

A Proof of Lemma 1: Objective Inconsistency in Quadratic Model

Formulation. Consider a simple setting where each local objective function is strongly convex and defined as follows:

$$F_i(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H}_i \mathbf{x} - \mathbf{e}_i^\top \mathbf{x} + \frac{1}{2} \mathbf{e}_i^\top \mathbf{H}_i^{-1} \mathbf{e}_i \geq 0 \quad (13)$$

where $\mathbf{H}_i \in \mathbb{R}^{d \times d}$ is an invertible matrix and $\mathbf{e}_i \in \mathbb{R}^d$ is an arbitrary vector. It is easy to show that the optimum of the i -th local function is $\mathbf{x}_i^* = \mathbf{H}_i^{-1} \mathbf{e}_i$. Without loss of generality, we assume the global objective function to be a weighted average across all local functions, that is:

$$F(\mathbf{x}) = \sum_{i=1}^m p_i F_i(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \overline{\mathbf{H}} \mathbf{x} - \overline{\mathbf{e}}^\top \mathbf{x} + \frac{1}{2} \sum_{i=1}^m p_i \mathbf{e}_i^\top \mathbf{h}_i^{-1} \mathbf{e}_i \quad (14)$$

where $\overline{\mathbf{H}} = \sum_{i=1}^m p_i \mathbf{H}_i$ and $\overline{\mathbf{e}} = \sum_{i=1}^m p_i \mathbf{e}_i$. As a result, the global minimum is $\mathbf{x}^* = \overline{\mathbf{H}}^{-1} \overline{\mathbf{e}}$. Now, let us study whether previous federated optimization algorithms can converge to this global minimum.

Local Update Rule. The local update rule of FedProx for the i -th device can be written as follows:

$$\mathbf{x}_i^{(t,k+1)} = \mathbf{x}_i^{(t,k)} - \eta \left[\mathbf{H}_i \mathbf{x}_i^{(t,k)} - \mathbf{e}_i + \mu (\mathbf{x}_i^{(t,k)} - \mathbf{x}^{(t,0)}) \right] \quad (15)$$

$$= (\mathbf{I} - \eta \mu \mathbf{I} - \eta \mathbf{H}_i) \mathbf{x}_i^{(t,k)} + \eta \mathbf{e}_i + \eta \mu \mathbf{x}^{(t,0)} \quad (16)$$

where $\mathbf{x}_i^{(t,k)}$ denotes the local model parameters at the k -th local iteration after t communication rounds, η denotes the local learning rate and μ is a tunable hyper-parameter in FedProx. When $\mu = 0$, the algorithm will reduce to FedAvg. We omit the device index in $\mathbf{x}^{(t,0)}$, since it is synchronized and the same across all devices.

After minor arranging (16), we obtain

$$\mathbf{x}_i^{(t,k+1)} - \mathbf{c}_i^{(t)} = (\mathbf{I} - \eta \mu \mathbf{I} - \eta \mathbf{H}_i) \left(\mathbf{x}_i^{(t,k)} - \mathbf{c}_i^{(t)} \right). \quad (17)$$

where $\mathbf{c}_i^{(t)} = (\mathbf{H}_i + \mu \mathbf{I})^{-1} (\mathbf{e}_i + \mu \mathbf{x}^{(t,0)})$. Then, after performing τ_i steps of local updates, the local model becomes

$$\mathbf{x}_i^{(t,\tau_i)} = (\mathbf{I} - \eta \mu \mathbf{I} - \eta \mathbf{H}_i)^{\tau_i} \left(\mathbf{x}^{(t,0)} - \mathbf{c}_i^{(t)} \right) + \mathbf{c}_i^{(t)}, \quad (18)$$

$$\mathbf{x}_i^{(t,\tau_i)} - \mathbf{x}^{(t,0)} = (\mathbf{I} - \eta \mu \mathbf{I} - \eta \mathbf{H}_i)^{\tau_i} \left(\mathbf{x}^{(t,0)} - \mathbf{c}_i^{(t)} \right) + \mathbf{c}_i^{(t)} - \mathbf{x}^{(t,0)} \quad (19)$$

$$= [(\mathbf{I} - \eta \mu \mathbf{I} - \eta \mathbf{H}_i)^{\tau_i} - \mathbf{I}] \left(\mathbf{x}^{(t,0)} - \mathbf{c}_i^{(t)} \right) \quad (20)$$

$$= [\mathbf{I} - (\mathbf{I} - \eta \mu \mathbf{I} - \eta \mathbf{H}_i)^{\tau_i}] (\mathbf{H}_i + \mu \mathbf{I})^{-1} \left(\mathbf{e}_i - \mathbf{H}_i \mathbf{x}^{(t,0)} \right). \quad (21)$$

For the ease of writing, we define $\mathbf{K}_i(\eta, \mu) = [\mathbf{I} - (\mathbf{I} - \eta \mu \mathbf{I} - \eta \mathbf{H}_i)^{\tau_i}] (\mathbf{H}_i + \mu \mathbf{I})^{-1}$.

Server Aggregation. For simplicity, we only consider the case when all devices participate in the each round. In FedProx, the server averages all local models according to the sample size:

$$\mathbf{x}^{(t+1,0)} - \mathbf{x}^{(t,0)} = \sum_{i=1}^m p_i \left(\mathbf{x}_i^{(t,\tau_i)} - \mathbf{x}^{(t,0)} \right) \quad (22)$$

$$= \sum_{i=1}^m p_i \mathbf{K}_i(\eta, \mu) \left(\mathbf{e}_i - \mathbf{H}_i \mathbf{x}^{(t,0)} \right). \quad (23)$$

Accordingly, we get the following update rule for the central model:

$$\mathbf{x}^{(t+1,0)} = \left[\mathbf{I} - \sum_{i=1}^m p_i \mathbf{K}_i(\eta, \mu) \mathbf{H}_i \right] \mathbf{x}^{(t,0)} + \sum_{i=1}^m p_i \mathbf{K}_i(\eta, \mu) \mathbf{e}_i. \quad (24)$$

It is equivalent to

$$\mathbf{x}^{(t+1,0)} - \tilde{\mathbf{x}} = \left[\mathbf{I} - \sum_{i=1}^m p_i \mathbf{K}_i(\eta, \mu) \mathbf{H}_i \right] \left[\mathbf{x}^{(t,0)} - \tilde{\mathbf{x}} \right]. \quad (25)$$

where

$$\tilde{\mathbf{x}} = \left(\sum_{i=1}^m p_i \mathbf{K}_i(\eta, \mu) \mathbf{H}_i \right)^{-1} \left(\sum_{i=1}^m p_i \mathbf{K}_i(\eta, \mu) \mathbf{e}_i \right). \quad (26)$$

After T communication rounds, one can get

$$\mathbf{x}^{(T,0)} = \left[\mathbf{I} - \sum_{i=1}^m p_i \mathbf{K}_i(\eta, \mu) \mathbf{H}_i \right]^T \left[\mathbf{x}^{(t,0)} - \tilde{\mathbf{x}} \right] + \tilde{\mathbf{x}}. \quad (27)$$

Accordingly, when $\|\mathbf{I} - \sum_{i=1}^m p_i \mathbf{K}_i(\eta, \mu) \mathbf{H}_i\|_2 < 1$, the iterates will converge to

$$\lim_{T \rightarrow \infty} \mathbf{x}^{(T,0)} = \tilde{\mathbf{x}} = \left(\sum_{i=1}^m p_i \mathbf{K}_i(\eta, \mu) \mathbf{H}_i \right)^{-1} \left(\sum_{i=1}^m p_i \mathbf{K}_i(\eta, \mu) \mathbf{e}_i \right). \quad (28)$$

Recall that $\mathbf{K}_i(\eta, \mu) = [\mathbf{I} - (\mathbf{I} - \eta\mu\mathbf{I} - \eta\mathbf{H}_i)^{\tau_i}] (\mathbf{H}_i + \mu\mathbf{I})^{-1}$.

Concrete Example in Lemma 1. Now let us focus on a concrete example where $p_1 = p_2 = \dots = p_m = 1/m$, $\mathbf{H}_1 = \mathbf{H}_2 = \dots = \mathbf{H}_m = \mathbf{I}$ and $\mu = 0$. Then, in this case, $\mathbf{K}_i = 1 - (1 - \eta)^{\tau_i}$. As a result, we have

$$\lim_{T \rightarrow \infty} \mathbf{x}^{(T,0)} = \frac{\sum_{i=1}^m [1 - (1 - \eta)^{\tau_i}] \mathbf{e}_i}{\sum_{i=1}^m [1 - (1 - \eta)^{\tau_i}]}. \quad (29)$$

Furthermore, when the learning rate is sufficiently small (*e.g.*, can be achieved by gradually decaying the learning rate), according to L'Hospital's rule, we obtain

$$\lim_{\eta \rightarrow 0} \lim_{T \rightarrow \infty} \mathbf{x}^{(T,0)} = \frac{\sum_{i=1}^m \tau_i \mathbf{e}_i}{\sum_{i=1}^m \tau_i}. \quad (30)$$

Here, we complete the proof of Lemma 1.

B Detailed Derivations for Various Local Solvers

In this section, we will derive the specific expression of the vector \mathbf{a}_i when using different local solvers. Recall that the local change at client i is $\Delta_i^{(t)} = -\eta \mathbf{G}_i^{(t)} \mathbf{a}_i$ where $\mathbf{G}_i^{(t)}$ stacks all stochastic gradients in the current round and \mathbf{a} is a non-negative vector.

B.1 SGD with Proximal Updates

In this case, we can write the update rule of local models as follows:

$$\mathbf{x}_i^{(t, \tau_i)} = \mathbf{x}_i^{(t, \tau_i - 1)} - \eta \left[g_i(\mathbf{x}_i^{(t, \tau_i - 1)}) + \mu \left(\mathbf{x}_i^{(t, \tau_i - 1)} - \mathbf{x}^{(t, 0)} \right) \right]. \quad (31)$$

Subtracting $\mathbf{x}_i^{(t, 0)}$ on both sides, we obtain

$$\mathbf{x}_i^{(t, \tau_i)} - \mathbf{x}^{(t, 0)} = \mathbf{x}_i^{(t, \tau_i - 1)} - \mathbf{x}^{(t, 0)} - \eta \left[g_i(\mathbf{x}_i^{(t, \tau_i - 1)}) + \mu \left(\mathbf{x}_i^{(t, \tau_i - 1)} - \mathbf{x}^{(t, 0)} \right) \right] \quad (32)$$

$$= (1 - \eta\mu) \left(\mathbf{x}_i^{(t, \tau_i - 1)} - \mathbf{x}^{(t, 0)} \right) - \eta g_i(\mathbf{x}_i^{(t, \tau_i - 1)}). \quad (33)$$

Repeating the above procedure, it follows that

$$\Delta_i^{(t)} = \mathbf{x}_i^{(t, \tau_i)} - \mathbf{x}^{(t, 0)} = -\eta \sum_{k=0}^{\tau_i - 1} (1 - \eta\mu)^{\tau_i - 1 - k} g_i(\mathbf{x}_i^{(t, k)}). \quad (34)$$

According to the definition, we have $\mathbf{a}_i = [(1 - \alpha)^{\tau_i - 1}, (1 - \alpha)^{\tau_i - 2}, \dots, (1 - \alpha), 1]$ where $\alpha = \eta\mu$.

B.2 SGD with Local Momentum

Let us firstly write down the update rule of the local models. Suppose that ρ denotes the local momentum factor and \mathbf{u}_i is the local momentum buffer at client i . Then, the update rule of local momentum SGD is:

$$\mathbf{u}_i^{(t, \tau_i)} = \rho \mathbf{u}_i^{(t, \tau_i - 1)} + g_i(\mathbf{x}_i^{(t, \tau_i - 1)}), \quad (35)$$

$$\mathbf{x}_i^{(t, \tau_i)} = \mathbf{x}_i^{(t, \tau_i - 1)} - \eta \mathbf{u}_i^{(t, \tau_i)}. \quad (36)$$

One can expand the expression of local momentum buffer as follows:

$$\mathbf{u}_i^{(t, \tau_i)} = \rho \mathbf{u}_i^{(t, \tau_i - 1)} + g_i(\mathbf{x}_i^{(t, \tau_i - 1)}) \quad (37)$$

$$= \rho^2 \mathbf{u}_i^{(t, \tau_i - 2)} + \rho g_i(\mathbf{x}_i^{(t, \tau_i - 2)}) + g_i(\mathbf{x}_i^{(t, \tau_i - 1)}) \quad (38)$$

$$= \sum_{k=0}^{\tau_i - 1} \rho^{\tau_i - 1 - k} g_i(\mathbf{x}_i^{(t, k)}) \quad (39)$$

where the last equation comes from the fact $\mathbf{u}_i^{(t, 0)} = 0$. Substituting (39) into (36), we have

$$\mathbf{x}_i^{(t, \tau_i)} = \mathbf{x}_i^{(t, \tau_i - 1)} - \eta \sum_{k=0}^{\tau_i - 1} \rho^{\tau_i - 1 - k} g_i(\mathbf{x}_i^{(t, k)}) \quad (40)$$

$$= \mathbf{x}_i^{(t, \tau_i - 2)} - \eta \sum_{k=0}^{\tau_i - 2} \rho^{\tau_i - 2 - k} g_i(\mathbf{x}_i^{(t, k)}) - \eta \sum_{k=0}^{\tau_i - 1} \rho^{\tau_i - 1 - k} g_i(\mathbf{x}_i^{(t, k)}). \quad (41)$$

Repeating the above procedure, it follows that

$$\mathbf{x}_i^{(t, \tau_i)} - \mathbf{x}_i^{(t, 0)} = -\eta \sum_{s=0}^{\tau_i - 1} \sum_{k=0}^s \rho^{s-k} g_i(\mathbf{x}_i^{(t, k)}) \quad (42)$$

Then, the coefficient of $g_i(\mathbf{x}_i^{(t, k)})$ is

$$\sum_{s \geq k}^{\tau_i - 1} \rho^{s-k} = 1 + \rho + \rho^2 + \dots + \rho^{\tau_i - 1 - k} = \frac{1 - \rho^{\tau_i - k}}{1 - \rho}. \quad (43)$$

That is, $\mathbf{a}_i = [1 - \rho^{\tau_i}, 1 - \rho^{\tau_i - 1}, \dots, 1 - \rho] / (1 - \rho)$. In this case, the ℓ_1 norm of \mathbf{a}_i is

$$\|\mathbf{a}_i\|_1 = \frac{1}{1 - \rho} \sum_{k=0}^{\tau_i - 1} (1 - \rho^{\tau_i - k}) = \frac{1}{1 - \rho} \left(\tau_i - \sum_{k=0}^{\tau_i - 1} \rho^{\tau_i - k} \right) \quad (44)$$

$$= \frac{1}{1 - \rho} \left[\tau_i - \frac{\rho(1 - \rho^{\tau_i})}{1 - \rho} \right]. \quad (45)$$

C Proof of Theorem 1: Convergence of Surrogate Objective

C.1 Preliminaries

For the ease of writing, let us define a surrogate objective function $\tilde{F}(\mathbf{x}) = \sum_{i=1}^m w_i F_i(\mathbf{x})$, where $\sum_{i=1}^m w_i = 1$, and define the following auxiliary variables

$$\text{Normalized Stochastic Gradient: } \mathbf{d}_i^{(t)} = \frac{1}{a_i} \sum_{k=0}^{\tau_i-1} a_i^{(k)} g_i(\mathbf{x}_i^{(t,k)}), \quad (46)$$

$$\text{Normalized Gradient: } \mathbf{h}_i^{(t)} = \frac{1}{a_i} \sum_{k=0}^{\tau_i-1} a_i^{(k)} \nabla F_i(\mathbf{x}_i^{(t,k)}) \quad (47)$$

where $a_i^{(k)} \geq 0$ is an arbitrary scalar, $\mathbf{a}_i = [a_i^{(0)}, \dots, a_i^{(\tau_i-1)}]^\top$, and $a_i = \|\mathbf{a}_i\|_1$. Besides, one can show that $\mathbb{E}[\mathbf{d}_i^{(t)} - \mathbf{h}_i^{(t)}] = \mathbf{0}$. In addition, since workers are independent to each other, we have $\mathbb{E} \langle \mathbf{d}_i^{(t)} - \mathbf{h}_i^{(t)}, \mathbf{d}_j^{(t)} - \mathbf{h}_j^{(t)} \rangle = 0, \forall i \neq j$. Recall that the update rule of the global model can be written as follows:

$$\mathbf{x}^{(t+1,0)} - \mathbf{x}^{(t,0)} = -\tau_{\text{eff}} \eta \sum_{i=1}^m w_i \mathbf{d}_i^{(t)}. \quad (48)$$

According to the Lipschitz-smooth assumption, it follows that

$$\begin{aligned} & \mathbb{E} \left[\tilde{F}(\mathbf{x}^{(t+1,0)}) \right] - \tilde{F}(\mathbf{x}^{(t,0)}) \\ & \leq \underbrace{-\tau_{\text{eff}} \eta \mathbb{E} \left[\left\langle \nabla \tilde{F}(\mathbf{x}^{(t,0)}), \sum_{i=1}^m w_i \mathbf{d}_i^{(t)} \right\rangle \right]}_{T_1} + \underbrace{\frac{\tau_{\text{eff}}^2 \eta^2 L}{2} \mathbb{E} \left[\left\| \sum_{i=1}^m w_i \mathbf{d}_i^{(t)} \right\|_F^2 \right]}_{T_2} \end{aligned} \quad (49)$$

where the expectation is taken over mini-batches $\xi_i^{(t,k)}, \forall i \in \{1, 2, \dots, m\}, k \in \{0, 1, \dots, \tau_i - 1\}$. Before diving into the detailed bounds for T_1 and T_2 , we would like to firstly introduce several useful lemmas.

Lemma 2. Suppose $\{A_k\}_{k=1}^T$ is a sequence of random matrices and $\mathbb{E}[A_k | A_{k-1}, A_{k-2}, \dots, A_1] = \mathbf{0}, \forall k$. Then,

$$\mathbb{E} \left[\left\| \sum_{k=1}^T A_k \right\|_F^2 \right] = \sum_{k=1}^T \mathbb{E} \left[\|A_k\|_F^2 \right]. \quad (50)$$

Proof.

$$\mathbb{E} \left[\left\| \sum_{k=1}^T A_k \right\|_F^2 \right] = \sum_{k=1}^T \mathbb{E} \left[\|A_k\|_F^2 \right] + \sum_{i=1}^T \sum_{j=1, j \neq i}^T \mathbb{E} \left[\text{Tr} \{ A_i^\top A_j \} \right] \quad (51)$$

$$= \sum_{k=1}^T \mathbb{E} \left[\|A_k\|_F^2 \right] + \sum_{i=1}^T \sum_{j=1, j \neq i}^T \text{Tr} \{ \mathbb{E} [A_i^\top A_j] \} \quad (52)$$

Assume $i < j$. Then, using the law of total expectation,

$$\mathbb{E} [A_i^\top A_j] = \mathbb{E} [A_i^\top \mathbb{E} [A_j | A_i, \dots, A_1]] = \mathbf{0}. \quad (53)$$

□

C.2 Bounding First term in (49)

For the first term on the right hand side (RHS) in (49), we have

$$T_1 = \mathbb{E} \left[\left\langle \nabla \tilde{F}(\mathbf{x}^{(t,0)}), \sum_{i=1}^m w_i (\mathbf{d}_i^{(t)} - \mathbf{h}_i^{(t)}) \right\rangle \right] + \mathbb{E} \left[\left\langle \nabla \tilde{F}(\mathbf{x}^{(t,0)}), \sum_{i=1}^m w_i \mathbf{h}_i^{(t)} \right\rangle \right] \quad (54)$$

$$= \mathbb{E} \left[\left\langle \nabla \tilde{F}(\mathbf{x}^{(t,0)}), \sum_{i=1}^m w_i \mathbf{h}_i^{(t)} \right\rangle \right] \quad (55)$$

$$= \frac{1}{2} \left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 + \frac{1}{2} \mathbb{E} \left[\left\| \sum_{i=1}^m w_i \mathbf{h}_i^{(t)} \right\|^2 \right] - \frac{1}{2} \mathbb{E} \left[\left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) - \sum_{i=1}^m w_i \mathbf{h}_i^{(t)} \right\|^2 \right] \quad (56)$$

where the last equation uses the fact: $2 \langle a, b \rangle = \|a\|^2 + \|b\|^2 - \|a - b\|^2$.

C.3 Bounding Second term in (49)

For the second term on the right hand side (RHS) in (49), we have

$$T_2 = \mathbb{E} \left[\left\| \sum_{i=1}^m w_i (\mathbf{d}_i^{(t)} - \mathbf{h}_i^{(t)}) + \sum_{i=1}^m w_i \mathbf{h}_i^{(t)} \right\|^2 \right] \quad (57)$$

$$\leq 2 \mathbb{E} \left[\left\| \sum_{i=1}^m w_i (\mathbf{d}_i^{(t)} - \mathbf{h}_i^{(t)}) \right\|^2 \right] + 2 \mathbb{E} \left[\left\| \sum_{i=1}^m w_i \mathbf{h}_i^{(t)} \right\|^2 \right] \quad (58)$$

$$= 2 \sum_{i=1}^m w_i^2 \mathbb{E} \left[\left\| \mathbf{d}_i^{(t)} - \mathbf{h}_i^{(t)} \right\|^2 \right] + 2 \mathbb{E} \left[\left\| \sum_{i=1}^m w_i \mathbf{h}_i^{(t)} \right\|^2 \right] \quad (59)$$

where (58) follows the fact: $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$ and (59) uses the special property of $\mathbf{d}_i^{(t)}, \mathbf{h}_i^{(t)}$, that is, $\mathbb{E} \langle \mathbf{d}_i^{(t)} - \mathbf{h}_i^{(t)}, \mathbf{d}_j^{(t)} - \mathbf{h}_j^{(t)} \rangle = 0, \forall i \neq j$. Then, let us expand the expression of $\mathbf{d}_i^{(t)}$ and $\mathbf{h}_i^{(t)}$, we obtain that

$$T_2 \leq \sum_{i=1}^m \frac{2w_i^2}{a_i^2} \sum_{k=0}^{\tau_i-1} [a_i^{(k)}]^2 \mathbb{E} \left[\left\| g_i(\mathbf{x}_i^{(t,k)}) - \nabla F_i(\mathbf{x}_i^{(t,k)}) \right\|^2 \right] + 2 \mathbb{E} \left[\left\| \sum_{i=1}^m w_i \mathbf{h}_i^{(t)} \right\|^2 \right] \quad (60)$$

$$\leq 2\sigma^2 \sum_{i=1}^m \frac{w_i^2 \|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} + 2 \mathbb{E} \left[\left\| \sum_{i=1}^m w_i \mathbf{h}_i^{(t)} \right\|^2 \right] \quad (61)$$

where (60) is derived using Lemma 2, (61) follows Assumption 2.

C.4 Intermediate Result

Plugging (56) and (61) back into (49), we have

$$\begin{aligned} \mathbb{E} \left[\tilde{F}(\mathbf{x}^{(t+1,0)}) \right] - \tilde{F}(\mathbf{x}^{(t,0)}) &\leq -\frac{\tau_{\text{eff}} \eta}{2} \left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 - \frac{\tau_{\text{eff}} \eta}{2} (1 - 2\tau_{\text{eff}} \eta L) \mathbb{E} \left[\left\| \sum_{i=1}^m w_i \mathbf{h}_i^{(t)} \right\|^2 \right] \\ &\quad + \tau_{\text{eff}}^2 \eta^2 L \sigma^2 \sum_{i=1}^m \frac{w_i^2 \|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} + \frac{\tau_{\text{eff}} \eta}{2} \mathbb{E} \left[\left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) - \sum_{i=1}^m w_i \mathbf{h}_i^{(t)} \right\|^2 \right] \end{aligned} \quad (62)$$

When $\tau_{\text{eff}}\eta L \leq 1/2$, it follows that

$$\begin{aligned} \frac{\mathbb{E} [\tilde{F}(\mathbf{x}^{(t+1,0)})] - \tilde{F}(\mathbf{x}^{(t,0)})}{\eta\tau_{\text{eff}}} &\leq -\frac{1}{2} \left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 + \tau_{\text{eff}}\eta L\sigma^2 \sum_{i=1}^m \frac{w_i^2 \|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} \\ &\quad + \frac{1}{2} \mathbb{E} \left[\left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) - \sum_{i=1}^m w_i \mathbf{h}_i^{(t)} \right\|^2 \right] \end{aligned} \quad (63)$$

$$\begin{aligned} &\leq -\frac{1}{2} \left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 + \tau_{\text{eff}}\eta L\sigma^2 \sum_{i=1}^m \frac{w_i^2 \|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} \\ &\quad + \frac{1}{2} \sum_{i=1}^m w_i \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) - \mathbf{h}_i^{(t)} \right\|^2 \right] \end{aligned} \quad (64)$$

where the last inequality uses the fact $\tilde{F}(\mathbf{x}) = \sum_{i=1}^m w_i F_i(\mathbf{x})$ and Jensen's Inequality: $\left\| \sum_{i=1}^m w_i z_i \right\|^2 \leq \sum_{i=1}^m w_i \|z_i\|^2$. Next, we will focus on bounding the last term in (64).

C.5 Bounding the Difference Between Server Gradient and Normalized Gradient

Recall the definition of $\mathbf{h}_i^{(t)}$, one can derive that

$$\mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) - \mathbf{h}_i^{(t)} \right\|^2 \right] = \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) - \frac{1}{a_i} \sum_{k=0}^{\tau_i-1} a_i^{(k)} \nabla F_i(\mathbf{x}_i^{(t,k)}) \right\|^2 \right] \quad (65)$$

$$= \mathbb{E} \left[\left\| \frac{1}{a_i} \sum_{k=0}^{\tau_i-1} a_i^{(k)} \left(\nabla F_i(\mathbf{x}^{(t,0)}) - \nabla F_i(\mathbf{x}_i^{(t,k)}) \right) \right\|^2 \right] \quad (66)$$

$$\leq \frac{1}{a_i} \sum_{k=0}^{\tau_i-1} \left\{ a_i^{(k)} \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) - \nabla F_i(\mathbf{x}_i^{(t,k)}) \right\|^2 \right] \right\} \quad (67)$$

$$\leq \frac{L^2}{a_i} \sum_{k=0}^{\tau_i-1} \left\{ a_i^{(k)} \mathbb{E} \left[\left\| \mathbf{x}^{(t,0)} - \mathbf{x}_i^{(t,k)} \right\|^2 \right] \right\} \quad (68)$$

where (67) uses Jensen's Inequality again: $\left\| \sum_{i=1}^m w_i z_i \right\|^2 \leq \sum_{i=1}^m w_i \|z_i\|^2$, and (68) follows Assumption 1. Now, we turn to bounding the difference between the server model $\mathbf{x}^{(t,0)}$ and the local model $\mathbf{x}_i^{(t,k)}$. Plugging into the local update rule and using the fact $\|a+b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$,

$$\mathbb{E} \left[\left\| \mathbf{x}^{(t,0)} - \mathbf{x}_i^{(t,k)} \right\|^2 \right] = \eta^2 \cdot \mathbb{E} \left[\left\| \sum_{s=0}^{k-1} a_i^{(s)} g_i(\mathbf{x}_i^{(t,s)}) \right\|^2 \right] \quad (69)$$

$$\leq 2\eta^2 \mathbb{E} \left[\left\| \sum_{s=0}^{k-1} a_i^{(s)} \left(g_i(\mathbf{x}_i^{(t,s)}) - \nabla F_i(\mathbf{x}_i^{(t,s)}) \right) \right\|^2 \right] \quad (70)$$

$$+ 2\eta^2 \mathbb{E} \left[\left\| \sum_{s=0}^{k-1} a_i^{(s)} \nabla F_i(\mathbf{x}_i^{(t,s)}) \right\|^2 \right] \quad (71)$$

Applying Lemma 2 to the first term,

$$\begin{aligned} \mathbb{E} \left[\left\| \mathbf{x}^{(t,0)} - \mathbf{x}_i^{(t,k)} \right\|^2 \right] &= 2\eta^2 \sum_{s=0}^{k-1} [a_i^{(s)}]^2 \mathbb{E} \left[\left\| g_i(\mathbf{x}_i^{(t,s)}) - \nabla F_i(\mathbf{x}_i^{(t,s)}) \right\|^2 \right] \\ &\quad + 2\eta^2 \mathbb{E} \left[\left\| \sum_{s=0}^{k-1} a_i^{(s)} \nabla F_i(\mathbf{x}_i^{(t,s)}) \right\|^2 \right] \end{aligned} \quad (72)$$

$$\leq 2\eta^2 \sigma^2 \sum_{s=0}^{k-1} [a_i^{(s)}]^2 + 2\eta^2 \mathbb{E} \left[\left\| \sum_{s=0}^{k-1} a_i^{(s)} \nabla F_i(\mathbf{x}_i^{(t,s)}) \right\|^2 \right] \quad (73)$$

$$\leq 2\eta^2 \sigma^2 \sum_{s=0}^{k-1} [a_i^{(s)}]^2 + 2\eta^2 \left[\sum_{s=0}^{k-1} a_i^{(s)} \right] \sum_{s=0}^{k-1} a_i^{(s)} \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}_i^{(t,s)}) \right\|^2 \right] \quad (74)$$

$$\leq 2\eta^2 \sigma^2 \sum_{s=0}^{k-1} [a_i^{(s)}]^2 + 2\eta^2 \left[\sum_{s=0}^{k-1} a_i^{(s)} \right] \sum_{s=0}^{\tau_i-1} a_i^{(s)} \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}_i^{(t,s)}) \right\|^2 \right] \quad (75)$$

where (74) follows from Jensen's Inequality. Furthermore, note that

$$\frac{1}{\|\mathbf{a}_i\|_1} \sum_{k=0}^{\tau_i-1} a_i^{(k)} \left[\sum_{s=0}^{k-1} [a_i^{(s)}]^2 \right] \leq \frac{1}{a_i} \sum_{k=0}^{\tau_i-1} a_i^{(k)} \left[\sum_{s=0}^{\tau_i-2} [a_i^{(s)}]^2 \right] \quad (76)$$

$$= \sum_{s=0}^{\tau_i-2} [a_i^{(s)}]^2 = \|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2, \quad (77)$$

$$\frac{1}{\|\mathbf{a}_i\|_1} \sum_{k=0}^{\tau_i-1} a_i^{(k)} \left[\sum_{s=0}^{k-1} [a_i^{(s)}] \right] \leq \frac{1}{a_i} \sum_{k=0}^{\tau_i-1} a_i^{(k)} \left[\sum_{s=0}^{\tau_i-2} [a_i^{(s)}] \right] \quad (78)$$

$$= \sum_{s=0}^{\tau_i-2} [a_i^{(s)}] = \|\mathbf{a}_i\|_1 - a_{i,-1} \quad (79)$$

where $a_{i,-1}$ is the last element in the vector \mathbf{a}_i . As a result, we have

$$\begin{aligned} \frac{1}{\|\mathbf{a}_i\|_1} \sum_{k=0}^{\tau_i-1} a_i^{(k)} \mathbb{E} \left[\left\| \mathbf{x}^{(t,0)} - \mathbf{x}_i^{(t,k)} \right\|^2 \right] &\leq 2\eta^2 \sigma^2 \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2 \right) \\ &\quad + 2\eta^2 \left(\|\mathbf{a}_i\|_1 - a_{i,-1} \right) \sum_{k=0}^{\tau_i-1} a_i^{(s)} \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}_i^{(t,k)}) \right\|^2 \right] \end{aligned} \quad (80)$$

In addition, we can bound the second term using the following inequality:

$$\mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}_i^{(t,k)}) \right\|^2 \right] \leq 2\mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}_i^{(t,k)}) - \nabla F_i(\mathbf{x}^{(t,0)}) \right\|^2 \right] + 2\mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) \right\|^2 \right] \quad (81)$$

$$\leq 2L^2 \mathbb{E} \left[\left\| \mathbf{x}^{(t,0)} - \mathbf{x}_i^{(t,k)} \right\|^2 \right] + 2\mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) \right\|^2 \right]. \quad (82)$$

Substituting (82) into (75), we get

$$\begin{aligned} &\frac{1}{\|\mathbf{a}_i\|_1} \sum_{k=0}^{\tau_i-1} a_i^{(k)} \mathbb{E} \left[\left\| \mathbf{x}^{(t,0)} - \mathbf{x}_i^{(t,k)} \right\|^2 \right] \\ &\leq 2\eta^2 \sigma^2 \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2 \right) + 4\eta^2 L^2 \left(\|\mathbf{a}_i\|_1 - a_{i,-1} \right) \sum_{k=0}^{\tau_i-1} a_i^{(k)} \mathbb{E} \left[\left\| \mathbf{x}^{(t,0)} - \mathbf{x}_i^{(t,k)} \right\|^2 \right] \\ &\quad + 4\eta^2 \left(\|\mathbf{a}_i\|_1 - a_{i,-1} \right) \sum_{k=0}^{\tau_i-1} a_i^{(k)} \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}_i^{(t,0)}) \right\|^2 \right] \end{aligned} \quad (83)$$

After minor rearranging, it follows that

$$\begin{aligned} \frac{1}{\|\mathbf{a}_i\|_1} \sum_{k=0}^{\tau_i-1} a_i^{(k)} \mathbb{E} \left[\left\| \mathbf{x}^{(t,0)} - \mathbf{x}_i^{(t,k)} \right\|^2 \right] &\leq \frac{2\eta^2 \sigma^2}{1 - 4\eta^2 L^2 \|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1})} \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2 \right) \\ &\quad + \frac{4\eta^2 \|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1})}{1 - 4\eta^2 L^2 \|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1})} \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) \right\|^2 \right] \end{aligned} \quad (84)$$

Define $D = 4\eta^2 L^2 \max_i \{ \|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1}) \} < 1$. We can simplify (84) as follows

$$\frac{L^2}{a_i} \sum_{k=0}^{\tau_i-1} a_i^{(k)} \mathbb{E} \left[\left\| \mathbf{x}^{(t,0)} - \mathbf{x}_i^{(t,k)} \right\|^2 \right] \leq \frac{2\eta^2 L^2 \sigma^2}{1 - D} \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2 \right) + \frac{D}{1 - D} \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) \right\|^2 \right]. \quad (85)$$

Taking the average across all workers and applying Assumption 3, one can obtain

$$\begin{aligned} \frac{1}{2} \sum_{i=1}^m w_i \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) - \mathbf{h}_i^{(t)} \right\|^2 \right] &\leq \frac{\eta^2 L^2 \sigma^2}{1 - D} \sum_{i=1}^m w_i \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2 \right) \\ &\quad + \frac{D}{2(1 - D)} \sum_{i=1}^m w_i \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) \right\|^2 \right] \quad (86) \\ &\leq \frac{\eta^2 L^2 \sigma^2}{1 - D} \sum_{i=1}^m w_i \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2 \right) \\ &\quad + \frac{D\beta^2}{2(1 - D)} \mathbb{E} \left[\left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 \right] + \frac{D\kappa^2}{2(1 - D)}. \quad (87) \end{aligned}$$

Now, we are ready to derive the final result.

C.6 Final Results

Plugging (87) back into (64), we have

$$\begin{aligned} \frac{\mathbb{E} [\tilde{F}(\mathbf{x}^{(t+1,0)})] - \tilde{F}(\mathbf{x}^{(t,0)})}{\eta\tau_{\text{eff}}} &\leq -\frac{1}{2} \left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 + \tau_{\text{eff}} \eta L \sigma^2 \sum_{i=1}^m \frac{w_i^2 \|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} \\ &\quad + \frac{\eta^2 L^2 \sigma^2}{1 - D} \sum_{i=1}^m w_i \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2 \right) \\ &\quad + \frac{D\kappa^2}{2(1 - D)} + \frac{D\beta^2}{2(1 - D)} \mathbb{E} \left[\left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 \right] \quad (88) \end{aligned}$$

$$\begin{aligned} &= -\frac{1}{2} \left(\frac{1 - D(1 + \beta^2)}{1 - D} \right) \left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 + \tau_{\text{eff}} \eta L \sigma^2 \sum_{i=1}^m \frac{w_i^2 \|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} + \\ &\quad + \frac{\eta^2 L^2 \sigma^2}{1 - D} \sum_{i=1}^m w_i \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2 \right) + \frac{D\kappa^2}{2(1 - D)}. \quad (89) \end{aligned}$$

If $D \leq \frac{1}{2\beta^2+1}$, then it follows that $\frac{1}{1-D} \leq 1 + \frac{1}{2\beta^2}$ and $\frac{D\beta^2}{1-D} \leq \frac{1}{2}$. These facts can help us further simplify inequality (89).

$$\begin{aligned}
\frac{\mathbb{E} [\tilde{F}(\mathbf{x}^{(t+1,0)})] - \tilde{F}(\mathbf{x}^{(t,0)})}{\eta\tau_{\text{eff}}} &\leq -\frac{1}{4} \left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 + \tau_{\text{eff}}\eta L\sigma^2 \sum_{i=1}^m \frac{w_i^2 \|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} \\
&\quad + \eta^2 L^2 \sigma^2 \left(1 + \frac{1}{2\beta^2}\right) \sum_{i=1}^m w_i \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2\right) \\
&\quad + 2\eta^2 L^2 \max_i \{\|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1})\} \kappa^2 \left(1 + \frac{1}{2\beta^2}\right) \quad (90) \\
&\leq -\frac{1}{4} \left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 + \tau_{\text{eff}}\eta L\sigma^2 \sum_{i=1}^m \frac{w_i^2 \|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} \\
&\quad + \frac{3}{2}\eta^2 L^2 \sigma^2 \sum_{i=1}^m w_i \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2\right) \\
&\quad + 3\eta^2 L^2 \kappa^2 \max_i \{\|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1})\} \quad (91)
\end{aligned}$$

Taking the average across all rounds, we get

$$\begin{aligned}
\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[\left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 \right] &\leq \frac{4 [\tilde{F}(\mathbf{x}^{(0,0)}) - \tilde{F}_{\text{inf}}]}{\eta\tau_{\text{eff}}T} + 4\tau_{\text{eff}}\eta L\sigma^2 \sum_{i=1}^m \frac{w_i^2 \|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} \\
&\quad + 6\eta^2 L^2 \sigma^2 \sum_{i=1}^m w_i \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2\right) \\
&\quad + 12\eta^2 L^2 \kappa^2 \max_i \{\|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1})\}. \quad (92)
\end{aligned}$$

For the ease of writing, we define the following auxiliary variables:

$$A = m\tau_{\text{eff}} \sum_{i=1}^m \frac{w_i^2 \|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2}, \quad (93)$$

$$B = \sum_{i=1}^m w_i \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2\right), \quad (94)$$

$$C = \max_i \{\|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1})\}. \quad (95)$$

It follows that

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[\left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 \right] \leq \frac{4 [\tilde{F}(\mathbf{x}^{(0,0)}) - \tilde{F}_{\text{inf}}]}{\eta\tau_{\text{eff}}T} + \frac{4\eta L\sigma^2 A}{m} + 6\eta^2 L^2 \sigma^2 B + 12\eta^2 L^2 \kappa^2 C \quad (96)$$

Since $\min_{t \in [T]} \mathbb{E} \left[\left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 \right] \leq \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[\left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 \right]$, we have

$$\min_{t \in [T]} \mathbb{E} \left[\left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 \right] \leq \frac{4 [\tilde{F}(\mathbf{x}^{(0,0)}) - \tilde{F}_{\text{inf}}]}{\eta\tau_{\text{eff}}T} + \frac{4\eta L\sigma^2 A}{m} + 6\eta^2 L^2 \sigma^2 B + 12\eta^2 L^2 \kappa^2 C. \quad (97)$$

C.7 Constraint on Local Learning Rate

Here, let us summarize the constraints on local learning rate:

$$\eta L \leq \frac{1}{2\tau_{\text{eff}}}, \quad (98)$$

$$4\eta^2 L^2 \max_i \{\|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1})\} \leq \frac{1}{2\beta^2 + 1}. \quad (99)$$

For the second constraint, we can further tighten it as follows:

$$4\eta^2 L^2 \max_i \{\|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1})\} \leq 4\eta^2 L^2 \max_i \|\mathbf{a}_i\|_1^2 \leq \frac{1}{2\beta^2 + 1} \quad (100)$$

That is,

$$\eta L \leq \frac{1}{2} \min \left\{ \frac{1}{\max_i \|\mathbf{a}_i\|_1 \sqrt{2\beta^2 + 1}}, \frac{1}{\tau_{\text{eff}}} \right\}. \quad (101)$$

C.8 Further Optimizing the Bound

By setting $\eta = \sqrt{\frac{m}{\bar{\tau}T}}$ where $\bar{\tau} = \frac{1}{m} \sum_{i=1}^m \tau_i$, we have

$$\min_{t \in [T]} \mathbb{E} \left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 \leq \mathcal{O} \left(\frac{\bar{\tau}/\tau_{\text{eff}}}{\sqrt{m\bar{\tau}T}} \right) + \mathcal{O} \left(\frac{A\sigma^2}{\sqrt{m\bar{\tau}T}} \right) + \mathcal{O} \left(\frac{mB\sigma^2}{\bar{\tau}T} \right) + \mathcal{O} \left(\frac{mC\kappa^2}{\bar{\tau}T} \right). \quad (102)$$

Here, we complete the proof of Theorem 1.

D Proof of Theorem 2: Including Bias in the Error Bound

Lemma 3. *For any model parameter \mathbf{x} , the difference between the gradients of $F(\mathbf{x})$ and $\tilde{F}(\mathbf{x})$ can be bounded as follows:*

$$\|\nabla F(\mathbf{x}) - \nabla \tilde{F}(\mathbf{x})\|^2 \leq \chi_{\mathbf{p}\|\mathbf{w}}^2 \left[(\beta^2 - 1) \|\nabla \tilde{F}(\mathbf{x})\|^2 + \kappa^2 \right] \quad (103)$$

where $\chi_{\mathbf{p}\|\mathbf{w}}^2$ denotes the chi-square distance between \mathbf{p} and \mathbf{w} , i.e., $\chi_{\mathbf{p}\|\mathbf{w}}^2 = \sum_{i=1}^m (p_i - w_i)^2 / w_i$.

Proof. According to the definition of $F(\mathbf{x})$ and $\tilde{F}(\mathbf{x})$, we have

$$\nabla F(\mathbf{x}) - \nabla \tilde{F}(\mathbf{x}) = \sum_{i=1}^m (p_i - w_i) \nabla F_i(\mathbf{x}) \quad (104)$$

$$= \sum_{i=1}^m (p_i - w_i) (\nabla F_i(\mathbf{x}) - \nabla \tilde{F}(\mathbf{x})) \quad (105)$$

$$= \sum_{i=1}^m \frac{p_i - w_i}{\sqrt{w_i}} \cdot \sqrt{w_i} (\nabla F_i(\mathbf{x}) - \nabla \tilde{F}(\mathbf{x})). \quad (106)$$

Applying Cauchy–Schwarz inequality, it follows that

$$\|\nabla F(\mathbf{x}) - \nabla \tilde{F}(\mathbf{x})\|^2 \leq \left[\sum_{i=1}^m \frac{(p_i - w_i)^2}{w_i} \right] \left[\sum_{i=1}^m w_i \|\nabla F_i(\mathbf{x}) - \nabla \tilde{F}(\mathbf{x})\|^2 \right] \quad (107)$$

$$\leq \chi_{\mathbf{p}\|\mathbf{w}}^2 \left[(\beta^2 - 1) \|\nabla \tilde{F}(\mathbf{x})\|^2 + \kappa^2 \right]. \quad (108)$$

where the last inequality uses Assumption 3. \square

Note that

$$\|\nabla F(\mathbf{x})\|^2 \leq 2 \|\nabla F(\mathbf{x}) - \nabla \tilde{F}(\mathbf{x})\|^2 + 2 \|\nabla \tilde{F}(\mathbf{x})\|^2 \quad (109)$$

$$\leq 2 \left[\chi_{\mathbf{p}\|\mathbf{w}}^2 (\beta^2 - 1) + 1 \right] \|\nabla \tilde{F}(\mathbf{x})\|^2 + 2\chi_{\mathbf{p}\|\mathbf{w}}^2 \kappa^2. \quad (110)$$

As a result, we obtain

$$\min_{t \in [T]} \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \leq \frac{1}{T} \sum_{t=0}^{T-1} \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \quad (111)$$

$$\leq 2 \left[\chi_{\mathbf{p}\|\mathbf{w}}^2 (\beta^2 - 1) + 1 \right] \frac{1}{T} \sum_{t=0}^{T-1} \left\| \nabla \tilde{F}(\mathbf{x}^{(t,0)}) \right\|^2 + 2\chi_{\mathbf{p}\|\mathbf{w}}^2 \kappa^2 \quad (112)$$

$$\leq 2 \left[\chi_{\mathbf{p}\|\mathbf{w}}^2 (\beta^2 - 1) + 1 \right] \epsilon_{\text{opt}} + 2\chi_{\mathbf{p}\|\mathbf{w}}^2 \kappa^2 \quad (113)$$

where ϵ_{opt} denotes the optimization error.

D.1 Constructing a Lower Bound

In this subsection, we are going to construct a lower bound of $\mathbb{E} \left\| \nabla F(\mathbf{x})^{(t,0)} \right\|^2$, showing that (10) is tight and the non-vanishing error term in Theorem 2 is not an artifact of our analysis.

Lemma 4. *One can manually construct a strongly convex objective function such that FedAvg with heterogeneous local updates cannot converge to its global optimum. In particular, the gradient norm of the objective function does not vanish as learning rate approaches to zero. We have the following lower bound:*

$$\lim_{T \rightarrow \infty} \mathbb{E} \left\| \nabla F(\mathbf{x}^{(T,0)}) \right\|^2 = \Omega(\chi_{\mathbf{p}\|\mathbf{w}}^2 \kappa^2) \quad (114)$$

where $\chi_{\mathbf{p}\|\mathbf{w}}^2$ denotes the chi-square divergence between weight vectors and κ^2 quantifies the dissimilarities among local objective functions and is defined in Assumption 3.

Proof. Suppose that there are only two clients with local objectives $F_1(x) = \frac{1}{2}(x - a)^2$ and $F_2(x) = \frac{1}{2}(x + a)^2$. The global objective is defined as $F(x) = \frac{1}{2}F_1(x) + \frac{1}{2}F_2(x)$. For any set of weights $w_1, w_2, w_1 + w_2 = 1$, we define the surrogate objective function as $\tilde{F}(x) = w_1F_1(x) + w_2F_2(x)$. As a consequence, we have

$$\begin{aligned} & \sum_{i=1}^m w_i \left\| \nabla F_i(x) - \nabla \tilde{F}(x) \right\|^2 \\ &= w_1[(x - a) - [x - (w_1 - w_2)a]]^2 + w_2[(x + a) - [x - (w_1 - w_2)a]]^2 \end{aligned} \quad (115)$$

$$= w_1[2w_2a]^2 + w_2[2w_1a]^2 = 2(w_1 + w_2)(w_1w_2a^2) = 2w_1w_2a^2 \quad (116)$$

Comparing with Assumption 3, we can define $\kappa^2 = 2w_1w_2a^2$ and $\beta^2 = 1$ in this case. Furthermore, according to the derivations in Appendix A, the iterate of FedAvg can be written as follows:

$$\lim_{T \rightarrow \infty} x^{(T,0)} = \frac{\tau_1 a - \tau_2 a}{\tau_1 + \tau_2}. \quad (117)$$

As a results, we have

$$\lim_{T \rightarrow \infty} \left\| \nabla F(x^{(T,0)}) \right\|^2 = \lim_{T \rightarrow \infty} \left[\frac{1}{2}(x^{(T,0)} - a) + \frac{1}{2}(x^{(T,0)} + a) \right]^2 \quad (118)$$

$$= \lim_{T \rightarrow \infty} \left[x^{(T,0)} \right]^2 \quad (119)$$

$$= \left(\frac{\tau_1 - \tau_2}{\tau_1 + \tau_2} \right)^2 a^2 = \frac{(\tau_2 - \tau_1)^2}{2\tau_1\tau_2} \kappa^2 = \Omega(\chi_{\mathbf{p}\|\mathbf{w}}^2 \kappa^2). \quad (120)$$

where $\chi_{\mathbf{p}\|\mathbf{w}}^2 = \sum_{i=1}^m (p_i - w_i)^2 / w_i = (w_1 - 1/2)^2 / w_1 + (w_2 - 1/2)^2 / w_2$. \square

E Special Cases of Theorem 1

Here, we provide several instantiations of Theorem 1 and check its consistency with previous results.

E.1 FedAvg

In FedAvg, $\mathbf{a}_i = [1, 1, \dots, 1]^\top \in \mathbb{R}^{\tau_i}$, $\|\mathbf{a}_i\|_2^2 = \tau_i$, and $\|\mathbf{a}_i\|_1 = \tau_i$. In addition, we have $w_i = p_i \tau_i / (\sum_{i=1}^m p_i \tau_i)$. Accordingly, we get the closed-form expressions of the following quantities:

$$\tau_{\text{eff}} = \sum_{i=1}^m p_i \tau_i = \mathbb{E}_{\mathbf{p}}[\boldsymbol{\tau}], \quad (121)$$

$$A_{\text{FedAvg}} = m \tau_{\text{eff}} \sum_{i=1}^m \frac{w_i^2 \|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} = \frac{m \sum_{i=1}^m p_i^2 \tau_i}{\sum_{i=1}^m p_i \tau_i}, \quad (122)$$

$$B_{\text{FedAvg}} = \sum_{i=1}^m w_i \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2 \right) = \frac{\sum_{i=1}^m p_i \tau_i (\tau_i - 1)}{\sum_{i=1}^m p_i \tau_i} = \mathbb{E}_{\mathbf{p}}[\boldsymbol{\tau}] - 1 + \frac{\text{var}_{\mathbf{p}}[\boldsymbol{\tau}]}{\mathbb{E}_{\mathbf{p}}[\boldsymbol{\tau}]}, \quad (123)$$

$$C_{\text{FedAvg}} = \max_i \{ \|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1}) \} = \tau_{\max} (\tau_{\max} - 1). \quad (124)$$

In the case where all clients have the same local dataset size, *i.e.*, $p_i = 1/m, \forall i$. It follows that

$$\tau_{\text{eff}} = \bar{\tau}, A_{\text{FedAvg}} = 1, B_{\text{FedAvg}} = \bar{\tau} - 1 + \frac{\text{var}[\boldsymbol{\tau}]}{\bar{\tau}}, C_{\text{FedAvg}} = \tau_{\max}(\tau_{\max} - 1). \quad (125)$$

Substituting (125) into Theorem 1, we get the convergence guarantee for FedAvg. We formally state it in the following corollary.

Corollary 1 (Convergence of FedAvg). *Under the same conditions as Theorem 1, if $p_i = 1/m$, then FedAvg algorithm (vanilla SGD with fixed local learning rate as local solver) will converge to the stationary point of a surrogate objective $\tilde{F}(\mathbf{x}) = \sum_{i=1}^m \tau_i F_i(\mathbf{x}) / \sum_{i=1}^m \tau_i$. The optimization error will be bounded as follows:*

$$\min_{t \in [T]} \mathbb{E} \|\nabla \tilde{F}(\mathbf{x})\|^2 \leq \mathcal{O} \left(\frac{1 + \sigma^2}{\sqrt{m\bar{\tau}T}} \right) + \mathcal{O} \left(\frac{m\sigma^2(\bar{\tau} - 1 + \text{var}[\boldsymbol{\tau}]/\bar{\tau})}{\bar{\tau}T} \right) + \mathcal{O} \left(\frac{m\kappa^2\tau_{\max}(\tau_{\max} - 1)}{\bar{\tau}T} \right) \quad (126)$$

where \mathcal{O} swallows all constants (including L), and $\text{var}[\boldsymbol{\tau}] = \sum_{i=1}^m \tau_i^2/m - \bar{\tau}^2$ denotes the variance of local steps.

Consistent with Previous Results. When all clients perform the same local steps, *i.e.*, $\tau_i = \tau$, then $\text{var}[\boldsymbol{\tau}] = 0$ and the above error bound (126) recovers previous results [8, 24, 20]. When $\tau_i = 1$, then FedAvg reduces to fully synchronous SGD and the error bound (126) becomes $1/\sqrt{mT}$, which is the same as standard SGD convergence rate [51].

E.2 FedProx

In FedProx, we have $\mathbf{a}_i = [(1 - \alpha)^{\tau_i - 1}, \dots, (1 - \alpha), 1]^\top \in \mathbb{R}^{\tau_i}$. Accordingly, the norms of \mathbf{a}_i can be written as:

$$\|\mathbf{a}_i\|_2^2 = \frac{1 - (1 - \alpha)^{2\tau_i}}{1 - (1 - \alpha)^2}, \|\mathbf{a}_i\|_1 = \frac{1 - (1 - \alpha)^{\tau_i}}{\alpha}, w_i = \frac{p_i[1 - (1 - \alpha)^{\tau_i}]}{\sum_{i=1}^m p_i[1 - (1 - \alpha)^{\tau_i}]} \quad (127)$$

As a consequence, we can derive the closed-form expression of $\tau_{\text{eff}}, A, B, C$ as follows:

$$\tau_{\text{eff}} = \frac{1}{\alpha} \sum_{i=1}^m p_i[1 - (1 - \alpha)^{\tau_i}], \quad (128)$$

$$A_{\text{FedProx}} = \frac{m\alpha}{\sum_{i=1}^m p_i(1 - (1 - \alpha)^{\tau_i})} \sum_{i=1}^m p_i^2 \frac{1 - (1 - \alpha)^{2\tau_i}}{1 - (1 - \alpha)^2}, \quad (129)$$

$$B_{\text{FedProx}} = \sum_{i=1}^m \frac{p_i[1 - (1 - \alpha)^{\tau_i}]}{\sum_{i=1}^m p_i[1 - (1 - \alpha)^{\tau_i}]} \left[\frac{1 - (1 - \alpha)^{2\tau_i}}{1 - (1 - \alpha)^2} - 1 \right], \quad (130)$$

$$C_{\text{FedProx}} = \frac{1 - (1 - \alpha)^{\tau_{\max}}}{\alpha} \left(\frac{1 - (1 - \alpha)^{\tau_{\max}}}{\alpha} - 1 \right). \quad (131)$$

Substituting $A_{\text{FedProx}}, B_{\text{FedProx}}, C_{\text{FedProx}}$ back into Theorem 1, one can obtain the convergence guarantee for FedProx. Again, it will converge to the stationary points of a surrogate objective due to $w_i \neq p_i$.

Consistency with FedAvg. From the update rule of FedProx, we know that when $\mu = 0$ (or $\alpha = 0$), FedProx is equivalent to FedAvg. This can also be validated from the expressions of $A_{\text{FedProx}}, B_{\text{FedProx}}, C_{\text{FedProx}}$. Using L'Hospital law, it is easy to show that

$$\lim_{\alpha \rightarrow 0} A_{\text{FedProx}} = A_{\text{FedAvg}}, \lim_{\alpha \rightarrow 0} B_{\text{FedProx}} = B_{\text{FedAvg}}, \lim_{\alpha \rightarrow 0} C_{\text{FedProx}} = C_{\text{FedAvg}}. \quad (132)$$

Best value of α in FedProx. Given the expressions of τ_{eff} and A, B, C , we can further select a best value of α that optimizes the error bound of FedProx, as stated in the following corollary.

Corollary 2. *Under the same conditions as Theorem 1 and suppose $p_i = 1/m$ and $\tau_i \gg 1$, then $\alpha = \mathcal{O}(m^{\frac{1}{2}}/\bar{\tau}^{\frac{1}{2}}T^{\frac{1}{6}})$ minimizes the optimization error bound of FedProx in terms of converging to the stationary points of the surrogate objective. In particular, we have*

$$\min_{t \in [T]} \mathbb{E} \|\nabla \tilde{F}(\mathbf{x})\|^2 \leq \mathcal{O} \left(\frac{1}{\sqrt{m\bar{\tau}T}} \right) + \mathcal{O} \left(\frac{1}{T^{\frac{2}{3}}} \right) \quad (133)$$

where \mathcal{O} swallows all other constants. Furthermore, if we define $K = \bar{\tau}T$ the average gradient evaluations at clients and let $\bar{\tau} \leq \mathcal{O}(K^{\frac{1}{4}}m^{-\frac{3}{4}})$ (which is equivalent to $T \geq \mathcal{O}(K^{\frac{3}{4}}m^{\frac{3}{4}})$), then it follows that $\min_{t \in [T]} \mathbb{E} \|\nabla \tilde{F}(\mathbf{x})\|^2 \leq \mathcal{O}(1/\sqrt{mK})$.

Discussion: Corollary 2 shows that there exists a non-zero value of α that optimizes the error upper bound of FedProx. That is to say, FedProx ($\alpha > 0$) is better than FedAvg ($\alpha = 0$) by a constant in terms of error upper bound. However, on the other hand, it is worth noting that the minimal communication rounds of FedProx to achieve $1/\sqrt{mK}$ rate, given by Corollary 2, is exactly the same as FedAvg [24]. In this sense, FedProx has the same convergence rate as FedAvg and cannot further reduce the communication overhead.

Proof. First of all, let us relax the error terms of FedProx. Under the assumption of $\tau_i \gg 1$, the quantities A, B, C can be bounded or approximated as follows:

$$\tau_{\text{eff}} \simeq \frac{1}{\alpha}, \quad (134)$$

$$A_{\text{FedProx}} \simeq m\alpha \sum_{i=1}^m \frac{p_i^2}{(2-\alpha)\alpha} = \frac{m \sum_{i=1}^m p_i^2}{2-\alpha} \leq m \sum_{i=1}^m p_i^2 = 1, \quad (135)$$

$$B_{\text{FedProx}} \leq \frac{1 - (1-\alpha)^{2\tau_i}}{1 - (1-\alpha)^2} - 1 \leq \frac{1}{\alpha(2-\alpha)} \leq \frac{1}{\alpha} \leq \frac{1}{\alpha^2}, \quad (136)$$

$$C_{\text{FedProx}} \leq \frac{1}{\alpha^2}. \quad (137)$$

Accordingly, the error upper bound of FedProx can be rewritten as follows:

$$\min_{t \in [T]} \mathbb{E} \|\nabla \tilde{F}(\mathbf{x})\|^2 \leq \mathcal{O}\left(\frac{\alpha\bar{\tau}}{\sqrt{m\bar{\tau}T}}\right) + \mathcal{O}\left(\frac{1}{\sqrt{m\bar{\tau}T}}\right) + \mathcal{O}\left(\frac{m}{\alpha^2\bar{\tau}T}\right). \quad (138)$$

In order to optimize the above bound, we can simply take the derivative with respect to α . When the derivative equals to zero, we get

$$\frac{\bar{\tau}}{\sqrt{m\bar{\tau}T}} = \frac{m}{\alpha^3\bar{\tau}T} \implies \alpha = \mathcal{O}\left(\frac{m^{\frac{1}{2}}}{\bar{\tau}^{\frac{1}{2}}T^{\frac{1}{6}}}\right). \quad (139)$$

Plugging the expression of best α into (138), we have

$$\min_{t \in [T]} \mathbb{E} \|\nabla \tilde{F}(\mathbf{x})\|^2 \leq \mathcal{O}\left(\frac{1}{\sqrt{m\bar{\tau}T}}\right) + \mathcal{O}\left(\frac{1}{T^{\frac{2}{3}}}\right) = \mathcal{O}\left(\frac{1}{\sqrt{mK}}\right) + \mathcal{O}\left(\frac{\tau^{\frac{2}{3}}}{K^{\frac{2}{3}}}\right) \quad (140)$$

where $K = \bar{\tau}T$ denotes the average total gradient steps at clients. In order to let the first term dominates the convergence rate, it requires that

$$\frac{1}{\sqrt{mK}} \geq \frac{\bar{\tau}^{\frac{2}{3}}}{K^{\frac{2}{3}}} \implies \bar{\tau} \leq \mathcal{O}\left(K^{\frac{1}{4}}m^{-\frac{3}{4}}\right). \quad (141)$$

As a results, the total communication rounds $T = K/\bar{\tau}$ should be greater than $\mathcal{O}(K^{\frac{3}{4}}m^{\frac{3}{4}})$. \square

F Proof of Theorem 3

In the case of FedNova, the aggregated weights w_i equals to p_i . Therefore, the surrogate objective $\tilde{F}(\mathbf{x}) = \sum_{i=1}^m w_i F_i(\mathbf{x})$ is the same as the original objective function $F(\mathbf{x}) = \sum_{i=1}^m p_i F_i(\mathbf{x})$. We can directly reuse the intermediate results in the proof of Theorem 1. According to (91), we have

$$\frac{\mathbb{E}[F(\mathbf{x}^{(t+1,0)})] - F(\mathbf{x}^{(t,0)})}{\eta\tau_{\text{eff}}} \leq -\frac{1}{4} \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 + \frac{\eta L \sigma^2 A^{(t)}}{m} + \frac{3}{2} \eta^2 L^2 \sigma^2 B^{(t)} + 3\eta^2 L^2 \kappa^2 C^{(t)} \quad (142)$$

where quantities $A^{(t)}, B^{(t)}, C^{(t)}$ are defined as follows:

$$A^{(t)} = m\tau_{\text{eff}} \sum_{i=1}^m \frac{w_i^2 \left\| \mathbf{a}_i^{(t)} \right\|_2^2}{\left\| \mathbf{a}_i^{(t)} \right\|_1^2}, \quad (143)$$

$$B^{(t)} = \sum_{i=1}^m p_i \left(\left\| \mathbf{a}_i^{(t)} \right\|_2^2 - [a_{i,-1}^{(t)}]^2 \right), \quad (144)$$

$$C^{(t)} = \max_i \left\{ \left\| \mathbf{a}_i^{(t)} \right\|_1 \left(\left\| \mathbf{a}_i^{(t)} \right\|_1 - a_{i,-1}^{(t)} \right) \right\}. \quad (145)$$

Taking the total expectation and averaging over all rounds, it follows that

$$\begin{aligned} \frac{\mathbb{E}[F(\mathbf{x}^{(T,0)})] - F(\mathbf{x}^{(0,0)})}{\eta\tau_{\text{eff}}T} &\leq -\frac{1}{4T} \sum_{t=0}^{T-1} \mathbb{E} \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 + \frac{\eta L \sigma^2 \tilde{A}}{m} \\ &\quad + \frac{3}{2} \eta^2 L^2 \sigma^2 \tilde{B} + 3\eta^2 L^2 \kappa^2 \tilde{C} \end{aligned} \quad (146)$$

where $\tilde{A} = \sum_{t=0}^{T-1} A^{(t)}/T$, $\tilde{B} = \sum_{t=0}^{T-1} B^{(t)}/T$, and $\tilde{C} = \sum_{t=0}^{T-1} C^{(t)}/T$. After minor rearranging, we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[\left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \right] \leq \frac{4[F(\mathbf{x}^{(0,0)}) - F_{\text{inf}}]}{\eta\tau_{\text{eff}}T} + \frac{4\eta L \sigma^2 \tilde{A}}{m} + 6\eta^2 L^2 \sigma^2 \tilde{B} + 12\eta^2 L^2 \kappa^2 \tilde{C}. \quad (147)$$

By setting $\eta = \sqrt{\frac{m}{\tilde{\tau}T}}$ where $\tilde{\tau} = \sum_{t=0}^{T-1} \tilde{\tau}^{(t)}/T$, the above upper bound can be further optimized as follows:

$$\begin{aligned} \min_{t \in [T]} \mathbb{E} \left[\left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \right] &\leq \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[\left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \right] \\ &\leq \frac{4\tilde{\tau}/\tau_{\text{eff}} \cdot [F(\mathbf{x}^{(0,0)}) - F_{\text{inf}}]}{\sqrt{m\tilde{\tau}T}} + \frac{4L\sigma^2\tilde{A}}{\sqrt{m\tilde{\tau}T}} + \frac{6mL^2\sigma^2\tilde{B}}{\tilde{\tau}T} + \frac{12mL^2\kappa^2\tilde{C}}{\tilde{\tau}T} \\ &= \mathcal{O} \left(\frac{\tilde{\tau}/\tau_{\text{eff}}}{\sqrt{m\tilde{\tau}T}} \right) + \mathcal{O} \left(\frac{\tilde{A}\sigma^2}{\sqrt{m\tilde{\tau}T}} \right) + \mathcal{O} \left(\frac{m\tilde{B}\sigma^2}{\tilde{\tau}T} \right) + \mathcal{O} \left(\frac{m\tilde{C}\kappa^2}{\tilde{\tau}T} \right). \end{aligned} \quad (148)$$

$$(149)$$

$$(150)$$

Here, we complete the proof of Theorem 3.

Moreover, it is worth mentioning the constraints on the local learning rate. Recall that, at the t -th round, we have the following constraint:

$$\eta L \leq \frac{1}{2} \min \left\{ \frac{1}{\max_i \left\| \mathbf{a}_i^{(t)} \right\|_1 \sqrt{2\beta^2 + 1}}, \frac{1}{\tau_{\text{eff}}} \right\}. \quad (151)$$

In order to guarantee the convergence, the above inequality should hold in every round. That is to say,

$$\eta L \leq \frac{1}{2} \min \left\{ \frac{1}{\max_{i \in [m], t \in [T]} \left\| \mathbf{a}_i^{(t)} \right\|_1 \sqrt{2\beta^2 + 1}}, \frac{1}{\tau_{\text{eff}}} \right\}. \quad (152)$$

G Extension: Incorporating Client Sampling

In this section, we extend the convergence guarantee of FedNova to the case of client sampling. Following previous works [38, 12, 20, 15], we assume the sampling scheme guarantees that the update rule (11) hold in expectation. This can be achieved by sampling with replacement from $\{1, 2, \dots, m\}$ with probabilities $\{p_i\}$, and averaging local updates from selected clients with equal weights. Specifically, we have

$$\mathbf{x}^{(t+1,0)} - \mathbf{x}^{(t,0)} = -\tau_{\text{eff}} \sum_{j=1}^q \frac{1}{q} \cdot \eta \mathbf{d}_{l_j}^{(t)} \quad \text{where } \mathbf{d}_{l_j}^{(t)} = \mathbf{G}_{l_j}^{(t)} \mathbf{a}_{l_j} / \|\mathbf{a}_{l_j}\|_1 \quad (153)$$

where q is the number of selected clients per round, and l_j is a random index sampled from $\{1, 2, \dots, m\}$ satisfying $\mathbb{P}(l_j = i) = p_i$. Recall that $p_i = n_i/n$ is the relative sample size at client i . For the ease of presentation, let \mathbf{a}_i to be fixed across rounds. One can directly validate that

$$\mathbb{E}_S \left[\frac{1}{q} \sum_{j=1}^q \mathbf{d}_{l_j}^{(t)} \right] = \frac{1}{q} \sum_{j=1}^q \mathbb{E}_S \left[\mathbf{d}_{l_j}^{(t)} \right] = \mathbb{E}_S \left[\mathbf{d}_{l_j}^{(t)} \right] = \sum_{i=1}^m p_i \mathbf{d}_i^{(t)} \quad (154)$$

where \mathbb{E}_S represents the expectation over random indices at current round.

Corollary 3. *Under the same condition as Theorem 1, suppose at each round, the server randomly selects $q (\leq m)$ clients with replacement to perform local computation. The probability of choosing the i -th client is $p_i = n_i/n$. In this case, FedNova will converge to the stationary points of the global objective $F(\mathbf{x})$. If we set $\eta = \sqrt{q/\tilde{\tau}T}$ where $\tilde{\tau}$ is the average local updates across all rounds, then the expected gradient norm is bounded as follows:*

$$\min_{t \in [T]} \mathbb{E} \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \leq \mathcal{O} \left(\frac{\tilde{\tau}/\tau_{\text{eff}}}{\sqrt{q\tilde{\tau}T}} \right) + \mathcal{O} \left(\frac{\tau_{\text{eff}}/\tilde{\tau}}{\sqrt{q\tilde{\tau}T}} \right) + \mathcal{O} \left(\frac{q(B+C)}{\tilde{\tau}T} \right) \quad (155)$$

where \mathcal{O} swallows all other constants (including L, σ^2, κ^2).

Proof. According to the Lipschitz-smooth assumption, it follows that

$$\mathbb{E} \left[F(\mathbf{x}^{(t+1,0)}) \right] - F(\mathbf{x}^{(t,0)}) \leq \underbrace{-\tau_{\text{eff}}\eta \mathbb{E} \left[\left\langle \nabla F(\mathbf{x}^{(t,0)}), \sum_{j=1}^q \frac{\mathbf{d}_{l_j}^{(t)}}{q} \right\rangle \right]}_{T_3} + \underbrace{\frac{\tau_{\text{eff}}^2 \eta^2 L}{2} \mathbb{E} \left[\left\| \sum_{j=1}^q \frac{\mathbf{d}_{l_j}^{(t)}}{q} \right\|^2 \right]}_{T_4} \quad (156)$$

where the expectation is taken over randomly selected indices $\{l_j\}$ as well as mini-batches $\xi_i^{(t,k)}, \forall i \in \{1, 2, \dots, m\}, k \in \{0, 1, \dots, \tau_i - 1\}$.

For the first term in (156), we can first take the expectation over indices and obtain

$$T_3 = \mathbb{E} \left[\left\langle \nabla F(\mathbf{x}^{(t,0)}), \mathbb{E}_S \left[\sum_{j=1}^q \frac{\mathbf{d}_{l_j}^{(t)}}{q} \right] \right\rangle \right] \quad (157)$$

$$= \mathbb{E} \left[\left\langle \nabla F(\mathbf{x}^{(t,0)}), \sum_{i=1}^m p_i \mathbf{d}_i^{(t)} \right\rangle \right]. \quad (158)$$

This term is exactly the same as the first term in (49). We can directly reuse previous results in the proof of Theorem 1. Comparing with (56), we have

$$T_3 = \frac{1}{2} \left\| \nabla F(\mathbf{x}^{(t)}) \right\|^2 + \frac{1}{2} \mathbb{E} \left[\left\| \sum_{i=1}^m p_i \mathbf{h}_i^{(t)} \right\|^2 \right] - \frac{1}{2} \mathbb{E} \left[\left\| \nabla F(\mathbf{x}^{(t,0)}) - \sum_{i=1}^m p_i \mathbf{h}_i^{(t)} \right\|^2 \right] \quad (159)$$

$$\geq \frac{1}{2} \left\| \nabla F(\mathbf{x}^{(t)}) \right\|^2 + \frac{1}{2} \mathbb{E} \left[\left\| \sum_{i=1}^m p_i \mathbf{h}_i^{(t)} \right\|^2 \right] - \frac{1}{2} \sum_{i=1}^m p_i \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) - \mathbf{h}_i^{(t)} \right\|^2 \right]. \quad (160)$$

For the second term in (156),

$$T_4 \leq 2\mathbb{E} \left[\left\| \frac{1}{q} \sum_{j=1}^q (\mathbf{d}_{l_j}^{(t)} - \mathbf{h}_{l_j}^{(t)}) \right\|^2 \right] + 2\mathbb{E} \left[\left\| \frac{1}{q} \sum_{j=1}^q \mathbf{h}_{l_j}^{(t)} \right\|^2 \right] \quad (161)$$

$$= \frac{1}{q} \sum_{i=1}^m p_i \mathbb{E} \left[\left\| \mathbf{d}_i^{(t)} - \mathbf{h}_i^{(t)} \right\|^2 \right] + 2\mathbb{E} \left[\left\| \frac{1}{q} \sum_{j=1}^q \mathbf{h}_{l_j}^{(t)} \right\|^2 \right] \quad (162)$$

$$\leq \frac{2\sigma^2}{q} \sum_{i=1}^m p_i \frac{\|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} + 2\mathbb{E} \left[\left\| \frac{1}{q} \sum_{j=1}^q \mathbf{h}_{l_j}^{(t)} \right\|^2 \right] \quad (163)$$

$$\leq \frac{2\sigma^2}{q} \sum_{i=1}^m p_i \frac{\|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} + 6 \sum_{i=1}^m p_i \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) - \mathbf{h}_i^{(t)} \right\|^2 \right] + \frac{6}{q} \left(\beta^2 \|\nabla F(\mathbf{x}^{(t,0)})\|^2 + \kappa^2 \right) + 6 \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \quad (164)$$

where the last inequality comes from Lemma 5, stated below.

Lemma 5. *Suppose we are given $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_m, \mathbf{x} \in \mathbb{R}^d$ and let l_1, l_2, \dots, l_q be i.i.d. sampled from a multinomial distribution \mathcal{D} supported on $\{1, 2, \dots, m\}$ satisfying $\mathbb{P}(l = i) = p_i$ and $\sum_{i=1}^m p_i = 1$. We have*

$$\mathbb{E} \left[\frac{1}{q} \sum_{j=1}^q \mathbf{z}_{l_j} \right] = \sum_{i=1}^m p_i \mathbf{z}_i, \quad (165)$$

$$\mathbb{E} \left[\left\| \frac{1}{q} \sum_{j=1}^q \mathbf{z}_{l_j} \right\|^2 \right] \leq 3 \sum_{i=1}^m p_i \|\mathbf{z}_i - \nabla F_i(\mathbf{x})\|^2 + 3 \|\nabla F(\mathbf{x})\|^2 + \frac{3}{q} (\beta^2 \|\nabla F(\mathbf{x})\|^2 + \kappa^2). \quad (166)$$

Proof. First, we have

$$\begin{aligned} & \mathbb{E} \left[\left\| \frac{1}{q} \sum_{j=1}^q \mathbf{z}_{l_j} \right\|^2 \right] \\ &= \mathbb{E} \left[\left\| \left(\frac{1}{q} \sum_{j=1}^q \mathbf{z}_{l_j} - \frac{1}{q} \sum_{j=1}^q \nabla F_{l_j}(\mathbf{x}) \right) + \left(\frac{1}{q} \sum_{j=1}^q \nabla F_{l_j}(\mathbf{x}) - \nabla F(\mathbf{x}) \right) + \nabla F(\mathbf{x}) \right\|^2 \right] \quad (167) \\ &\leq 3\mathbb{E} \left[\left\| \frac{1}{q} \sum_{j=1}^q \mathbf{z}_{l_j} - \frac{1}{q} \sum_{j=1}^q \nabla F_{l_j}(\mathbf{x}) \right\|^2 \right] + 3\mathbb{E} \left[\left\| \frac{1}{q} \sum_{j=1}^q \nabla F_{l_j}(\mathbf{x}) - \nabla F(\mathbf{x}) \right\|^2 \right] + 3 \|\nabla F(\mathbf{x})\|^2. \end{aligned} \quad (168)$$

For the first term, by Cauchy-Schwarz inequality, we have

$$\mathbb{E} \left[\left\| \frac{1}{q} \sum_{j=1}^q \mathbf{z}_{l_j} - \frac{1}{q} \sum_{j=1}^q \nabla F_{l_j}(\mathbf{x}) \right\|^2 \right] \leq \frac{1}{q} \sum_{j=1}^q \mathbb{E}_{l_j \sim \mathcal{D}} \left[\|\mathbf{z}_{l_j} - \nabla F_{l_j}(\mathbf{x})\|^2 \right] = \sum_{i=1}^m p_i \|\mathbf{z}_i - \nabla F_i(\mathbf{x})\|^2. \quad (169)$$

The second term can be bounded as following

$$\mathbb{E} \left[\left\| \frac{1}{q} \sum_{j=1}^q \nabla F_{l_j}(\mathbf{x}) - \nabla F(\mathbf{x}) \right\|^2 \right] = \frac{1}{q} \mathbb{E}_{l_j \sim \mathcal{D}} \left[\|\nabla F_{l_j}(\mathbf{x}) - \nabla F(\mathbf{x})\|^2 \right] \quad (170)$$

$$= \frac{1}{q} \sum_{i=1}^m p_i \|\nabla F_i(\mathbf{x}) - \nabla F(\mathbf{x})\|^2 \quad (171)$$

$$\leq \frac{1}{q} [(\beta^2 - 1) \|\nabla F(\mathbf{x})\|^2 + \kappa^2]. \quad (172)$$

where the first identity follows from $\mathbb{E}_{i \sim \mathcal{D}}[F_i(\mathbf{x})] = \nabla F(\mathbf{x})$ and the independence between l_1, \dots, l_q , and the last inequality is a direct application of Assumption 3.

Substituting (169) and (170) into (167) completes the proof. \square

Substituting (160) and (164) into (156), we have

$$\begin{aligned}
\frac{\mathbb{E}[F(\mathbf{x}^{(t+1,0)})] - F(\mathbf{x}^{(t,0)})}{\eta\tau_{\text{eff}}} &\leq -\frac{1}{2}(1 - 6\tau_{\text{eff}}\eta L) \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \\
&\quad + \left(\frac{1}{2} + 3\tau_{\text{eff}}\eta L \right) \sum_{i=1}^m p_i \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) - \mathbf{h}_i^{(t)} \right\|^2 \right] \\
&\quad + \frac{\tau_{\text{eff}}\eta L \sigma^2}{q} \sum_{i=1}^m p_i \frac{\|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} + \frac{3\tau_{\text{eff}}\eta L}{q} \left(\beta^2 \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 + \kappa^2 \right) \\
&\hspace{15em} (173) \\
&= -\frac{1}{2} \left(1 - 6\tau_{\text{eff}}\eta L - \frac{6\tau_{\text{eff}}\eta L \beta^2}{q} \right) \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \\
&\quad + \frac{\tau_{\text{eff}}\eta L \sigma^2}{q} \sum_{i=1}^m p_i \frac{\|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} \\
&\quad + \left(\frac{1}{2} + 2\tau_{\text{eff}}\eta L \right) \sum_{i=1}^m p_i \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) - \mathbf{h}_i^{(t)} \right\|^2 \right] + \frac{3\tau_{\text{eff}}\eta L \kappa^2}{q}. \\
&\hspace{15em} (174)
\end{aligned}$$

When $\eta L \leq 1/(2\tau_{\text{eff}})$ and $6\tau_{\text{eff}}\eta L + 6\tau_{\text{eff}}\eta L \beta^2/q \leq \frac{1}{2}$, it follows that

$$\begin{aligned}
\frac{\mathbb{E}[F(\mathbf{x}^{(t+1,0)})] - F(\mathbf{x}^{(t,0)})}{\eta\tau_{\text{eff}}} &\leq -\frac{1}{4} \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 + \frac{\tau_{\text{eff}}\eta L \sigma^2}{q} \sum_{i=1}^m p_i \frac{\|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} \\
&\quad + \frac{3}{2} \sum_{i=1}^m p_i \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) - \mathbf{h}_i^{(t)} \right\|^2 \right] + \frac{3\tau_{\text{eff}}\eta L \kappa^2}{q}. \quad (175)
\end{aligned}$$

Recall that the third term in (175) can be bounded as follows (see (87)):

$$\begin{aligned}
\frac{1}{2} \sum_{i=1}^m p_i \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) - \mathbf{h}_i^{(t)} \right\|^2 \right] &\leq \frac{\eta^2 L^2 \sigma^2}{1-D} \sum_{i=1}^m p_i \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2 \right) \\
&\quad + \frac{D\beta^2}{2(1-D)} \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 + \frac{D\kappa^2}{2(1-D)} \quad (176)
\end{aligned}$$

where $D = 4\eta^2 L^2 \max_i \{ \|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1}) \} < 1$. If $D \leq \frac{1}{12\beta^2+1}$, then it follows that $\frac{1}{1-D} \leq 1 + \frac{1}{12\beta^2} \leq 2$ and $\frac{3D\beta^2}{1-D} \leq \frac{1}{4}$. These facts can help us further simplify inequality (176). One can obtain

$$\begin{aligned}
\frac{3}{2} \sum_{i=1}^m p_i \mathbb{E} \left[\left\| \nabla F_i(\mathbf{x}^{(t,0)}) - \mathbf{h}_i^{(t)} \right\|^2 \right] &\leq 6\eta^2 L^2 \sigma^2 \sum_{i=1}^m p_i \left(\|\mathbf{a}_i\|_2^2 - [a_{i,-1}]^2 \right) + \frac{1}{8} \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \\
&\quad + 12\eta^2 L^2 \kappa^2 \max_i \{ \|\mathbf{a}_i\|_1 (\|\mathbf{a}_i\|_1 - a_{i,-1}) \} \quad (177)
\end{aligned}$$

$$= 6\eta^2 L^2 \sigma^2 B + \frac{1}{8} \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 + 12\eta^2 L^2 \kappa^2 C \quad (178)$$

Substituting (178) into (175), we have

$$\begin{aligned}
\frac{\mathbb{E}[F(\mathbf{x}^{(t+1,0)})] - F(\mathbf{x}^{(t,0)})}{\eta\tau_{\text{eff}}} &\leq -\frac{1}{8} \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 + \frac{\tau_{\text{eff}}\eta L \sigma^2}{q} \sum_{i=1}^m p_i \frac{\|\mathbf{a}_i\|_2^2}{\|\mathbf{a}_i\|_1^2} + \frac{3\tau_{\text{eff}}\eta L \kappa^2}{q} \\
&\quad + 6\eta^2 L^2 \sigma^2 B + 12\eta^2 L^2 \kappa^2 C \quad (179)
\end{aligned}$$

$$\begin{aligned}
&\leq -\frac{1}{8} \left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 + \frac{\tau_{\text{eff}}\eta L \sigma^2}{q} + \frac{3\tau_{\text{eff}}\eta L \kappa^2}{q} \\
&\quad + 6\eta^2 L^2 \sigma^2 B + 12\eta^2 L^2 \kappa^2 C \quad (180)
\end{aligned}$$

where the last inequality uses the fact that $\|\mathbf{a}\|_2 \leq \|\mathbf{a}\|_1$, for any vector \mathbf{a} . Taking the total expectation and averaging all rounds, one can obtain

$$\begin{aligned} \frac{\mathbb{E} [F(\mathbf{x}^{(T,0)})] - F(\mathbf{x}^{(0,0)})}{\eta\tau_{\text{eff}}T} &\leq -\frac{1}{8T} \sum_{t=0}^{T-1} \mathbb{E} \left[\left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \right] + \frac{\tau_{\text{eff}}\eta L(\sigma^2 + 3\kappa^2)}{q} \\ &\quad + 6\eta^2 L^2 \sigma^2 B + 12\eta^2 L^2 \kappa^2 C. \end{aligned} \quad (181)$$

After minor rearranging, the above inequality is equivalent to

$$\begin{aligned} &\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[\left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \right] \\ &\leq \frac{8 [F(\mathbf{x}^{(0,0)}) - F_{\text{inf}}]}{\eta\tau_{\text{eff}}T} + \frac{8\tau_{\text{eff}}\eta L(\sigma^2 + 3\kappa^2)}{q} + 48\eta^2 L^2 \sigma^2 B + 96\eta^2 L^2 \kappa^2 C. \end{aligned} \quad (182)$$

If we set the learning rate to be small enough, *i.e.*, $\eta = \sqrt{\frac{q}{\tilde{\tau}T}}$ where $\tilde{\tau} = \sum_{t=0}^{T-1} \tilde{\tau}/T$, then we get

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[\left\| \nabla F(\mathbf{x}^{(t,0)}) \right\|^2 \right] \leq \mathcal{O} \left(\frac{\tilde{\tau}/\tau_{\text{eff}}}{\sqrt{q\tilde{\tau}T}} \right) + \mathcal{O} \left(\frac{\tau_{\text{eff}}/\tilde{\tau}}{\sqrt{q\tilde{\tau}T}} \right) + \mathcal{O} \left(\frac{q(B+C)}{\tilde{\tau}T} \right) \quad (183)$$

where \mathcal{O} swallows all other constants. \square

H Pseudo-code of FedNova

Here we provide a pseudo-code of FedNova (see Algorithm 1) as a general algorithmic framework. Then, as an example, we show the pseudo-code of a special case of FedNova, where the local solver is specified as momentum SGD with cross-client variance reduction [21, 20] (see Algorithm 2). Note that when the server updates the global model, we set τ_{eff} to be the same as FedAvg, *i.e.*, $\tau_{\text{eff}} = \sum_{i \in \mathcal{S}_t} p_i \|\mathbf{a}_i^{(t)}\|_1$ where \mathcal{S}_t denotes the randomly selected subset of clients. Alternatively, the server can also choose other values of τ_{eff} .

Algorithm 1: FedNova Framework

Input: Client learning rate η ; Client momentum factor ρ .

- 1 **for** $t \in \{0, 1, \dots, T - 1\}$ **do**
- 2 Randomly sample a subset of clients \mathcal{S}_t
- 3 **Communication:** Broadcast global model $\mathbf{x}^{(t,0)}$ to selected clients
- 4 Clients perform local updates
- 5 **Communication:** Receive $\|\mathbf{a}_i^{(t)}\|_1$ and $\mathbf{d}_i^{(t)}$ from clients
- 6 Update global model: $\mathbf{x}^{(t+1,0)} = \mathbf{x}^{(t,0)} - \frac{\sum_{i \in \mathcal{S}_t} p_i \|\mathbf{a}_i^{(t)}\|_1}{\sum_{i \in \mathcal{S}_t} p_i} \sum_{i \in \mathcal{S}_t} \frac{\eta p_i \mathbf{d}_i^{(t)}}{\sum_{i \in \mathcal{S}_t} p_i}$
- 7 **end**

Algorithm 2: FedNova with Client-side Momentum SGD + Cross-client Variance Reduction

Input: Client learning rate η ; Client momentum factor ρ .

- 1 **for** $t \in \{0, 1, \dots, T - 1\}$ **at client** i **in parallel do**
- 2 Zero client optimizer buffers $\mathbf{u}_i^{(t,0)} = 0$
- 3 **Communication:** Receive $\mathbf{x}^{(t,0)} = \mathbf{x}^{(t-1,0)} - (\sum_{i=1}^m p_i a_i) \eta \sum_{i=1}^m p_i \mathbf{d}_i^{(t-1)}$ from server
- 4 **Communication:** Receive $\sum_{i=1}^m p_i \mathbf{d}_i^{(t-1)}$ from server
- 5 Update gradient correction term: $\mathbf{c}_i^{(t)} = -\mathbf{d}_i^{(t-1)} + \sum_{i=1}^m p_i \mathbf{d}_i^{(t-1)}$
- 6 **for** $k \in \{0, 1, \dots, \tau_i - 1\}$ **do**
- 7 Compute: $\tilde{g}_i(\mathbf{x}^{(t,k)}) = g_i(\mathbf{x}^{(t,k)}) + \mathbf{c}_i^{(t)}$
- 8 Update momentum buffer: $\mathbf{u}_i^{(t,k)} = \rho \mathbf{u}_i^{(t,k-1)} + \tilde{g}_i(\mathbf{x}^{(t,k)})$
- 9 Update local model: $\mathbf{x}_i^{(t,k)} = \mathbf{x}_i^{(t,k-1)} - \eta \mathbf{u}_i^{(t,k)}$
- 10 **end**
- 11 Compute: $a_i = [\tau_i - \rho(1 - \rho^{\tau_i}) / (1 - \rho)] / (1 - \rho)$
- 12 Compute normalized gradient: $\mathbf{d}_i^{(t)} = (\mathbf{x}^{(t,0)} - \mathbf{x}^{(t,\tau_i)}) / (\eta a_i)$
- 13 **Communication:** Send $p_i a_i$ and $p_i \mathbf{d}_i^{(t)}$ to the server
- 14 **end**

I More Experiments Details

Platform. All experiments in this paper are conducted on a cluster of 16 machines, each of which is equipped with one NVIDIA TitanX GPU. The machines communicate (*i.e.*, transfer model parameters) with each other via Ethernet. We treat each machine as one client in the federated learning setting. The algorithms are implemented by PyTorch. We run each experiments for 3 times with different random seeds.

Hyper-parameter Choices. On non-IID CIFAR10 dataset, we fix the mini-batch size per client as 32. When clients use momentum SGD as the local solver, the momentum factor is 0.9; when clients use proximal SGD, the proximal parameter μ is selected from $\{0.0005, 0.001, 0.005, 0.01\}$. It turns out that when $E_i = 2$, $\mu = 0.005$ is the best and when $E_i(t) \sim \mathcal{U}(2, 5)$, $\mu = 0.001$ is the best. The client learning rate η is tuned from $\{0.005, 0.01, 0.02, 0.05, 0.08\}$ for FedAvg with each local solver separately. When using the same local solver, FedNova uses the same client learning rate as FedAvg.

Specifically, if the local solver is momentum SGD, then we set $\eta = 0.02$. In other cases, $\eta = 0.05$ consistently performs the best. On the synthetic dataset, the mini-batch size per client is 20 and the client learning rate is 0.02.

Training Curves on Non-IID CIFAR10. The training curves of FedAvg and FedNova are presented in Figure 6. Observe that FedNova (red curve) outperforms FedAvg (blue curve) by a large margin. FedNova only requires about half of the total rounds to achieve the same test accuracy as FedAvg. Besides, note that in [54], the test accuracy of FedAvg is higher than ours. This is because the authors of [54] let clients to perform 20 local epochs per round, which is 10 times more than our setting. In [54], after 100 communication rounds, FedAvg equivalently runs $100 \times 20 = 2000$ epochs.

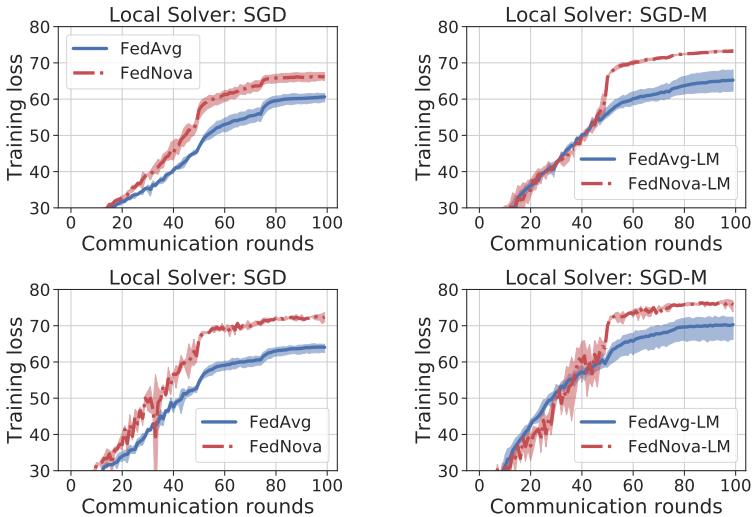


Figure 6: Training curves on non-IID partitioned CIFAR10 dataset. In these curves, the only difference between FedAvg and FedNova is the weights when aggregating normalized gradients. ‘LM’ represents for local momentum. **First row**: All clients perform $E_i = 2$ local epochs; **Second row**: All clients perform random and time-varying local epochs $E_i(t) \sim \mathcal{U}(2, 5)$.

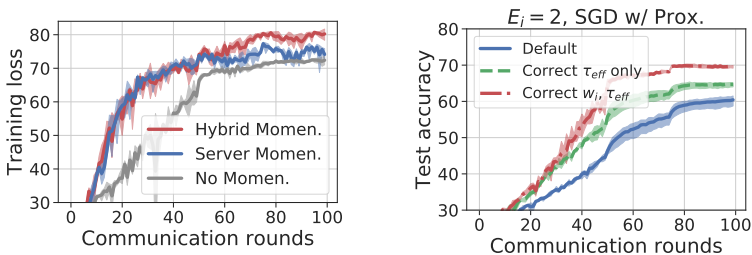


Figure 7: **Left**: Comparison of different momentum schemes in FedNova. ‘Hybrid momentum’ corresponds to the combination of server momentum and client momentum. **Right**: How FedNova-prox outperform vanilla FedProx (blue curve). By setting $\tau_{\text{eff}} = \sum_{i=1}^m p_i \tau_i$ instead of its default value, the accuracy of FedProx can be improved by 5% (see the green curve). By further correcting the aggregated weights, FedNova-prox (red curves) achieves around 10% higher accuracy than FedProx.