

Appendix

Proof of Eq. 2

Proof. Let Δ denote the local update before clipping and adding noise, $\tilde{\Delta}_i^t = \Delta_i^t \cdot \min(1, \frac{S}{\|\Delta_i^t\|})$ denote the local update before clipping but after adding noise, $\bar{\Delta}_i^t = \tilde{\Delta}_i^t + \mathcal{N}(0, \frac{\sigma^2 S^2}{|\mathcal{P}_t|})$ denote the local update after clipping and adding noise. Then we have

$$\begin{aligned}
& \mathbb{E} \left[\frac{1}{d} \left\| \bar{\Delta}_i^t - \Delta_i^t \right\|_2^2 \right] \\
& \leq \frac{1}{d} \mathbb{E} \left[\left\| \tilde{\Delta}_i^t - \Delta_i^t \right\|_2^2 \right] + \frac{1}{d} \mathbb{E} \left[\left\| \bar{\Delta}_i^t - \tilde{\Delta}_i^t \right\|_2^2 \right] \\
& = \frac{1}{d} \mathbb{E} \left[\left\| \Delta_i^t \cdot \min(1, \frac{S}{\|\Delta_i^t\|}) - \Delta_i^t \right\|_2^2 \right] + \frac{1}{d} \mathbb{E} \left[\left\| \mathcal{N}(0, \frac{\sigma^2 S^2}{|\mathcal{P}_t|} \mathbf{I}_d) \right\|_2^2 \right] \\
& = \frac{1}{d} \mathbb{E} \left[\left\| \frac{\Delta_i^t}{\|\Delta_i^t\|} \cdot \min(0, S - \|\Delta_i^t\|) \right\|_2^2 \right] + \frac{1}{d} \mathbb{E} \left[d \cdot \mathcal{N}(0, \frac{\sigma^2 S^2}{|\mathcal{P}_t|})^2 \right] \\
& = \frac{1}{d} \min(0, S - \|\Delta_i^t\|)^2 + \mathbb{E} \left[\mathcal{N}(0, \frac{\sigma^2 S^2}{|\mathcal{P}_t|})^2 \right] \\
& = \frac{1}{d} \max(0, \|\Delta_i^t\| - S)^2 + \frac{\sigma^2 S^2}{|\mathcal{P}_t|} \quad (14)
\end{aligned}$$

Proof of Lemma 2

Proof. According to Eq. 5 and Eq. 7, the local update at step Q of agent i at round t can be written as

$$\begin{aligned}
\mathbf{w}_i^{t,Q} &= \mathbf{w}_i^{t,Q-1} - \eta_l \left(\mathbf{g}_i^{t,Q} + \lambda(\mathbf{w}_i^{t,Q-1} - \mathbf{w}^t) \right) \\
&= (1 - \lambda\eta_l) \mathbf{w}_i^{t,Q-1} - \eta_l \mathbf{g}_i^{t,Q} + \lambda\eta_l \mathbf{w}^t \quad (15)
\end{aligned}$$

where $\mathbf{g}_i^{t,q} = \frac{1}{|\mathcal{B}_q|} \sum_{(x,y) \in \mathcal{B}_q} \nabla f_i(\mathbf{w}_i^{t,q-1}, x, y)$. Unrolling Eq. 15, we have

$$\begin{aligned}
\mathbf{w}_i^{t,Q} &= -\eta \sum_{q=0}^{Q-1} \left[(1 - \lambda\eta_l)^q \mathbf{g}_i^{t,q} \right] \\
&\quad + \underbrace{\left[(1 - \lambda\eta_l)^Q + \sum_{q=0}^{Q-1} (1 - \lambda\eta_l)^q \lambda\eta_l \right]}_{A_1} \cdot \mathbf{w}^t \quad (16)
\end{aligned}$$

Simplifying A_1 , we have

$$\begin{aligned}
A_1 &= (1 - \lambda\eta_l)^Q + \sum_{q=0}^{Q-1} (1 - \lambda\eta_l)^q \lambda\eta_l \\
&= (1 - \lambda\eta_l)^Q + (1 - \lambda\eta_l)^{Q-1} \lambda\eta_l + \sum_{q=0}^{Q-2} (1 - \lambda\eta_l)^q \lambda\eta_l \\
&= (1 - \lambda\eta_l)^{Q-1} (1 - \lambda\eta_l + \lambda\eta_l) + \sum_{q=0}^{Q-2} (1 - \lambda\eta_l)^q \lambda\eta_l \\
&= (1 - \lambda\eta_l)^{Q-1} + \sum_{q=0}^{Q-2} (1 - \lambda\eta_l)^q \lambda\eta_l \\
&= (1 - \lambda\eta_l)^{Q-2} + \sum_{q=0}^{Q-3} (1 - \lambda\eta_l)^q \lambda\eta_l \\
&\dots \\
&= (1 - \lambda\eta_l)^1 + \sum_{q=0}^0 (1 - \lambda\eta_l)^q \lambda\eta_l \\
&= 1 \quad (17)
\end{aligned}$$

Replacing A_1 in Eq. 16 yields the desired result.

Proof of Theorem 1

Our Theorem 1 is directly obtained by applying [[2], Theorem 1] to the global learning process of our protocol, with sampling probability $q = \frac{P}{N}$ and global steps T . Referring to [2] for more details.

Proof of Theorem 2

Proof. For convenience, we first define the following notations:

$$\begin{aligned}
\Delta_i^t &:= -\eta_l \sum_{q=0}^{Q-1} \gamma_i^{t,q} \mathbf{g}_i^{t,q} \cdot \alpha_i^t, \\
\bar{\Delta}_i^t &:= -\eta_l \sum_{q=0}^{Q-1} \gamma_i^{t,q} \mathbf{g}_i^{t,q} \cdot \bar{\alpha}^t, \\
\check{\Delta}_i^t &:= -\eta_l \sum_{q=0}^{Q-1} \gamma_i^{t,q} \nabla f_i(x_i^{t,q}) \cdot \bar{\alpha}^t
\end{aligned}$$

where

$$\begin{aligned}
\alpha_i^t &:= \min \left(1, \frac{S}{\eta_l \beta_i^t \left\| \sum_{q=0}^{Q-1} \gamma_i^{t,q} \mathbf{g}_i^{t,q} \right\|} \right), \bar{\alpha}^t := \frac{1}{N} \sum_{i=1}^N \alpha_i^t, \\
\beta_i^t &= \frac{\|M_i^t \circ \sum_{q=0}^{Q-1} \gamma_i^{t,q} \mathbf{g}_i^{t,q}\|}{\left\| \sum_{q=0}^{Q-1} \gamma_i^{t,q} \mathbf{g}_i^{t,q} \right\|}
\end{aligned}$$

By Lipschitz smoothness, we have

$$\begin{aligned}
& \mathbb{E}[f(x_{t+1})] \\
& \leq f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{L}{2} \|x_{t+1} - x_t\|^2 \\
& = f(x_t) + \eta_g \left\langle \nabla f(x_t), \mathbb{E} \left[\frac{1}{P} \sum_{i \in \mathcal{P}_t} \Delta_i^t + z_i^t \right] \right\rangle \\
& \quad + \frac{L}{2} \eta_g^2 \mathbb{E} \left[\left\| \frac{1}{P} \sum_{i \in \mathcal{P}_t} \Delta_i^t + z_i^t \right\|^2 \right] \\
& = f(x_t) + \eta_g \left\langle \nabla f(x_t), \mathbb{E} \left[\frac{1}{P} \sum_{i \in \mathcal{P}_t} \Delta_i^t \right] \right\rangle \\
& \quad + \frac{L}{2} \eta_g^2 \mathbb{E} \left[\left\| \frac{1}{P} \sum_{i \in \mathcal{P}_t} \Delta_i^t \right\|^2 \right] + \frac{\eta_g^2 L \sigma^2 S^2 d}{2P^2} \quad (18)
\end{aligned}$$

where d represents dimension of $x_i^{t,q}$; in the last equation we use the fact that z_i^t is zero mean. Next, we will analyse the bias caused by clipping, through analyzing the first order term in the above expression. Towards this end, we have the following series of relations:

$$\begin{aligned}
& \left\langle \nabla f(x_t), \mathbb{E} \left[\frac{1}{P} \sum_{i \in \mathcal{P}_t} \Delta_i^t \right] \right\rangle \\
& = \left\langle \nabla f(x_t), \mathbb{E} \left[\frac{1}{P} \mathbb{E}_i \left[\sum_{i \in \mathcal{P}_t} \Delta_i^t \right] \right] \right\rangle \\
& = \left\langle \nabla f(x_t), \frac{1}{P} P \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \Delta_i^t \right] \right\rangle \\
& = \left\langle \nabla f(x_t), \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \Delta_i^t - \bar{\Delta}_i^t \right] \right\rangle \\
& \quad + \left\langle \nabla f(x_t), \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \bar{\Delta}_i^t \right] \right\rangle \quad (19)
\end{aligned}$$

Then we bound the two terms in the above expression,

respectively. To bound the first term, we have:

$$\begin{aligned}
& \left\langle \nabla f(x_t), \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \Delta_i^t - \bar{\Delta}_i^t \right] \right\rangle \\
& = \left\langle \nabla f(x_t), \mathbb{E} \left[-\frac{1}{N} \sum_{i=1}^N \sum_{q=0}^{Q-1} \eta_l (\alpha_i^t - \bar{\alpha}^t) \gamma_i^{t,q} \mathbf{g}_i^{t,q} \right] \right\rangle \\
& \leq \left\langle \nabla f(x_t), \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \sum_{q=0}^{Q-1} \eta_l |\alpha_i^t - \bar{\alpha}^t| \gamma_i^{t,q} \mathbf{g}_i^{t,q} \right] \right\rangle \\
& \leq \left\langle \nabla f(x_t), \eta_l Q \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N |\alpha_i^t - \bar{\alpha}^t| \mathbf{g}_i^{t,q} \right] \right\rangle \\
& = \frac{\eta_l Q}{N} \sum_{i=1}^N |\alpha_i^t - \bar{\alpha}^t| \langle \nabla f(x_t), \mathbf{g}_i^{t,q} \rangle \\
& = \frac{\eta_l Q}{N} \sum_{i=1}^N |\alpha_i^t - \bar{\alpha}^t| \left\langle \frac{1}{N} \sum_{j=1}^N \nabla f_j(x_t), \mathbf{g}_i^{t,q} \right\rangle \\
& = \frac{\eta_l Q}{N} \sum_{i=1}^N |\alpha_i^t - \bar{\alpha}^t| \left(\frac{1}{N} \langle \nabla f_i(x_t), \mathbf{g}_i^{t,q} \rangle \right) \\
& \leq \frac{\eta_l Q G^2}{N} \sum_{i=1}^N |\alpha_i^t - \bar{\alpha}^t| \\
& = \eta_l \tilde{\alpha}^t Q G^2 \quad (20)
\end{aligned}$$

where $\tilde{\alpha}^t = \frac{1}{N} \sum_{i=1}^N |\alpha_i^t - \bar{\alpha}^t|$. To bound the second term, we have:

$$\begin{aligned}
& \left\langle \nabla f(x_t), \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \bar{\Delta}_i^t \right] \right\rangle = \mathbb{E} \left[\left\langle \nabla f(x_t), \frac{1}{N} \sum_{i=1}^N \check{\Delta}_i^t \right\rangle \right] \\
& = \frac{-\eta_l \bar{\alpha}^t Q}{2} \|\nabla f(x_t)\|^2 - \frac{\eta_l \bar{\alpha}^t}{2Q} \mathbb{E} \left[\left\| \frac{1}{\eta_l \bar{\alpha}^t N} \sum_{i=1}^N \check{\Delta}_i^t \right\|^2 \right] \\
& \quad + \underbrace{\frac{\eta_l \bar{\alpha}^t}{2} \mathbb{E} \left[\left\| \sqrt{Q} \nabla f(x_t) - \frac{1}{\eta_l \bar{\alpha}^t N \sqrt{Q}} \sum_{i=1}^N \check{\Delta}_i^t \right\|^2 \right]}_{A_2} \quad (21)
\end{aligned}$$

where the second equation is because $\langle a, b \rangle = -\frac{1}{2} \|a\|^2 - \frac{1}{2} \|b\|^2 + \frac{1}{2} \|a - b\|^2$ holds true for any vector a and b . Next,

we bound A_2 as follows:

$$\begin{aligned}
A &= Q\mathbb{E} \left[\left\| \nabla f(x_t) - \frac{1}{QN} \sum_{i=1}^N \sum_{q=0}^{Q-1} \gamma_i^{t,q} \nabla f_i(x_i^{t,q}) \right\|^2 \right] \\
&= Q\mathbb{E} \left[\left\| \frac{1}{QN} \sum_{i=1}^N \sum_{q=0}^{Q-1} \nabla f_i(x^t) - \gamma_i^{t,q} \nabla f_i(x_i^{t,q}) \right\|^2 \right] \\
&\leq \frac{1}{N} \sum_{i=1}^N \sum_{q=0}^{Q-1} \mathbb{E} \left[\left\| \nabla f_i(x^t) - \gamma_i^{t,q} \nabla f_i(x_i^{t,q}) \right\|^2 \right] \\
&\leq \frac{1}{N} \sum_{i=1}^N \sum_{q=0}^{Q-1} \mathbb{E} \left[\left\| \nabla f_i(x^t) - \nabla f_i(x_i^{t,q}) \right\|^2 \right] \\
&+ \frac{1}{N} \sum_{i=1}^N \sum_{q=0}^{Q-1} \mathbb{E} \left[\left\| \nabla f_i(x_i^{t,q}) - \gamma_i^{t,q} \nabla f_i(x_i^{t,q}) \right\|^2 \right] \\
&\leq \frac{1}{N} \sum_{i=1}^N \sum_{q=0}^{Q-1} \left(L^2 \mathbb{E} \left[\|x^t - x_i^{t,q}\|^2 \right] + \mathbb{E} \left[(1 - (\gamma_i^{t,q}))^2 G^2 \right] \right) \\
&\leq L^2 5Q^2 \eta_l^2 (\sigma_l^2 + 6Q\sigma_g^2) + L^2 30Q^3 \eta_l^2 \|\nabla f(x_t)\|^2 + QG^2 \quad (22)
\end{aligned}$$

where the first inequality comes from Jensen's inequality, the third inequality comes from L -smoothness, and the last inequality is due to [29, Lemma 3] which indicates that the following inequality holds for any q

$$\begin{aligned}
&\frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[\|x^t - x_i^{t,q}\|^2 \right] \\
&\leq 5Q\eta_l^2 (\sigma_l^2 + 6Q\sigma_g^2) + 30Q^2 \eta_l^2 \|\nabla f(x_t)\|^2 \quad (23)
\end{aligned}$$

Next, we turn to upper bounding the second order term in Eq. 18 as follows

$$\begin{aligned}
&\mathbb{E} \left[\left\| \frac{1}{P} \sum_{i \in \mathcal{P}_t} \Delta_i^t \right\|^2 \right] \\
&= \mathbb{E} \left[\left\| \frac{1}{P} \sum_{i \in \mathcal{P}_t} \sum_{q=0}^{Q-1} \eta_l \alpha_i^t \gamma_i^{t,q} \mathbf{g}_i^{t,q} \right\|^2 \right] \\
&= \mathbb{E} \left[\frac{1}{P} \sum_{i \in \mathcal{P}_t} \sum_{q=0}^{Q-1} (\eta_l \alpha_i^t \gamma_i^{t,q})^2 \|\mathbf{g}_i^{t,q}\|^2 \right] \\
&\leq \frac{\eta_l^2 G^2}{P^2} \mathbb{E} \left[\sum_{i \in \mathcal{P}_t} \sum_{q=0}^{Q-1} (\alpha_i^t \gamma_i^{t,q})^2 \right] \\
&\leq \frac{\eta_l^2 G^2 Q}{P} \quad (24)
\end{aligned}$$

where the first inequation is because the bounded gradient assumption.

Combining Eq. 18-24, we have

$$\begin{aligned}
\mathbb{E}[f(x_{t+1})] &\leq f(x_t) + \eta_g \eta_l \bar{\alpha}^t QG^2 \\
&- \frac{\eta_g \eta_l \bar{\alpha}^t Q}{2} \|\nabla f(x_t)\|^2 - \frac{\eta_g \eta_l \bar{\alpha}^t}{2Q} \mathbb{E} \left[\left\| \frac{1}{\eta_l N \bar{\alpha}^t} \sum_{i=1}^N \Delta_i^t \right\|^2 \right] \\
&+ \frac{\eta_g \eta_l \bar{\alpha}^t}{2} \left(5L^2 Q^2 \eta_l^2 (\sigma_l^2 + 6Q\sigma_g^2) + 30L^2 Q^3 \eta_l^2 \|\nabla f(x_t)\|^2 + QG^2 \right) \\
&+ \frac{\eta_l^2 \eta_g^2 LG^2 Q}{2P} + \frac{\eta_g^2 L \sigma^2 S^2 d}{2P^2} \quad (25)
\end{aligned}$$

When $\eta_l \leq \frac{1}{\sqrt{60QL}}$ the above inequality simplifies to

$$\begin{aligned}
\mathbb{E}[f(x_{t+1})] &\leq f(x_t) - \frac{\eta_g \eta_l \bar{\alpha}^t Q}{4} \|\nabla f(x_t)\|^2 + \eta_g \eta_l \bar{\alpha}^t QG^2 \\
&+ \frac{5\eta_g \eta_l^3 \bar{\alpha}^t L^2 Q^2}{2} (\sigma_l^2 + 6Q\sigma_g^2) + \frac{\eta_g \eta_l \bar{\alpha}^t QG^2}{2} \\
&+ \frac{\eta_l^2 \eta_g^2 LG^2 Q}{2P} + \frac{\eta_g^2 L \sigma^2 S^2 d}{2P^2} \quad (26)
\end{aligned}$$

Sum over t from 1 to T , divide both side by $\frac{\eta_g \eta_l QT}{4}$, and rearrange, we have

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\bar{\alpha}^t \|\nabla f(x_t)\|^2 \right] &\leq \frac{4}{\eta_g \eta_l QT} \mathbb{E}[f(x_1) - f(x_{T+1})] \\
&+ \frac{4G^2}{T} \sum_{t=1}^T \bar{\alpha}^t + \frac{10\eta_l^2 L^2 Q}{T} (\sigma_l^2 + 6Q\sigma_g^2) \sum_{t=1}^T \bar{\alpha}^t + \frac{2G^2}{T} \sum_{t=1}^T \bar{\alpha}^t \\
&+ \frac{2\eta_l \eta_g LG^2}{P} + \frac{2\eta_g L \sigma^2 S^2 d}{\eta_l QP^2} \quad (27)
\end{aligned}$$

As both $\frac{1}{T} \sum_{t=1}^T \bar{\alpha}^t$ and $\frac{1}{T} \sum_{t=1}^T \bar{\alpha}^t$ are bounded, the big- \mathcal{O} convergence about T, Q, P, η_l, η_g is

$$\begin{aligned}
&\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\bar{\alpha}^t \|\nabla f(x_t)\|^2 \right] \\
&\leq \mathcal{O} \left(\frac{1}{\eta_g \eta_l QT} + \eta_l^2 Q^2 + \frac{\eta_g \eta_l}{P} + \frac{\eta_g \sigma^2 S^2 d}{\eta_l QP^2} \right) \quad (28)
\end{aligned}$$

which yield the result.