q}'\|_1}}. \quad (196)$$

If we take $\mathbf{q} : \|\mathbf{q}\|_1 = 1, q_k = \alpha_k, \zeta = \frac{1}{2}, \mathbf{q}' = \frac{Q}{N}\mathbf{1}$, and we let $Q \rightarrow +\infty$, then the LHS in (196) converges to:

$$\lim_{Q \rightarrow +\infty} \frac{\left(\frac{\frac{1}{2}\alpha_k + \frac{1}{2}\frac{Q}{N} - \alpha_k}{\frac{1}{2}1 + \frac{1}{2}Q} - \alpha_k\right)^2}{\frac{\frac{1}{2}\alpha_k + \frac{1}{2}\frac{Q}{N}}{\frac{1}{2}1 + \frac{1}{2}Q}} = \frac{\left(\frac{1}{N} - \alpha_k\right)^2}{\frac{1}{N}}. \quad (197)$$

On the other hand, for the same choices of $q_k, \mathbf{q}, \mathbf{q}'$, and ζ , and if we let $Q \rightarrow +\infty$, the right-hand side (RHS) of (195) is:

$$\zeta (g_k \circ h_k)(\mathbf{q}) + (1 - \zeta) (g_k \circ h_k)(\mathbf{q}') = 0 + \frac{1}{2} \frac{\left(\frac{1}{N} - \alpha_k\right)^2}{\frac{1}{N}}. \quad (198)$$

Finally, comparing (197) and (198), we conclude that, for Q large enough, the LHS in (195) is larger than the RHS. \square

Proposition 2. *The function $\chi_{\alpha\|\mathbf{p}}^2(\mathbf{q})$ is convex over $\mathbb{R}_{>0}^N \cap \{\mathbf{q} : \|\mathbf{q}\|_1 = Q > 0\}$.*

Proof of Proposition 2. To verify the convexity of $\chi_{\alpha\|\mathbf{p}}^2(\mathbf{q}) = \sum_{k=1}^N (g_k \circ h_k)(\mathbf{q})$ over $\mathbb{R}_{>0}^N \cap \{\mathbf{q} : \|\mathbf{q}\|_1 = Q > 0\}$, one possible approach is to demonstrate the convexity of each function $(g_k \circ h_k)(\mathbf{q})$ over the set $\mathbb{R}_{>0}^N \cap \{\mathbf{q} : \|\mathbf{q}\|_1 = Q > 0\}$.

We prove this result for a more general case. We show that, if

$$\tilde{g} \text{ is a convex function over its domain } \mathcal{D}_g \quad (199)$$

and

$$\tilde{h}(\mathbf{q}) = \frac{\mathbf{A}\mathbf{q} + b}{\mathbf{c}^\top \mathbf{q} + d}, \quad (200)$$

then

$$\tilde{g} \circ \tilde{h} \text{ is convex over } \mathcal{D} = \mathbb{R}_{>0}^N \cap \{\mathbf{q} : \mathbf{c}^\top \mathbf{q} + d = Q > 0, \frac{\mathbf{A}\mathbf{q} + b}{\mathbf{c}^\top \mathbf{q} + d} \in \mathcal{D}_g\}. \quad (201)$$

It is then sufficient to apply this result to each pair (g_k, h_k) to conclude that $(g_k \circ h_k)$ is convex and then $\chi_{\alpha \|p}^2(\mathbf{q})$ is convex.

By direct inspection, for all $\mathbf{q}, \mathbf{q}' \in \mathcal{D}$, $\forall \zeta \in [0, 1]$, the following equality holds:

$$\left(\tilde{g} \circ \tilde{h}\right)(\zeta \mathbf{q} + (1 - \zeta)\mathbf{q}') = \tilde{g}\left(\tilde{h}(\zeta \mathbf{q} + (1 - \zeta)\mathbf{q}')\right) = \tilde{g}\left(\zeta' \frac{\mathbf{A}\mathbf{q} + b}{\mathbf{c}^\top \mathbf{q} + d} + (1 - \zeta') \frac{\mathbf{A}\mathbf{q}' + b}{\mathbf{c}^\top \mathbf{q}' + d}\right), \quad (202)$$

where:

$$\zeta' = \frac{\zeta(\mathbf{c}^\top \mathbf{q} + d)}{\zeta(\mathbf{c}^\top \mathbf{q} + d) + (1 - \zeta)(\mathbf{c}^\top \mathbf{q}' + d)} \in [0, 1]. \quad (203)$$

Applying the convexity of \tilde{g} , we bound Equation (202) as follows:

$$\tilde{g}\left(\zeta' \frac{\mathbf{A}\mathbf{q} + b}{\mathbf{c}^\top \mathbf{q} + d} + (1 - \zeta') \frac{\mathbf{A}\mathbf{q}' + b}{\mathbf{c}^\top \mathbf{q}' + d}\right) \stackrel{\text{convexity of } \tilde{g}}{\leq} \zeta' \tilde{g}\left(\frac{\mathbf{A}\mathbf{q} + b}{\mathbf{c}^\top \mathbf{q} + d}\right) + (1 - \zeta') \tilde{g}\left(\frac{\mathbf{A}\mathbf{q}' + b}{\mathbf{c}^\top \mathbf{q}' + d}\right) \quad (204)$$

$$= \zeta' \left(\tilde{g} \circ \tilde{h}\right)(\mathbf{q}) + (1 - \zeta') \left(\tilde{g} \circ \tilde{h}\right)(\mathbf{q}'). \quad (205)$$

Finally, to conclude the proof, we show that $\zeta' = \zeta$. This is true because, for any \mathbf{q} and $\mathbf{q}' \in \mathcal{D}$, $\mathbf{c}^\top \mathbf{q} + d = \mathbf{c}^\top \mathbf{q}' + d = Q > 0$. In fact, by using this condition in Equation (203), we have that:

$$\zeta' = \frac{\zeta Q}{\zeta Q + (1 - \zeta)Q} = \zeta, \quad (206)$$

which establishes the convexity of $\tilde{g} \circ \tilde{h}$ by definition. \square

APPENDIX E MINIMIZING $\bar{\epsilon}_{\text{OPT}}$

Equation (13) can be rewritten as:

$$\left(\sum_{t=1}^T \eta_t\right) \mathbb{E}[F_B(\bar{\mathbf{w}}_{T,0}) - F_B^*] \leq \frac{\frac{1}{2} \mathbf{q}^\top \Sigma \mathbf{q} + v}{\boldsymbol{\pi}^\top \mathbf{q}} + \psi + \frac{\phi}{\ln(1/\lambda(\mathbf{P}))} \quad (207)$$

$$= \frac{\frac{1}{2} \mathbf{q}^\top \mathbf{A} \mathbf{q} + B}{\boldsymbol{\pi}^\top \mathbf{q}} + C := J(\mathbf{q}), \quad (208)$$

where:

$$\mathbf{A} := \Sigma = \text{diag}\left(2(E+1)\pi_k \sigma_k^2 \sum_{t=1}^{+\infty} \eta_t^2\right); \quad (209)$$

$$B := v = \frac{2}{E} \text{diam}(W)^2 + \frac{MQ}{4} \sum_{t=1}^{+\infty} \left(\eta_t^2 + \frac{1}{t^2}\right); \quad (210)$$

$$C := \psi + \frac{\phi}{\ln(1/\lambda(\mathbf{P}))} = (4L(1 + EQ)\Gamma + 2E^2G^2) \left(\sum_{t=1}^{+\infty} \eta_t^2\right) + 2EDGQ \left(\sum_{t=1}^{+\infty} \mathcal{J}_t \cdot \eta_{t-\mathcal{J}_t}^2\right) + H \left(\sum_{t=1}^{T_P-1} \eta_t\right). \quad (211)$$

The minimization of (208), defines the following optimization problem:

$$\underset{\mathbf{q}}{\text{minimize}} \quad J(\mathbf{q}) := \frac{\frac{1}{2} \mathbf{q}^\top \mathbf{A} \mathbf{q} + B}{\boldsymbol{\pi}^\top \mathbf{q}} + C; \quad (212a)$$

$$\text{subject to} \quad \mathbf{q} \geq 0, \quad (212b)$$

$$\boldsymbol{\pi}^\top \mathbf{q} > 0, \quad (212c)$$

$$\|\mathbf{q}\|_1 = Q. \quad (212d)$$

Remark. In Problem (212a)–(212d), when setting some q_k to zero, we do not consider the possibility of redefining the Markov chain $(\mathcal{A}_t)_{t \geq 0}$ in Assumption 1 by considering the reduced state space of clients with $q_k > 0$. In this case, the redefined Markov chain would have a different transition matrix $\mathbf{P}' \neq \mathbf{P}$ with $\lambda(\mathbf{P}') \neq \lambda(\mathbf{P})$, resulting in C no longer being constant.

E1. The optimization problem in (212a)–(212d) is convex

Let us rewrite the problem by adding a variable $s := 1/\boldsymbol{\pi}^\top \mathbf{q}$ and then replacing $\mathbf{y} := s\mathbf{q}$. We have:

$$J(\mathbf{y}, s) = s \left(\frac{1}{2} \frac{\mathbf{y}^\top \mathbf{A} \mathbf{y}}{s} + B \right) + C = s \cdot K \left(\frac{\mathbf{y}}{s} \right) + C, \quad (213)$$

where $K : \mathbb{R}^N \rightarrow \mathbb{R}$, $K(\mathbf{q}) := \frac{1}{2} \mathbf{q}^\top \mathbf{A} \mathbf{q} + B$ is a (strictly) convex function, and:

$$\underset{\mathbf{y}, s}{\text{minimize}} \quad J(\mathbf{y}, s) = \frac{1}{2s} \mathbf{y}^\top \mathbf{A} \mathbf{y} + Bs + C \quad (214a)$$

$$\text{subject to} \quad \mathbf{y} \geq 0, \quad (214b)$$

$$s > 0, \quad (214c)$$

$$\boldsymbol{\pi}^\top \mathbf{y} = 1, \quad (214d)$$

$$\|\mathbf{y}\|_1 = Qs. \quad (214e)$$

Note that the objective function $J(\mathbf{y}, s) : \mathbb{R}^{N+1} \rightarrow \mathbb{R}$, $J(\mathbf{y}, s) = s \cdot K(\mathbf{y}/s) + C$ in (213) is the perspective of the convex function $K(\mathbf{q}) + C$, and is therefore convex [40, pp. 89–90]. Moreover, the constraints in (214b)–(214e) define a convex set, and then the optimization problem defined by (214a)–(214e) is convex. We solve it with the method of Lagrange multipliers.

E2. Support for Guideline A (Section III)

The Lagrangian function \mathcal{L} is as follows:

$$\mathcal{L}(\mathbf{y}, s, \iota, \theta, \boldsymbol{\omega}) = \frac{1}{2s} \mathbf{y}^\top \mathbf{A} \mathbf{y} + Bs + C + \iota(1 - \boldsymbol{\pi}^\top \mathbf{y}) + \theta(\|\mathbf{y}\|_1 - Qs) - \boldsymbol{\omega}^\top \mathbf{y}. \quad (215)$$

Since the constraint $s > 0$ defines an open set, the set defined by the constraints in (214b)–(214e) is not closed. However, the solution of the optimization problem defined by (214a)–(214e) is never on the boundary $s = 0$ because $\mathcal{L} \rightarrow +\infty$ as $s \rightarrow 0^+$, therefore we can consider $s \geq 0$. Moreover, strong duality holds for the Slater's constraint qualification for convex problems.

The KKT conditions read:

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial s}(\mathbf{y}^*, s^*, \iota^*, \theta^*, \boldsymbol{\omega}^*) = 0, & (216) \\ \nabla_{\mathbf{y}} \mathcal{L}(\mathbf{y}^*, s^*, \iota^*, \theta^*, \boldsymbol{\omega}^*) = 0, & (217) \\ \boldsymbol{\pi}^\top \mathbf{y}^* - 1 = 0, & (218) \\ \|\mathbf{y}^*\|_1 - Qs^* = 0, & (219) \\ \boldsymbol{\omega}^{*\top} \mathbf{y}^* = 0, & (220) \\ \mathbf{y}^*, \boldsymbol{\omega}^* \geq 0. & (221) \end{cases}$$

In particular, the KKT condition for \mathbf{y}^* read:

$$\nabla_{\mathbf{y}} \mathcal{L}(\mathbf{y}^*, s^*, \iota^*, \theta^*, \boldsymbol{\omega}^*) = \frac{1}{s^*} \mathbf{A} \mathbf{y}^* - \iota^* \boldsymbol{\pi} + \theta^* \mathbf{1} - \boldsymbol{\omega}^* = 0, \quad (222)$$

which is satisfied when:

$$\frac{\partial \mathcal{L}}{\partial y_k^*} = \frac{1}{s^*} A_{kk} y_k^* - \iota^* \pi_k + \theta^* - \omega_k^* = 0, \quad \forall k \in \mathcal{K}, \quad (223)$$

where A_{ij} denotes the element on the i -th row and the j -th column of matrix \mathbf{A} .

Furthermore, the Complementary Slackness conditions in (220) and (221) present two cases:

1) If $y_k^* > 0$ (and $q_k^* > 0$), then $\omega_k^* = 0$ and:

$$y_k^* = \frac{s^*}{A_{kk}} (\iota^* \pi_k - \theta^*), \quad q_k^* = \frac{1}{A_{kk}} (\iota^* \pi_k - \theta^*); \quad (224)$$

2) $y_k^* = q_k^* = 0$ otherwise.

By replacing the equality constraint (214d) in Problem (214a)–(214e) with the inequality constraint $\pi^\top \mathbf{y} \geq 1$, we establish an equivalent optimization problem. The equivalence holds because, for any feasible solution \mathbf{y}' with $\pi^\top \mathbf{y}' > 1$, we can consider the solution $\mathbf{y}'' = \frac{\mathbf{y}'}{\pi^\top \mathbf{y}'} < \mathbf{y}'$, leading to a lower objective function value. Additionally, the new problem states that the Lagrange multiplier (ι^*) associated with the inequality constraint must be non-negative. By considering $A_{kk} \geq 0$ and $\iota^* \geq 0$ in Equation (224), we conclude that q_k^* increases with π_k , providing analytical support for Guideline A.

E3. Closed-form solution of the optimization problem in (212a)–(212d)

The solution of the optimization problem in (212a)–(212d) is not of practical utility because its constants (e.g., L , ω , Γ , C_P) are in general problem-dependent and difficult to estimate during training. In particular, Γ poses particular difficulties as it is defined in terms of the minimizer of the target objective F , but the FL algorithm generally minimizes the biased function F_B . Nevertheless, we include the closed-form solution of the optimization problem in (212a)–(212d) for completeness.

We use the active-set method: let \mathcal{X} be the set of coordinates corresponding to the active inequalities, i.e., $\mathcal{X} = \{k \mid y_k^* = 0\}$.

From the KKT condition in (218), we derive a relation between ι^* and θ^* :

$$\pi^\top \mathbf{y}^* = \sum_{k \notin \mathcal{X}} \pi_k y_k^* = \sum_{k \notin \mathcal{X}} \pi_k \frac{s^*}{\mathbf{A}_{kk}} (\iota^* \pi_k - \theta^*) = \iota^* s^* \sum_{k \notin \mathcal{X}} \frac{\pi_k^2}{\mathbf{A}_{kk}} - \theta^* s^* \sum_{k \notin \mathcal{X}} \frac{\pi_k}{\mathbf{A}_{kk}} = 1. \quad (225)$$

We use the KKT condition in (219) to derive another relation between ι^* and θ^* :

$$\|\mathbf{y}^*\|_1 = \sum_{k \notin \mathcal{X}} y_k^* = \sum_{k \notin \mathcal{X}} \frac{s^*}{\mathbf{A}_{kk}} (\iota^* \pi_k - \theta^*) = Q s^* \Leftrightarrow \iota^* = \frac{Q + \theta^* \sum_{k \notin \mathcal{X}} \frac{1}{\mathbf{A}_{kk}}}{\sum_{k \notin \mathcal{X}} \frac{\pi_k}{\mathbf{A}_{kk}}}, \quad (226)$$

and, replacing (226) in (225), we derive the closed-form solution for θ^* :

$$\theta^* = \frac{\sum_{k \notin \mathcal{X}} \frac{\pi_k}{\mathbf{A}_{kk}} - Q s^* \sum_{k \notin \mathcal{X}} \frac{\pi_k^2}{\mathbf{A}_{kk}}}{s^* \left[\left(\sum_{k \notin \mathcal{X}} \frac{1}{\mathbf{A}_{kk}} \right) \cdot \left(\sum_{k \notin \mathcal{X}} \frac{\pi_k^2}{\mathbf{A}_{kk}} \right) - \left(\sum_{k \notin \mathcal{X}} \frac{\pi_k}{\mathbf{A}_{kk}} \right)^2 \right]}. \quad (227)$$

APPENDIX F BACKGROUND ON MARKOV CHAINS

F1. Markov Chain for the Analysis (Section III)

We recall some existing results [15], [31] for the Markov chain $(\mathcal{A}_t)_{t \geq 0}$ used in our analysis (Assumption 1).

Assumption 1. *The Markov chain $(\mathcal{A}_t)_{t \geq 0}$ on the M -finite state space \mathcal{M} is time-homogeneous, irreducible, and aperiodic. It has transition matrix \mathbf{P} , stationary distribution $\boldsymbol{\rho}$, and has state distribution $\boldsymbol{\rho}$ at time $t = 0$.*

Let $\boldsymbol{\rho}^{(t)} = [\rho_1^{(t)}, \rho_2^{(t)}, \dots, \rho_M^{(t)}]$, $\sum_{i=1}^M \rho_i^{(t)} = 1$ be the state probability distribution on the Markov chain $(\mathcal{A}_t)_{t \geq 0}$ at time step t . Assumption 1 guarantees the existence of a stationary distribution $\boldsymbol{\rho} = \lim_{t \rightarrow +\infty} \boldsymbol{\rho}^{(t)} = [\rho_1, \rho_2, \dots, \rho_M]$ with $\min_i \{\rho_i\} > 0$ and $\boldsymbol{\rho}^\top \mathbf{P} = \boldsymbol{\rho}^\top$. Then $\boldsymbol{\rho}$ is a left eigenvector relative to the eigenvalue 1, which is the largest eigenvalue of the matrix \mathbf{P} .

For the transition matrix \mathbf{P} , we label its eigenvalues in decreasing order:

$$1 = \lambda_1(\mathbf{P}) > \lambda_2(\mathbf{P}) \geq \dots \geq \lambda_M(\mathbf{P}). \quad (228)$$

We define:

$$\bar{\lambda}_2(\mathbf{P}) := \max \{|\lambda_2(\mathbf{P})|, |\lambda_M(\mathbf{P})|\} \quad \text{and} \quad \lambda(\mathbf{P}) := \frac{\bar{\lambda}_2(\mathbf{P}) + 1}{2}. \quad (229)$$

The second largest absolute eigenvalue $\bar{\lambda}_2(\mathbf{P})$ of the transition matrix \mathbf{P} characterizes the mixing time of a Markov chain. The absolute spectral gap $\gamma := 1 - \bar{\lambda}_2(\mathbf{P})$ and its reciprocal, the relaxation time $t_{\text{rel}} := \frac{1}{\gamma}$, play a role in this relationship. To quantify the convergence of the Markov chain towards stationarity, we use the parameter $d(t) := \max_{a \in \mathcal{M}} \|\mathbf{P}^t|_a - \boldsymbol{\rho}\|_{TV}$, which measures the maximum distance between the distribution $\mathbf{P}^t|_a$ and the stationary distribution $\boldsymbol{\rho}$ for all initial states $a \in \mathcal{M}$. The mixing time $t_{\text{mix}}(\varepsilon)$ is defined as the minimum time at which the distance $d(t)$ becomes less than or equal to a given threshold ε : $t_{\text{mix}}(\varepsilon) := \min \{t : d(t) \leq \varepsilon\}$. Upper and lower bounds exist for the mixing time based on the relaxation time and the stationary distribution: $(t_{\text{rel}} - 1) \log\left(\frac{1}{2\varepsilon}\right) \leq t_{\text{mix}}(\varepsilon) \leq \log\left(\frac{1}{\varepsilon \rho_{\min}}\right) t_{\text{rel}}$, where $\rho_{\min} := \min_{a \in \mathcal{M}} \rho_a$ [31, pp. 154–156].

F2. Markov Chain for Guideline B (Section IV)

In Section III-D (Guideline B), we examine a specific scenario where the availability of each client k follows an independent Markov chain $(\mathcal{A}_t^k)_{t \geq 0}$ with transition probability matrix \mathbf{P}_k . This setup allows us to model the aggregate process as a product of independent Markov chains, known as a Product Chain [31, Section 12.4].

Definition 2 (Product Chain). Let \mathbf{P}_1 and \mathbf{P}_2 be transition matrices on state spaces \mathcal{M}_1 and \mathcal{M}_2 respectively, with corresponding stationary distributions π_1 and π_2 . We consider a Markov Chain on the state space $\mathcal{M}_1 \times \mathcal{M}_2$ that moves independently in the first and second coordinates according to \mathbf{P}_1 and \mathbf{P}_2 respectively. The transition matrix of this Markov Chain is the Kronecker product $\tilde{\mathbf{P}} = \mathbf{P}_1 \otimes \mathbf{P}_2$, defined as:

$$\tilde{\mathbf{P}}((x, y), (z, w)) = \mathbf{P}_1(x, z)\mathbf{P}_2(y, w). \quad (230)$$

Proposition 3. The stationary distribution of the Markov chain defined by $\tilde{\mathbf{P}} = \mathbf{P}_1 \otimes \mathbf{P}_2$ is the Kronecker product $\tilde{\rho} = \pi_1 \otimes \pi_2$.

Proof. We can observe the following:

$$\tilde{\rho}^\top \tilde{\mathbf{P}} = (\pi_1 \otimes \pi_2)^\top \cdot (\mathbf{P}_1 \otimes \mathbf{P}_2) = (\pi_1^\top \mathbf{P}_1) \otimes (\pi_2^\top \mathbf{P}_2) = \pi_1^\top \otimes \pi_2^\top = \tilde{\rho}^\top, \quad (231)$$

where, in (231), we used the mixed-product property of the Kronecker product in the second step, and in the third step, we noted that π_1 and π_2 are the stationary distributions for \mathbf{P}_1 and \mathbf{P}_2 , respectively. For a comprehensive list of properties that the Kronecker product satisfies, please refer to [41, p. 597]. \square

Proposition 4 ([31, Exercise 12.6]). Let \mathbf{u} and \mathbf{v} be eigenvectors of \mathbf{P}_1 and \mathbf{P}_2 , respectively, with eigenvalues λ and μ . Then $\mathbf{u} \otimes \mathbf{v}$ is an eigenvector of $\mathbf{P}_1 \otimes \mathbf{P}_2$ with eigenvalue $\lambda\mu$.

Proof. We can verify the following:

$$(\mathbf{u} \otimes \mathbf{v})^\top (\mathbf{P}_1 \otimes \mathbf{P}_2) = (\mathbf{u}^\top \mathbf{P}_1) \otimes (\mathbf{v}^\top \mathbf{P}_2) = (\lambda \mathbf{u}^\top) \otimes (\mu \mathbf{v}^\top) = \lambda\mu (\mathbf{u} \otimes \mathbf{v})^\top. \quad (232)$$

In (232), we used the mixed-product property and the associativity of the scalar multiplication with the Kronecker product. \square

In general, let \mathbf{P}_1 be a $m \times m$ matrix with eigenvalues $\lambda_1, \dots, \lambda_m$, and \mathbf{P}_2 be a $n \times n$ matrix with eigenvalues μ_1, \dots, μ_n . The complete eigen-decomposition of $\mathbf{P}_1 \otimes \mathbf{P}_2$ depends on the Kronecker product structure and involves combinations of the eigenvalues and eigenvectors of \mathbf{P}_1 and \mathbf{P}_2 .

Proposition 5 (Spectrum of the Kronecker product, [41, Exercise 7.8.11]). Let the eigenvalues of $\mathbf{P}_1 \in \mathbb{R}^{m \times m}$ be denoted by λ_i and let the eigenvalues of $\mathbf{P}_2 \in \mathbb{R}^{n \times n}$ be denoted by μ_j . The eigenvalues of $\mathbf{P}_1 \otimes \mathbf{P}_2$ are the mn numbers $\{\lambda_i \mu_j\}_{i=1, j=1}^{m, n}$.

Proof. Let $\mathbf{J}_1 = \mathbf{A}_1^{-1} \mathbf{P}_1 \mathbf{A}_1$ and $\mathbf{J}_2 = \mathbf{A}_2^{-1} \mathbf{P}_2 \mathbf{A}_2$ be the respective Jordan forms for \mathbf{P}_1 and \mathbf{P}_2 . We use the mixed-product property and the inverse property of the Kronecker product to show that $\mathbf{P}_1 \otimes \mathbf{P}_2$ is similar to $\mathbf{J}_1 \otimes \mathbf{J}_2$:

$$\mathbf{J}_1 \otimes \mathbf{J}_2 = (\mathbf{A}_1^{-1} \mathbf{P}_1 \mathbf{A}_1) \otimes (\mathbf{A}_2^{-1} \mathbf{P}_2 \mathbf{A}_2) = (\mathbf{A}_1^{-1} \otimes \mathbf{A}_2^{-1}) (\mathbf{P}_1 \otimes \mathbf{P}_2) (\mathbf{A}_1 \otimes \mathbf{A}_2) = (\mathbf{A}_1 \otimes \mathbf{A}_2)^{-1} (\mathbf{P}_1 \otimes \mathbf{P}_2) (\mathbf{A}_1 \otimes \mathbf{A}_2). \quad (233)$$

Consequently, the eigenvalues of $\mathbf{P}_1 \otimes \mathbf{P}_2$ coincide with those of $\mathbf{J}_1 \otimes \mathbf{J}_2$. Since \mathbf{J}_1 and \mathbf{J}_2 are upper triangular with $\{\lambda_i\}_{i=1}^m$ and $\{\mu_j\}_{j=1}^n$ on the diagonals, respectively, $\mathbf{J}_1 \otimes \mathbf{J}_2$ is also upper triangular with diagonal entries given by $\{\lambda_i \mu_j\}_{i=1, j=1}^{m, n}$. \square

Proposition 6. Let $\bar{\lambda}_2(\mathbf{P}_k)$ denote the second largest eigenvalue in absolute value of the transition matrix \mathbf{P}_k associated with the k -th client, and define $\lambda(\mathbf{P}_k) := \frac{\bar{\lambda}_2(\mathbf{P}_k) + 1}{2}$. For the product chain defined by $\mathbf{P} = \bigotimes_{k \in \mathcal{K}} \mathbf{P}_k$, the second largest eigenvalue in absolute value $\bar{\lambda}_2(\mathbf{P})$ and $\lambda(\mathbf{P}) := \frac{\bar{\lambda}_2(\mathbf{P}) + 1}{2}$ satisfy:

$$\bar{\lambda}_2(\mathbf{P}) = \max_{k \in \mathcal{K}} \bar{\lambda}_2(\mathbf{P}_k) \quad \text{and} \quad \lambda(\mathbf{P}) = \max_{k \in \mathcal{K}} \lambda(\mathbf{P}_k). \quad (234)$$

The proof of Proposition 6 follows a similar structure to the one in [31, Corollary 12.13].

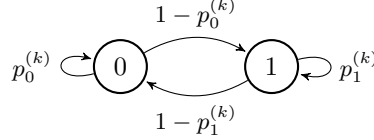
Proof. From Proposition 5, we know that the eigenvalues of $\mathbf{P} = \bigotimes_{k \in \mathcal{K}} \mathbf{P}_k$ are given by:

$$\left\{ \prod_{k \in \mathcal{K}} \lambda_i(\mathbf{P}_k) : \lambda_i(\mathbf{P}_k) \text{ an eigenvalue of } \mathbf{P}_k \right\}. \quad (235)$$

Recall that $\bar{\lambda}_2(\mathbf{P}_k)$ is the second largest eigenvalue of \mathbf{P}_k in absolute value. If k^* denotes the index such that $\bar{\lambda}_2(\mathbf{P}_{k^*}) = \max_{k \in \mathcal{K}} \bar{\lambda}_2(\mathbf{P}_k)$, the second largest eigenvalue in module of \mathbf{P} is the product of $\bar{\lambda}_2(\mathbf{P}_{k^*})$ for the k^* -th client and $\lambda_1(\mathbf{P}_j) = 1$ for the remaining clients $j \neq k^*$. The second result in (234) follows from the definitions of $\lambda(\mathbf{P})$ and $\lambda(\mathbf{P}_k)$. \square

F3. Markov Chain for the Experiments (Section V)

In the experiments (Section V-A), we consider a scenario where the activity of each client $k \in \mathcal{K}$ follows a two-state homogeneous Markov process. The state space \mathcal{M} consists of two states: “inactive” (with value 0) and “active” (with value 1):



We provide detailed expressions of the transition matrix \mathbf{P}_k , stationary distribution $\boldsymbol{\pi}^{(k)}$, and the second eigenvalue $\lambda_2(\mathbf{P}_k)$ used in the experiments for each client $k \in \mathcal{K}$:

$$\mathbf{P}_k = \begin{bmatrix} p_0^{(k)} & 1 - p_0^{(k)} \\ 1 - p_1^{(k)} & p_1^{(k)} \end{bmatrix} = \begin{bmatrix} 1 - (1 - \lambda_2(\mathbf{P}_k))\pi_k & (1 - \lambda_2(\mathbf{P}_k))\pi_k \\ (1 - \lambda_2(\mathbf{P}_k))(1 - \pi_k) & \lambda_2(\mathbf{P}_k) + (1 - \lambda_2(\mathbf{P}_k))\pi_k \end{bmatrix}. \quad (236)$$

$$\boldsymbol{\pi}^{(k)} = [1 - \pi_k, \pi_k] = \left[\frac{1 - p_1^{(k)}}{2 - p_0^{(k)} - p_1^{(k)}}, \frac{1 - p_0^{(k)}}{2 - p_0^{(k)} - p_1^{(k)}} \right]. \quad (237)$$

$$\lambda_2(\mathbf{P}_k) = p_0^{(k)} + p_1^{(k)} - 1. \quad (238)$$

APPENDIX G EXPERIMENTAL EVALUATION

G1. Details on Experimental Setup

A. Datasets and Models: In this section, we provide a detailed description of the datasets and models used in our experiments. We considered a total of $N = 100$ clients. We tested CA-Fed on the benchmark synthetic LEAF dataset [36] for regularized logistic regression tasks, which satisfy Assumptions 3-4. Additionally, we incorporated two “real-world” datasets: MNIST [37] for handwritten digit recognition and CIFAR-10 [38] for image recognition. Detailed descriptions of the datasets and the models used for each of them are provided below.

a) Synthetic LEAF dataset: Synthetic data provides us with precise control over heterogeneity. The Synthetic LEAF dataset achieves this by using parameters γ and δ , where γ determines the degree of variation among local models and δ determines the variability in the local data across different devices. The generation process follows the setup described in [23], [24]:

- 1) For each client $k \in \mathcal{K}$, sample the model parameters $\mathbf{W}_k \in \mathbb{R}^{10 \times 60}$ and $\mathbf{b}_k \in \mathbb{R}^{10}$ from a normal distribution with mean μ_k and standard deviation 1, where μ_k is sampled from $\mathcal{N}(0, \gamma)$.
- 2) For each client $k \in \mathcal{K}$, generate the client’s input data $\mathbf{X}_k \in \mathbb{R}^{n_k \times 60}$ as follows: sample each element $(x_k)_j$ from a normal distribution with mean v_k and standard deviation $\frac{1}{j^{1.2}}$, where v_k is sampled from $\mathcal{N}(B_k, 1)$ and B_k is sampled from $\mathcal{N}(0, \delta)$.
- 3) Generate synthetic samples $(\mathbf{X}_k, \mathbf{Y}_k)$, where $\mathbf{Y}_k \in \mathbb{R}^{n_k}$, according to the model $y = \arg \max(\text{softmax}(\mathbf{W}_k \mathbf{x} + \mathbf{b}_k))$, where $\mathbf{x} \in \mathbb{R}^{60}$.

The distribution of samples $n_k = |D_k|$ among the clients follows a power law, resulting in an imbalanced data distribution. We refer to the synthetic dataset with parameters γ and δ as synthetic(γ, δ). We set (γ, δ) values to $(0, 0)$, $(0.25, 0.25)$, $(0.5, 0.5)$, $(0.75, 0.75)$, and $(1, 1)$ to investigate various levels of heterogeneity in the data.

TABLE I: Average computation time and used CPU/GPU for each dataset.

Dataset	CPU/GPU	Simulation time
Binary Synthetic	Intel(R) Xeon(R) CPU	10min
Synthetic LEAF	Intel(R) Xeon(R) CPU	6min
MNIST [37]	GeForce GTX 1080 Ti	42min
CIFAR10 [38]	GeForce GTX 1080 Ti	2h37min

TABLE II: Learning rates η and $\bar{\eta}$ used for the experiments in Figure 1.

Dataset	Unbiased	More available	CA-Fed ($\bar{\kappa} = 1$)	AdaFed [20]	F3AST [19]
Synthetic LEAF	2.0/2.0	1.0/7.0	2.0/3.0	1.0/1.0	2.0/2.0
MNIST	0.03/1.0	0.1/4.0	0.1/1.0	0.03/1.0	0.1/0.3
CIFAR10	0.03/1.0	0.03/3.0	0.03/1.0	0.03/1.0	0.03/0.3

b) MNIST: To classify handwritten digits in the MNIST dataset, we employ multinomial logistic regression. The model takes a flattened 784-dimensional (28×28) image as input and predicts a class label from 0 to 9 as output. To introduce heterogeneity in the data distribution, we distribute the dataset among $N = 100$ clients using a Dirichlet allocation method [39] with parameter ς . This allocation scheme allows for varying proportions of the dataset to be assigned to each client, contributing to the heterogeneous nature of our experimental setting.

c) CIFAR-10: The CIFAR-10 dataset consists of 60,000 input images, sourced from a collection of 80 million tiny images, with 10 distinct labels. To partition the CIFAR-10 dataset among $N = 100$ clients, we employ a Dirichlet allocation [39] with parameter ς . For this particular dataset, we train a shallow neural network comprising two convolutional layers followed by one fully connected layer. This network architecture is designed to capture relevant features from the CIFAR-10 images and facilitate accurate classification.

B. Implementation Details:

a) Machines: The experiments were conducted on a CPU/GPU cluster, utilizing various available GPUs such as Nvidia Tesla V100, GeForce GTX 1080 Ti, and Quadro RTX 8000. The majority of experiments involving Synthetic datasets were executed on an Intel(R) Xeon(R) CPU E5-1660 v3 @ 3.00GHz. On the other hand, experiments involving MNIST and CIFAR-10 datasets were performed using GeForce GTX 1080 Ti cards. For each dataset, we conducted approximately 50 experiments, excluding the time dedicated to development and debugging. Due to the usage of a train batch size of 32 samples, the experiments with MNIST and CIFAR-10 datasets exhibited slower execution times. Table I provides the average duration required to execute one simulation for each dataset. The authors are grateful to the OPAL infrastructure from Université Côte d’Azur for providing resources and support.

b) Libraries: We extensively employed the PyTorch deep learning framework throughout our experiments. PyTorch provided us with a comprehensive set of tools and functionalities for model construction, training, and evaluation. It allowed us to efficiently implement and optimize various neural network architectures, including the multinomial logistic regression model for the MNIST dataset and the shallow neural network for the CIFAR-10 dataset. To simplify the data preparation process, we utilized Torchvision, a PyTorch package designed for computer vision tasks. Torchvision facilitated seamless dataset management, including the download and pre-processing of MNIST and CIFAR-10, enabling us to transform the raw image data into a suitable format for training and evaluation.

c) Hyper-parameters: For each method and task, we performed a grid search to determine the optimal learning rates η and $\bar{\eta}$. For the MNIST and CIFAR-10 datasets, we explored the grids $\eta = \{2.0, 1.0, 0.3, 0.1, 0.03, 0.01\}$ and $\bar{\eta} = \{5.0, 4.0, 3.0, 2.0, 1.0, 0.3, 0.1\}$. For the Synthetic LEAF dataset, we shifted the grid to $\bar{\eta} = \{8.0, 7.0, 6.0, 5.0, 4.0, 3.0, 2.0, 1.0\}$. Table II reports the learning rates η and $\bar{\eta}$ corresponding to the results in Figure 1 for each dataset and method. For CA-Fed, we use the hyper-parameters $\beta = \tau = 0$. In the case of AdaFed, we set full device participation, where the parameter server samples all active clients ($|\mathcal{S}_t| = |\mathcal{A}_t|$). To ensure a fair comparison, we set the number of clients sampled by F3AST to the average number of clients included by CA-Fed, which is 45 on average. Furthermore, we set the smoothness parameter β of F3AST to be $\mathcal{O}(1/T)$, as suggested by the authors in [19, Appendix D].

APPENDIX H
FURTHER DISCUSSION ABOUT CA-FED

H1. CA-Fed's computation/communication cost

CA-Fed aims to improve training convergence and not to reduce its computation and communication overhead. Nevertheless, excluding some available clients reduces the overall training cost, as we will discuss in this section referring, for the sake of concreteness, to neural networks' training.

In terms of computation, the available clients not selected for training are only requested to evaluate their local loss on the current model once on a single batch instead than performing E gradient updates, which would require roughly $2 \times E - 1$ more calculations (because of the forward and backward pass). The selected clients have no extra computation cost as computing the loss corresponds to the forward pass they should, in any case, perform during the first local gradient update.

In terms of communication, the excluded clients only transmit the loss, a single scalar, much smaller than the model update. Conversely, participating clients transmit the local loss and the model update. Still, this additional overhead is negligible and likely fully compensated by the communication savings for the excluded clients.

H2. CA-Fed and Client Sampling

In cross-device FL, a common practice is to employ client sampling, where a small subset of clients (denoted as \mathcal{S}_t) is uniformly selected at random from the set of active clients (\mathcal{A}_t) during each communication round of model training. This is primarily done to mitigate communication overhead and enhance scalability.

In our analysis, based on Assumption 1, we assume that spatial and temporal correlations primarily concern clients' availability dynamics and we consider, for simplicity, $\mathcal{S}_t = \mathcal{A}_t$. However, our findings have a noteworthy implication: while the set of available clients \mathcal{A}_t exhibits correlation, the client sampling in \mathcal{S}_t can be designed to make clients' participation dynamics independent over time and among clients. A promising direction for future research is to extend our work in this context and derive a refined bound similar to our result in Theorem 2 which quantifies the impact of client sampling on $\lambda(\mathbf{P})$.

Consistent with our analysis, we have designed our algorithm to align with the assumption $\mathcal{S}_t = \mathcal{A}_t$. By design, CA-Fed excludes clients with large temporal correlation and low availability and activates, in each communication round, only clients satisfying $\{k \in \mathcal{A}_t; q_k^{(t)} > 0\}$ (line 8 in Algorithm 1). However, when only a small fraction of clients is excluded, CA-Fed seamlessly integrates with client sampling. This only involves replacing \mathcal{A}_t with \mathcal{S}_t in Equation (17) and Algorithm 1 (server estimates for clients' local losses $(\hat{\mathbf{F}}^{(t)} = (\hat{F}_k^{(t)})_{k \in \mathcal{K}})$ are now updated from the sampled clients' losses $(\mathbf{F}^{(t)} = (F_k^{(t)})_{k \in \mathcal{S}_t})$).

H3. About CA-Fed's fairness

Strategies that exclude clients from the training phase, such as CA-Fed, may raise concerns about fairness. The concept of *fairness* in federated learning does not have a unified definition in the literature [42, Chapter 8]. Fairness goals can be established by appropriately selecting the target weights $\alpha = \{\alpha_k\}_{k \in \mathcal{K}}$ in the definition of the global target objective (1). For instance, *per-client fairness* can be achieved by setting α_k to be equal for every client (i.e., $\alpha_k = 1/N$), while *per-sample fairness* can be accomplished by setting α_k proportional to the local dataset size $|D_k|$ (i.e., $\alpha_k = |D_k|/|D|$).

Assuming that the global objective in (1) truly reflects fairness concerns, then CA-Fed can be considered intrinsically fair. This is because CA-Fed continually focuses on minimizing the total error $\epsilon := F(\mathbf{w}_T) - F^*$, which guarantees that the performance objective of the learned model is as close as possible to its optimal value at every time. Although CA-Fed occasionally excludes clients with low availability and high temporal correlation, the optimization problem (1) is carefully designed to ensure that the learned model performs well for these clients. As a result, CA-Fed effectively learns a model that is consistently accurate and fair across all clients, regardless of their availability or temporal correlation.