

## APPENDIX A

*Proof of lemma 1:* Firstly, we analyze the impact of quantized components on bitstream transmission. We omit the user subscript  $m$  when it does not cause ambiguity. Let  $\delta_i$  be represented as a bit sequence  $[b_1, b_2, \dots, b_{s(b)}]$ , where the error probability of each bit is  $P_e^1, P_e^2, \dots, P_e^{s(b)}$ . Then, the transmission error probability of  $q(\delta_i^m)$  satisfies

$$\begin{aligned}
 & \mathbb{E} \left[ (\hat{q}(\delta_i^m) - q(\delta_i^m))^2 \right] \\
 &= \mathbb{E} \left[ \mathbb{E} \left[ (\hat{q}(\delta_i^m) - q(\delta_i^m))^2 \middle| e_{b_1} \right] \right] \\
 &\leq P_e^1 \cdot (1 - (-1))^2 + \\
 &\quad (1 - P_e^1) \mathbb{E} \left[ (\hat{q}(\delta_i^m) - q(\delta_i^m))^2 \middle| e_{b_1} = 0 \right] \\
 &= P_e^1 \cdot 2^2 + \\
 &\quad (1 - P_e^1) \mathbb{E} \left[ \mathbb{E} \left[ (\hat{q}(\delta_i^m) - q(\delta_i^m))^2 \middle| e_{b_1} = 0, e_{b_2} \right] \right] \\
 &\leq P_e^1 \cdot 2^2 + (1 - P_e^1) (P_e^2 \cdot (2^{-1})^2 + \\
 &\quad \mathbb{E} \left[ (\hat{q}(\delta_i^m) - q(\delta_i^m))^2 \middle| e_{b_1} = 0, e_{b_2} = 0 \right]) \\
 &\leq P_e^1 \cdot 2^2 + (1 - P_e^1) \cdot \sum_{i=2}^{s(b)} P_e^i \cdot 2^{-2(i-1)}, \quad (22)
 \end{aligned}$$

where  $e_{b_i} = 1$  indicates that  $b_i$  bit is transmitted incorrectly.

Next, we analyze the bitstream transmission error of the parameter update vector  $\delta_t$  during the  $t$ -th round of upload

$$\begin{aligned}
 & \mathbb{E} \left\| \hat{\delta}_t - \delta_t \right\|^2 \\
 &= \mathbb{E} \left[ \left( \hat{\delta}_t - Q(\delta_t) \right) + \left( Q(\delta_t) - \delta_t \right) \right]^2 \\
 &= \mathbb{E} \left[ \left\| \hat{\delta}_t - Q(\delta_t) \right\|^2 \right] + \mathbb{E} \left[ \left\| Q(\delta_t) - \delta_t \right\|^2 \right], \quad (23)
 \end{aligned}$$

where (23) comes from the unbiasedness of the quantizer. For the first term in (23), we have

$$\begin{aligned}
 \mathbb{E} \left[ \left\| \hat{\delta}_t - Q(\delta_t) \right\|^2 \right] &= \|\delta_t\|^2 \mathbb{E} \left[ \|\hat{q}(\delta_t) - q(\delta_t)\|^2 \right] \\
 &= \|\delta_t\|^2 \sum_{i=1}^d \mathbb{E} \left[ (\hat{q}(\delta_{t,i}) - q(\delta_{t,i}))^2 \right] \\
 &\leq \|\delta_t\|^2 d P_e, \quad (24)
 \end{aligned}$$

where (24) comes from (22). For the second term in (23), we have [11]

$$\mathbb{E} \left[ \left\| Q(\delta_t) - \delta_t \right\|^2 \right] \leq \|\delta_t\|^2 \min \left( \frac{d}{s^2}, \frac{\sqrt{d}}{s} \right), \quad (25)$$

Plugging (24) and (25) into (23), we can obtain the result. ■

## APPENDIX B

*Proof of lemma 2:* After the  $t$ -th communication, the local stochastic gradient of the client  $m$  at round  $\tau = \{0, 1, \dots, K\}$  is  $\nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m)$ , and we have

$$\begin{aligned}
 f(\tilde{\mathbf{w}}_{t+1}) &= f \left( \sum_{m \in [M]} p_m \mathbf{w}_{t+1}^m \right) \\
 &= f \left( \sum_{m \in [M]} p_m \left( \mathbf{w}_t - \eta \sum_{\tau=0}^{K-1} \nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m) \right) \right) \\
 &= f \left( \mathbf{w}_t - \eta \sum_{m \in [M]} \sum_{\tau=0}^{K-1} p_m \nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m) \right) \\
 &\leq f(\mathbf{w}_t) - \left\langle \nabla f(\mathbf{w}_t), \eta \sum_{m \in [M]} \sum_{\tau=0}^{K-1} p_m \nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m) \right\rangle \\
 &\quad + \frac{L\eta^2}{2} \left\| \sum_{m \in [M]} \sum_{\tau=0}^{K-1} p_m \nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m) \right\|^2, \quad (26)
 \end{aligned}$$

where the inequality comes from the L-smooth property.

We consider the second term first and take the expectation of  $\nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m)$ . After transforming the coefficient, we have

$$\begin{aligned}
 & \mathbb{E} \left\langle \nabla f(\mathbf{w}_t), \eta \sum_{m \in [M]} \sum_{\tau=0}^{K-1} p_m \nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m) \right\rangle \\
 &= \eta \sum_{\tau=0}^{K-1} \mathbb{E} \left\langle \nabla f(\mathbf{w}_t), \sum_{m \in [M]} p_m \nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m) \right\rangle. \quad (27)
 \end{aligned}$$

Each term inside the expectation satisfies

$$\begin{aligned}
 & \left\langle \nabla f(\mathbf{w}_t), \sum_{m \in [M]} p_m \nabla f_m(\mathbf{w}_{t,\tau}^m) \right\rangle \\
 &= \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 + \frac{1}{2} \left\| \sum_{m \in [M]} p_m \nabla f_m(\mathbf{w}_{t,\tau}^m) \right\|^2 \\
 &\quad - \frac{1}{2} \left\| \nabla f(\mathbf{w}_t) - \sum_{m \in [M]} p_m \nabla f_m(\mathbf{w}_{t,\tau}^m) \right\|^2 \\
 &\geq \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 - \frac{1}{2} \left\| \nabla f(\mathbf{w}_t) - \sum_{m \in [M]} p_m \nabla f_m(\mathbf{w}_{t,\tau}^m) \right\|^2 \\
 &= \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 - \\
 &\quad \frac{1}{2} \left\| \sum_{m \in [M]} p_m (\nabla f_m(\mathbf{w}_t) - \nabla f_m(\mathbf{w}_{t,\tau}^m)) \right\|^2
 \end{aligned}$$

$$\begin{aligned}
&\geq \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 - \\
&\quad \frac{1}{2} M \sum_{m \in [M]} p_m^2 \|\nabla f_m(\mathbf{w}_t) - \nabla f_m(\mathbf{w}_{t,\tau}^m)\|^2 \quad (28) \\
&\geq \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 - \frac{1}{2} ML^2 \sum_{m \in [M]} p_m^2 \|\mathbf{w}_t - \mathbf{w}_{t,\tau}^m\|^2, \quad (29)
\end{aligned}$$

where (28) comes from the Cauchy inequality; (29) comes from the L-smooth condition in Assumption III-A. After taking expectation over a part of the second term on (29), we have

$$\begin{aligned}
&\mathbb{E} \left[ \|\mathbf{w}_t - \mathbf{w}_{t,\tau}^m\|^2 \right] \\
&= \mathbb{E} \left[ \left\| \eta \sum_{j=0}^{\tau-1} \nabla f_m(\mathbf{w}_{t,j}^m, \boldsymbol{\xi}_{t,j+1}^m) \right\|^2 \right] \\
&= \eta^2 \mathbb{E} \left[ \left\| \sum_{j=0}^{\tau-1} \nabla f_m(\mathbf{w}_{t,j}^m, \boldsymbol{\xi}_{t,j+1}^m) - \nabla f_m(\mathbf{w}_{t,j}^m) \right\|^2 \right] \\
&= \eta^2 \mathbb{E} \left[ \left\| \sum_{j=0}^{\tau-1} \nabla f_m(\mathbf{w}_{t,j}^m, \boldsymbol{\xi}_{t,j+1}^m) - \nabla f_m(\mathbf{w}_{t,j}^m) \right\|^2 \right] + \\
&\quad \eta^2 \mathbb{E} \left[ \left\| \sum_{j=0}^{\tau-1} \nabla f_m(\mathbf{w}_{t,j}^m) \right\|^2 \right] \\
&\leq \eta^2 \tau^2 \sigma^2 + \eta^2 \mathbb{E} \left[ \left\| \sum_{j=0}^{\tau-1} \nabla f_m(\mathbf{w}_{t,j}^m) \right\|^2 \right] \quad (30) \\
&= \eta^2 \tau^2 \sigma^2 + \eta^2 \mathbb{E} \left[ \left\| \sum_{j=0}^{\tau-1} \nabla f_m(\mathbf{w}_{t,j}^m) - \nabla f_m(\mathbf{w}_t) + \nabla f_m(\mathbf{w}_t) - \nabla f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t) \right\|^2 \right] \\
&\leq \eta^2 \tau^2 \sigma^2 + 3\eta^2 \left( \mathbb{E} \left[ \left\| \sum_{j=0}^{\tau-1} \nabla f_m(\mathbf{w}_{t,j}^m) - \nabla f_m(\mathbf{w}_t) \right\|^2 \right] \right. \\
&\quad \left. + \mathbb{E} \left[ \left\| \sum_{j=0}^{\tau-1} \nabla f_m(\mathbf{w}_t) - \nabla f(\mathbf{w}_t) \right\|^2 \right] + \mathbb{E} \left[ \left\| \sum_{j=0}^{\tau-1} \nabla f(\mathbf{w}_t) \right\|^2 \right] \right) \quad (31) \\
&\leq \eta^2 \tau^2 \sigma^2 + 3\eta^2 \left( L^2 \tau \sum_{j=0}^{\tau-1} \mathbb{E} [\|\mathbf{w}_t - \mathbf{w}_{t,j}^m\|^2] + \right. \\
&\quad \left. \tau^2 G_t^2 + \tau^2 \mathbb{E} [\|\nabla f(\mathbf{w}_t)\|^2] \right) \quad (32) \\
&\leq \eta^2 \tau^2 \sigma^2 + 3\eta^2 \left( L^2 \tau A_t + \tau^2 G_t^2 + \tau^2 \mathbb{E} [\|\nabla f(\mathbf{w}_t)\|^2] \right). \quad (33)
\end{aligned}$$

Here, (30) is derived from the gradient variance bound in hypothesis III-A, (31) is derived from Cauchy's inequality,

(32) is derived from the heterogeneous data property in hypothesis III-A, and (33) is derived from the definition  $A_t \triangleq \sum_{\tau=0}^{K-1} \mathbb{E} [\|\mathbf{w}_t - \mathbf{w}_{t,\tau}^m\|^2]$ . Adding both sides of inequality (33), we get

$$\begin{aligned}
A_t &\leq \sum_{\tau=0}^{K-1} \eta^2 \tau^2 \sigma^2 + 3\eta^2 \cdot \\
&\quad \left( L^2 \tau A_t + \tau^2 G_t^2 + \tau^2 \mathbb{E} [\|\nabla f(\mathbf{w}_t)\|^2] \right) \\
&= \eta^2 \left( \frac{K(K-1)(2K-1)}{6} (\sigma^2 + 3G_t^2 + \right. \\
&\quad \left. 3\mathbb{E} [\|\nabla f(\mathbf{w}_t)\|^2]) + \frac{3K(K-1)}{2} L^2 A_t \right). \quad (34)
\end{aligned}$$

Set  $0 < \eta < \frac{1}{\sqrt{3KL}}$ , then we have

$$\begin{aligned}
A_t &\leq \frac{\eta^2}{1 - \frac{3K(K-1)}{2} L^2 \eta^2} \frac{K(K-1)(2K-1)}{6} (\sigma^2 + \\
&\quad 3G_t^2 + 3\mathbb{E} [\|\nabla f(\mathbf{w}_t)\|^2]) \\
&\leq \frac{2\eta^2 K^3}{3(2 - 3K^2 L^2 \eta^2)} (\sigma^2 + 3G_t^2 + 3\mathbb{E} [\|\nabla f(\mathbf{w}_t)\|^2]).
\end{aligned}$$

Substituting the above results into (27), we get

$$\begin{aligned}
&\mathbb{E} \left\langle \nabla f(\mathbf{w}_t), \eta \sum_{m \in [M]} \sum_{\tau=0}^{K-1} p_m \nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m) \right\rangle \\
&= \eta \sum_{\tau=0}^{K-1} \mathbb{E} \left\langle \nabla f(\mathbf{w}_t), \sum_{m \in [M]} p_m \nabla f_m(\mathbf{w}_{t,\tau+1}^m) \right\rangle \\
&\geq \eta \sum_{\tau=0}^{K-1} \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 - \frac{1}{2} ML^2 \sum_{m \in [M]} p_m^2 \|\mathbf{w}_t - \mathbf{w}_{t,\tau}^m\|^2 \\
&= \frac{\eta K}{2} \|\nabla f(\mathbf{w}_t)\|^2 - \frac{ML^2}{2} \bar{p} A_t, \quad (35)
\end{aligned}$$

where  $\bar{p} \triangleq \sum_{m \in [M]} p_m^2$ .

Then we consider the third term of (26) and obtain

$$\begin{aligned}
&\frac{L\eta^2}{2} \left\| \sum_{m \in [M]} \sum_{\tau=0}^{K-1} p_m \nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m) \right\|^2 \\
&= \frac{L\eta^2}{2} \left\| \sum_{m \in [M]} p_m \sum_{\tau=0}^{K-1} \nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m) \right\|^2 \\
&\leq \frac{L\eta^2 M}{2} \sum_{m \in [M]} p_m^2 \left\| \sum_{\tau=0}^{K-1} \nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m) \right\|^2 \\
&\leq \frac{LM}{2} \bar{p} A_t. \quad (36)
\end{aligned}$$

Substituting (36) and (35) into (26), we get

$$\begin{aligned}
&\mathbb{E} f(\tilde{\mathbf{w}}_{t+1}) \\
&\leq \mathbb{E} f(\mathbf{w}_t) - \left( \frac{\eta K}{2} \|\nabla f(\mathbf{w}_t)\|^2 - \frac{ML^2}{2} \bar{p} A_t \right) + \frac{LM}{2} \bar{p} A_t
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E}f(\mathbf{w}_t) - \frac{\eta K}{2} \|\nabla f(\mathbf{w}_t)\|^2 + \frac{\bar{p}ML(L+1)}{2} A_t \\
&\leq \mathbb{E}f(\mathbf{w}_t) - \frac{\eta K}{2} \|\nabla f(\mathbf{w}_t)\|^2 + \frac{\eta^2 \bar{p}ML(L+1)K^3}{3(2-3K^2L^2\eta^2)} \\
&\quad \left( \sigma^2 + 3G_t^2 + 3\mathbb{E}[\|\nabla f(\mathbf{w}_t)\|^2] \right) \\
&= \mathbb{E}f(\mathbf{w}_t) + \left( \frac{\eta^2 \bar{p}ML(L+1)K^3}{2-3K^2L^2\eta^2} - \frac{\eta K}{2} \right) \|\nabla f(\mathbf{w}_t)\|^2 \\
&\quad + \frac{\eta^2 \bar{p}ML(L+1)K^3}{3(2-3K^2L^2\eta^2)} (\sigma^2 + 3G_t^2).
\end{aligned}$$

#### APPENDIX C

*Proof of Lemma 3:*

$$\begin{aligned}
&\mathbb{E}\|\mathbf{w}_{t+1} - \tilde{\mathbf{w}}_{t+1}\|^2 \\
&= \mathbb{E}\left\| \mathbf{w}_t + \sum_{m \in [M]} p_m \hat{\delta}_t^m - \sum_{m \in [M]} p_m \mathbf{w}_{t+1}^m \right\|^2 \\
&= \mathbb{E}\left\| \sum_{m \in [M]} p_m [\hat{\delta}_t^m - (\mathbf{w}_{t+1}^m - \mathbf{w}_t)] \right\|^2 \\
&\leq M \sum_{m \in [M]} p_m^2 \mathbb{E}\left[\|\hat{\delta}_t^m - \delta_t^m\|^2\right] \\
&\leq M \sum_{m \in [M]} p_m^2 \left\| \eta \sum_{\tau=0}^{K-1} \nabla f_m(\mathbf{w}_{t,\tau}^m, \boldsymbol{\xi}_{t,\tau+1}^m) \right\|^2 \alpha_t \quad (37) \\
&\leq M \bar{p} A_t \alpha_t \quad (38) \\
&\leq \frac{2\eta^2 K^3 M \bar{p} \alpha_t}{3(2-3K^2L^2\eta^2)} (\sigma^2 + 3G_t^2 + 3\mathbb{E}[\|\nabla f(\mathbf{w}_t)\|^2]).
\end{aligned}$$

Here, (37) is derived from Lemma 1 and (38) is derived from the definition of  $A_t$ .  $\blacksquare$

#### APPENDIX D

*Proof of Theorem 1:* We use the L-smooth condition to reveal the effect of quantized transmission in the training process:

$$\begin{aligned}
\mathbb{E}f(\mathbf{w}_{t+1}) &= \mathbb{E}f(\tilde{\mathbf{w}}_{t+1} + \mathbf{w}_{t+1} - \tilde{\mathbf{w}}_{t+1}) \\
&\leq \mathbb{E}[f(\tilde{\mathbf{w}}_{t+1}) + \langle \nabla f(\tilde{\mathbf{w}}_{t+1}), \mathbf{w}_{t+1} - \tilde{\mathbf{w}}_{t+1} \rangle \\
&\quad + \frac{L}{2} \|\mathbf{w}_{t+1} - \tilde{\mathbf{w}}_{t+1}\|^2] \quad (39)
\end{aligned}$$

$$\approx \mathbb{E}f(\tilde{\mathbf{w}}_{t+1}) + \frac{L}{2} \mathbb{E}\|\mathbf{w}_{t+1} - \tilde{\mathbf{w}}_{t+1}\|^2, \quad (40)$$

where (39) comes from the quadratic upper bound of L-smooth functions  $f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$ ; (40) comes from the nearly unbiasedness of quantized transmission in Lemma 1. Using Lemma 2 and Lemma 3, we get

$$\begin{aligned}
&\mathbb{E}f(\mathbf{w}_{t+1}) \\
&\leq \mathbb{E}[f(\mathbf{w}_t)] + \left( \frac{\eta^2 \bar{p}ML(L+1)K^3}{2-3K^2L^2\eta^2} - \frac{\eta K}{2} \right).
\end{aligned}$$

$$\begin{aligned}
&\mathbb{E}[\|\nabla f(\mathbf{w}_t)\|^2] + \frac{\eta^2 \bar{p}ML(L+1)K^3}{3(2-3K^2L^2\eta^2)} (\sigma^2 + 3G_t^2) \\
&\quad + \frac{2\eta^2 K^3 M \bar{p} \alpha_t}{3(2-3K^2L^2\eta^2)} (\sigma^2 + 3G_t^2 + 3\mathbb{E}[\|\nabla f(\mathbf{w}_t)\|^2]) \\
&= \mathbb{E}[f(\mathbf{w}_t)] + \left( \frac{MK^3 \bar{p} \eta^2}{2-3K^2L^2\eta^2} \right) (L^2 + L + 2\alpha_t) - \\
&\quad \frac{\eta K}{2} \mathbb{E}[\|\nabla f(\mathbf{w}_t)\|^2] + \frac{MK^3 \bar{p} \eta^2}{3(2-3K^2L^2\eta^2)} \\
&\quad (L^2 + L + 2\alpha_t) (\sigma^2 + 3G_t^2) \\
&= \mathbb{E}[f(\mathbf{w}_t)] + \left( \kappa (L^2 + L + 2\alpha_t) - \frac{\eta K}{2} \right) \cdot \\
&\quad \mathbb{E}[\|\nabla f(\mathbf{w}_t)\|^2] + \kappa (L^2 + L + 2\alpha_t) (\sigma^2 + 3G_t^2), \quad (41)
\end{aligned}$$

where (41) comes from the definition  $\kappa \triangleq \frac{MK^3 \bar{p} \eta^2}{3(2-3K^2L^2\eta^2)}$ .

Accumulate over  $t = 0, \dots, T-1$  and rearrange terms, yielding

$$\begin{aligned}
&\sum_{t=0}^{T-1} \left( \frac{\eta K}{2} - \kappa (L^2 + L + 2\alpha_t) \right) \mathbb{E}[\|\nabla f(\mathbf{w}_t)\|^2] \\
&\leq f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T) + \sum_{t=0}^{T-1} \kappa (L^2 + L + 2\alpha_t) (\sigma^2 + 3G_t^2).
\end{aligned}$$

We next study the lower bound of the coefficient of  $\mathbb{E}[\|\nabla f(\mathbf{w}_t)\|^2]$ . Take  $\eta$  such that  $0 < \eta < \frac{3}{4MK^2 \bar{p}(L^2 + L + 4d + E)}$ , then we have

$$\begin{aligned}
&\frac{\eta K}{2} - \kappa (L^2 + L + 2\alpha_t) \\
&= \frac{\eta K}{2} - \frac{MK^3 \bar{p} \eta^2}{3(2-3K^2L^2\eta^2)} (L^2 + L + 2\alpha_t) \quad (42)
\end{aligned}$$

$$\begin{aligned}
&\geq \frac{\eta K}{2} - \frac{MK^3 \bar{p} \eta^2}{3(2-3K^2L^2\eta^2)} (L^2 + L + 4d + E) \\
&> \frac{\eta K}{2} - \frac{MK^3 \bar{p} \eta^2}{3} (L^2 + L + 4d + E) \quad (43)
\end{aligned}$$

$$> \eta K \left( \frac{1}{2} - \frac{1}{4} \right) \quad (44)$$

$$= \frac{\eta K}{4}, \quad (45)$$

where (42) comes from the upper bound of transmission error  $\alpha_t$ , i.e.,

$$\alpha_t = dD_t + E \leq 4d + E. \quad (46)$$

(43) and (44) come from the constraint on the value range of  $\eta$ . Therefore, we have

$$\begin{aligned}
&\frac{\eta K}{4} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{w}_t)\|^2] \\
&< \sum_{t=0}^{T-1} \left( \frac{\eta K}{2} - \kappa (L^2 + L + 2\alpha_t) \right) \mathbb{E}[\|\nabla f(\mathbf{w}_t)\|^2] \\
&\leq f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T) + \sum_{t=0}^{T-1} \kappa (L^2 + L + 2\alpha_t) (\sigma^2 + 3G_t^2).
\end{aligned}$$

Rearranging, we get

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \left\| \nabla f(\mathbf{w}_t) \right\|^2 \right] \leq \frac{4}{\eta TK} (f(\mathbf{w}_0) - \mathbb{E} f(\mathbf{w}_T) + \sum_{t=0}^{T-1} \kappa (L^2 + L + 2\alpha_t) (\sigma^2 + 3G_t^2)) .$$

■