

A. Proof of Theorem 1

Theorem 1 (Global Loss Reduction) With Assumption on Lipschitz property, for an arbitrary set of clients $C_t \subseteq C$ selected by the server in round t , the reduction of global loss $F(w)$ is bounded by:

$$\underbrace{F(w_{\text{fed}}^{t-1}) - F(w_{\text{fed}}^t)}_{\text{global loss reduction}} \geq \underbrace{\sum_{c \in C_t} \sum_{i=0}^{m-1} \sum_{(x,y) \in B_c} [-\alpha_c \|\nabla_w l(w_c^{t,i}, x, y)\|_2^2]}_{\text{term 1}} + \underbrace{\beta_c \langle \nabla_w l(w_c^{t,i}, x, y), \nabla_w F(w_{\text{fed}}^{t-1}) \rangle}_{\text{term 2}}, \quad (15)$$

where $\alpha_c = \frac{L\zeta_c^t}{2} \cdot \left(\frac{\eta}{|B_c|}\right)^2$ and $\beta_c = \zeta_c^t \cdot \left(\frac{\eta}{|B_c|}\right)$.

Proof. From the L -Lipschitz continuity of global loss function $F(w)$, we have

$$F(w_{\text{fed}}^t) - F(w_{\text{fed}}^{t-1}) \leq \langle \nabla_w F(w_{\text{fed}}^{t-1}), w_{\text{fed}}^t - w_{\text{fed}}^{t-1} \rangle + \frac{L}{2} \|w_{\text{fed}}^t - w_{\text{fed}}^{t-1}\|^2, \quad (16)$$

where the global model can be further decomposed into the weighted sum of the updated models of participating clients:

$$\begin{aligned} & (16) \\ &= \langle \nabla_w F(w_{\text{fed}}^{t-1}), \sum_{c \in C_t} \zeta_c^t w_c^{t,m} - w_{\text{fed}}^{t-1} \rangle + \frac{L}{2} \left\| \sum_{c \in C_t} \zeta_c^t w_c^{t,m} - w_{\text{fed}}^{t-1} \right\|^2 \\ &\leq \sum_{c \in C_t} \zeta_c^t \langle \nabla_w F(w_{\text{fed}}^{t-1}), w_c^{t,m} - w_{\text{fed}}^{t-1} \rangle + \frac{L}{2} \sum_{c \in C_t} \zeta_c^t \|w_c^{t,m} - w_{\text{fed}}^{t-1}\|^2. \end{aligned} \quad (17)$$

We can further express the locally updated model $w_c^{t,m}$ of each client $c \in C_t$ as the difference between the original global model w_{fed}^{t-1} and m local model updates:

$$\begin{aligned} (17) &\leq \sum_{c \in C_t} \zeta_c^t \left[\langle \nabla_w F(w_{\text{fed}}^{t-1}), -\eta \sum_{i=0}^{m-1} \nabla_w \tilde{F}_c(w_c^{t,i}) \rangle \right. \\ &\quad \left. + \frac{L}{2} \left\| \sum_{i=0}^{m-1} \eta \nabla_w \tilde{F}_c(w_c^{t,i}) \right\|^2 \right] \\ &= \sum_{c \in C_t} \zeta_c^t \left[\frac{L}{2} \left\| \sum_{i=0}^{m-1} \eta \nabla_w \tilde{F}_c(w_c^{t,i}) \right\|^2 \right. \\ &\quad \left. - \eta \sum_{i=0}^{m-1} \langle \nabla_w F(w_{\text{fed}}^{t-1}), \nabla_w \tilde{F}_c(w_c^{t,i}) \rangle \right], \end{aligned} \quad (18)$$

where the i -th local model update of each client c can be expressed as the average gradient of the stored data samples:

$$\begin{aligned} (18) &= \sum_{c \in C_t} \left\| \sum_{i=0}^{m-1} \eta \sum_{x,y \in B_c} \frac{1}{|B_c|} \nabla_w l(w_c^{t,i}, x, y) \right\|^2 - \\ &\quad \frac{\eta}{|B_c|} \sum_{i=0}^{m-1} \sum_{(x,y) \in B_c} \langle \nabla_w F(w_{\text{fed}}^{t-1}), \nabla_w l(w_c^{t,i}, x, y) \rangle \\ &\leq \sum_{c \in C_t} \left[\sum_{i=0}^{m-1} \sum_{x,y \in B_c} \left(\alpha_c \|\nabla_w l(w_c^{t,i}, x, y)\|^2 \right. \right. \\ &\quad \left. \left. - \beta_c \langle \nabla_w F(w_{\text{fed}}^{t-1}), \nabla_w l(w_c^{t,i}, x, y) \rangle \right) \right]. \end{aligned}$$

where $\alpha_c = \frac{L\zeta_c^t}{2} \cdot \left(\frac{\eta}{|B_c|}\right)^2$ and $\beta_c = \zeta_c^t \cdot \left(\frac{\eta}{|B_c|}\right)$. \square

B. Proof of Theorem 2

Theorem 2 (Model Weight Divergence) With Assumption on Lipschitz property, for an arbitrary participating client set C_t , we have the following inequality for the weight divergence between the models trained through FL and CL after the t -th training round.

$$\begin{aligned} \|w_{\text{fed}}^t - w_{\text{cen}}^{mt}\|_2 &\leq (1 + \eta L)^m \|w_{\text{fed}}^{t-1} - w_{\text{cen}}^{m(t-1)}\|_2 \\ &\quad + \sum_{c \in C_t} \zeta_c^t \left[\eta \sum_{i=0}^{m-1} (1 + \eta L)^{m-1-i} G_c(w_c^{t,i}) \right], \end{aligned} \quad (19)$$

where $G_c(w) = \|\nabla_w \tilde{F}_c(w) - \nabla_w F(w)\|_2$.

Proof. According to the aggregation formula of Fed-Avg and the model update formula (1), we have:

$$\begin{aligned} \|w_{\text{fed}}^t - w_{\text{cen}}^{mt}\| &= \left\| \sum_{c \in C_t} \zeta_c^t w_c^{t,m} - w_{\text{cen}}^{mt} \right\| \\ &= \left\| \sum_{c \in C_t} \zeta_c^t \left[w_c^{t,m-1} - \eta \nabla_w \tilde{F}_c(w_c^{t,m-1}) \right] \right. \\ &\quad \left. - w_{\text{cen}}^{mt-1} + \eta \nabla_w F(w_{\text{cen}}^{mt-1}) \right\| \\ &\leq \left\| \sum_{c \in C_t} \zeta_c^t w_c^{t,m-1} - w_{\text{cen}}^{mt-1} \right\| \\ &\quad + \eta \left\| \sum_{c \in C_t} \zeta_c^t \nabla_w \tilde{F}_c(w_c^{t,m-1}) - \nabla_w F(w_{\text{cen}}^{mt-1}) \right\|. \end{aligned} \quad (20)$$

Then, we leverage the triangle inequality to upper bound (20):

$$\begin{aligned} (20) &\leq \sum_{c \in C_t} \zeta_c^t \|w_c^{t,m-1} - w_{\text{cen}}^{mt-1}\| + \eta \left\| \sum_{c \in C_t} \zeta_c^t \left[\nabla_w \tilde{F}_c(w_c^{t,m-1}) \right. \right. \\ &\quad \left. \left. - \nabla_w F(w_c^{t,m-1}) + \nabla_w F(w_c^{t,m-1}) - \nabla_w F(w_{\text{cen}}^{mt-1}) \right] \right\| \\ &= \sum_{c \in C_t} \zeta_c^t \|w_c^{t,m-1} - w_{\text{cen}}^{mt-1}\| \\ &\quad + \eta \left\| \sum_{c \in C_t} \zeta_c^t \left[\nabla_w \tilde{F}_c(w_c^{t,m-1}) - \nabla_w F(w_c^{t,m-1}) \right] \right\| \\ &\quad + \eta \left\| \sum_{c \in C_t} \zeta_c^t \left[\nabla_w F(w_c^{t,m-1}) - \nabla_w F(w_{\text{cen}}^{mt-1}) \right] \right\|. \end{aligned} \quad (21)$$

According to the Lipschitz continuity of global model F in Assumption 1, we can further upper bound the Eq. (21):

$$\begin{aligned} (21) &\leq \sum_{c \in C_t} \zeta_c^t (1 + \eta L) \|w_c^{t,m-1} - w_{\text{cen}}^{mt-1}\| + \\ &\quad \eta \left\| \sum_{c \in C_t} \zeta_c^t \left[\nabla_w \tilde{F}_c(w_c^{t,m-1}) - \nabla_w F(w_c^{t,m-1}) \right] \right\| \\ &\leq \sum_{c \in C_t} \zeta_c^t \left[(1 + \eta L) \|w_c^{t,m-1} - w_{\text{cen}}^{mt-1}\| + \eta G_c(w_c^{t,m-1}) \right], \end{aligned} \quad (22)$$

where $G_c(w) = \|\nabla_w \tilde{F}_c(w) - \nabla_w F(w)\|$ intuitively represents the divergence between empirical losses over the local data distribution and global data distribution.

Next, we derive the upper bound of $\|w_c^{t,m-1} - w_{\text{cen}}^{mt-1}\|$ for each participating client $c \in C_t$:

$$\begin{aligned} &\|w_c^{t,m-1} - w_{\text{cen}}^{mt-1}\| \\ &= \|w_c^{t,m-2} - \eta \nabla_w \tilde{F}_c(w_c^{t,m-2}) - w_{\text{cen}}^{mt-2} + \eta \nabla_w F(w_{\text{cen}}^{mt-2})\| \\ &\leq \|w_c^{t,m-2} - w_{\text{cen}}^{mt-2}\| + \eta \|\nabla_w F(w_c^{t,m-2}) - \nabla_w F(w_{\text{cen}}^{mt-2})\| + \eta \|\nabla_w \tilde{F}_c(w_c^{t,m-2}) - \nabla_w F(w_c^{t,m-2})\|, \end{aligned} \quad (23)$$

which can be transformed using Assumption of Lipschitz:

$$\begin{aligned}
(23) &\leq (1 + \eta L) \|w_c^{t,m-2} - w_{\text{cen}}^{mt-2}\| + \eta G_c(w_c^{t,m-2}) \\
&\leq (1 + \eta L)^{m-1} \|w_c^{t,0} - w_{\text{cen}}^{(t-1)m}\| + \\
&\quad \eta \sum_{i=0}^{m-2} (1 + \eta L)^{m-2-i} G_c(w_c^{t,i})
\end{aligned} \tag{24}$$

Combining inequalities (22) and (24), we have:

$$\begin{aligned}
\|w_{\text{fed}}^t - w_{\text{cen}}^{mt}\| &\leq \sum_{c \in C_t} \zeta_c^t \left[(1 + \eta L)^m \|w_c^{t,0} - w_{\text{cen}}^{(t-1)m}\| \right. \\
&\quad \left. + \eta \sum_{i=0}^{m-1} (1 + \eta L)^{m-1-i} G_c(w_c^{t,i}) \right] \\
&= (1 + \eta L)^m \|w_{\text{fed}}^{t-1} - w_{\text{cen}}^{(t-1)m}\| \\
&\quad + \sum_{c \in C_t} \zeta_c^t \left[\eta \sum_{i=0}^{m-1} (1 + \eta L)^{m-1-i} G_c(w_c^{t,i}) \right].
\end{aligned}$$

Hence, we bound the divergence between models trained through FL and CL by the initial model difference and the additional divergence caused by m local model updates of heterogeneous participating clients. \square

C. Proof of Lemma 1

Lemma 1 (Gradient Divergence) For an arbitrary client $c \in C$, $G_c(w) = \|\nabla \tilde{F}_c(w) - \nabla F(w)\|_2$ is bounded by:

$$\begin{aligned}
G_c(w) &\leq \left[\underbrace{\|\nabla_w F(w)\|_2^2}_{\text{constant}} + \sum_{(x,y) \in B_c} \frac{1}{|B_c|} \left(\underbrace{\|\nabla_w l(w, x, y)\|_2^2}_{\text{term 1}} - \right. \right. \\
&\quad \left. \left. 2 \underbrace{\langle \nabla_w l(w, x, y), \nabla_w F(w) \rangle}_{\text{term 2}} \right) \right]^{1/2},
\end{aligned} \tag{25}$$

where $\delta =$ is a constant term for all data samples.

Proof. The main idea of the proof is to express the local gradient $\nabla_w \tilde{F}_c(w)$ as the average gradient of the locally stored data samples in B_c :

$$\begin{aligned}
G_c(w) &= \|\nabla_w \tilde{F}_c(w) - \nabla_w F(w)\| \\
&= \left[\|\nabla_w F(w)\|^2 + \|\nabla_w \tilde{F}_c(w)\|^2 - 2\langle \nabla_w F(w), \nabla_w \tilde{F}_c(w) \rangle \right]^{1/2} \\
&= \left[\|\nabla_w F(w)\|^2 + \left\| \sum_{(x,y) \in B_c} \frac{\nabla_w l(w, x, y)}{|B_c|} \right\|^2 \right. \\
&\quad \left. - 2\langle \nabla_w F(w), \sum_{(x,y) \in B_c} \frac{\nabla_w l(w, x, y)}{|B_c|} \rangle \right]^{1/2} \\
&\leq \left[\|\nabla_w F(w)\|^2 + \sum_{(x,y) \in B_c} \frac{1}{|B_c|} \left(\|\nabla_w l(w, x, y)\|^2 \right. \right. \\
&\quad \left. \left. - 2\langle \nabla_w F(w), \nabla_w l(w, x, y) \rangle \right) \right]^{1/2},
\end{aligned}$$

where $\|\nabla_w F(w)\|^2$ is a constant term for all clients and local data samples. \square