

Supplementary Material for the Paper: “Wireless Personalized Federated Fine-Tuning of Large Language Models via Low-Rank Adaptation”

Haofeng Sun, Hui Tian, Jingheng Zheng, and Wanli Ni

APPENDIX A PROOF OF THEOREM 1

With Assumptions 1, 2, 3, we investigate the convergence behavior of global loss function (2) composed by the desired local loss function $\nabla F_k(\mathbf{W}_{t,n})$ in the following theorem.

Theorem 2: Suppose that Assumption 1 holds. The convergence behavior of the global loss function with respect to multiple tasks in AirPFed-LoRA framework between two consecutive fine-tuning rounds can be expressed by

$$\begin{aligned} & \sum_{n \in \mathcal{N}} \frac{D_n}{D} (F(\{\mathbf{W}_{t+1,n}\}) - F(\{\mathbf{W}_{t,n}\})) \\ & \leq \underbrace{\text{Vec}(\mathbf{W}_{t+1}^s - \mathbf{W}_t^s)^T \text{Vec}(\nabla \mathbf{W}_t^s F(\{\mathbf{W}_{t,n}\}))}_{E_{t,1}} \\ & + \underbrace{\sum_{n \in \mathcal{N}} \frac{D_n}{D} \text{Vec}(\mathbf{W}_{t+1,n}^p - \mathbf{W}_{t,n}^p)^T \text{Vec}(\nabla \mathbf{W}_{t,n}^p F_n^E(\mathbf{W}_{t,n}))}_{E_{t,2}} \\ & + \underbrace{\frac{\beta}{2} \|\mathbf{W}_{t+1}^s - \mathbf{W}_t^s\|_F^2}_{E_{t,3}} + \underbrace{\frac{\beta}{2} \sum_{n \in \mathcal{N}} \frac{D_n}{D} \|\mathbf{W}_{t+1,n}^p - \mathbf{W}_{t,n}^p\|_F^2}_{E_{t,4}}, \end{aligned} \quad (38)$$

where $\mathbf{W}_t^s = [\mathbf{B}_t^s; \mathbf{A}_t^s]$ denotes the s -adapter matrix, $\mathbf{W}_{t,n}^p = [\mathbf{B}_{t,n}^p; \mathbf{A}_{t,n}^p]$ denotes the p -adapter matrix of task n , and $F_n^E(\mathbf{W}_{t,n}) = \sum_{k \in \mathcal{K}_n} \frac{D_k}{D_n} F_k(\mathbf{W}_{t,n})$ denotes the desired edge loss function of task n .

Theorem 2 demonstrates that the convergence behavior of the global loss function is affected by four terms, i.e., $E_{t,1}$, $E_{t,2}$, $E_{t,3}$ and $E_{t,4}$. With Assumptions 2 and 3, we further bound the expectations of $E_{t,1} + E_{t,2}$ and $E_{t,3} + E_{t,4}$ in Lemma 1 and 2, respectively.

Lemma 1: The upper bound of $\mathbb{E}[E_{t,1} + E_{t,2}]$ is given by

$$\begin{aligned} \mathbb{E}[E_{t,1} + E_{t,2}] & \leq -\frac{\eta}{2} \left\| \sum_{n \in \mathcal{N}} \sum_{k \in \mathcal{K}_n} \frac{D_k}{D} \nabla F_k(\mathbf{W}_{t,n}) \right\|_F^2 \\ & - \frac{\eta}{2} \|\nabla \mathbf{W}_t^s F(\{\mathbf{W}_{t,n}\})\|_F^2 - \frac{\eta}{2} \sum_{n \in \mathcal{N}} \frac{D_n}{D} \|\nabla \mathbf{W}_{t,n}^p F_n^E(\mathbf{W}_{t,n})\|_F^2. \end{aligned} \quad (39)$$

H. Sun, H. Tian, and J. Zheng are with the State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications, Beijing 100876, China (e-mail: sunhaofeng@bupt.edu.cn; tianhui@bupt.edu.cn; zhengjh@bupt.edu.cn).

W. Ni is with the Department of Electronic Engineering, Tsinghua University, Beijing 100084, China (e-mail: niwanli@tsinghua.edu.cn).

Lemma 2: Suppose that Assumptions 2 and 3 hold. The expectation of $E_{t,3} + E_{t,4}$ can be bounded by

$$\begin{aligned} \mathbb{E}[E_{t,3} + E_{t,4}] & \leq \lambda \eta^2 (C + 1) \left\| \sum_{n \in \mathcal{N}} \sum_{k \in \mathcal{K}_n} \frac{D_k}{D} \nabla F_k(\mathbf{W}_{t,n}) \right\|_F^2 \\ & + \eta^2 (C + 1) \sigma_{\text{dif}}^2 + \eta^2 \sigma_{\text{SGD}}^2 + \eta^2 \left\| \hat{\mathbf{G}}_t^{\mathbf{W}_t^s} - \mathbf{G}_t^{\mathbf{W}_t^s} \right\|_F^2 \\ & + \eta^2 \sum_{n \in \mathcal{N}} \frac{D_n}{D} \left\| \hat{\mathbf{G}}_{t,n}^{\mathbf{W}_{t,n}^p} - \mathbf{G}_{t,n}^{\mathbf{W}_{t,n}^p} \right\|_F^2 \end{aligned} \quad (40)$$

By leveraging Theorem 2, Lemma 1 and 2, we derive an optimality gap of the AirPFed-LoRA framework and establish an upper bound for the variations of s - and p -adapters' matrices in Theorem 3.

Theorem 3: Building upon Theorem 2, Lemma 1 and 2, while setting $\eta \leq \frac{1}{\lambda(C+1)\beta}$, the optimality gap of the AirPFed-LoRA framework between two consecutive fine-tuning rounds can be bounded by

$$\begin{aligned} \sum_{n \in \mathcal{N}} \frac{D_n}{D} (F(\{\mathbf{W}_{t+1,n}\}) - F(\{\mathbf{W}_{t,n}\})) & \leq -\frac{\eta}{2} \Lambda_t \\ & + \frac{\beta \eta^2}{2} \sigma_{\text{SGD}}^2 + \frac{\beta \eta^2}{2} (C + 1) \sigma_{\text{dif}}^2 + \frac{\beta \eta^2 (L + Q) R}{2} \sigma_t^2 \Pi_t \end{aligned} \quad (41)$$

where $\Pi_t = \frac{R^s}{R} \text{MSE}_t^s + \frac{R^p}{R} \sum_{n \in \mathcal{N}} \frac{D_n}{D} \text{MSE}_{t,n}^p$ denotes the combined MSE of the edge and global aggregated gradient signals in AirComp, $R = R^s + R^p$ denotes the total rank of s - and p -adapters, Λ_t is the F -norm of gradient matrices which captures the variation of s - and p -adapters, given by

$$\Lambda_t = \|\nabla \mathbf{W}_t^s F(\{\mathbf{W}_{t,n}\})\|_F^2 + \sum_{n \in \mathcal{N}} \frac{D_n}{D} \|\nabla \mathbf{W}_{t,n}^p F_n^E(\mathbf{W}_{t,n})\|_F^2. \quad (42)$$

Suppose that Assumptions 1, 2, 3 hold while set the learning rate as $\eta = \frac{1}{\lambda(C+1)\beta}$. Then, the average gradients variation Λ_t can be bounded by

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \Lambda_t & \leq \frac{2 \sum_{n \in \mathcal{N}} \frac{D_n}{D} (F(\{\mathbf{W}_{1,n}\}) - F(\{\mathbf{W}_n^*\}))}{\eta T} \\ & + \frac{\sigma_{\text{SGD}}^2}{\lambda(C+1)} + \frac{\sigma_{\text{dif}}^2}{\lambda} + \frac{(L+Q)R}{\lambda(C+1)T} \sum_{t=1}^T \sigma_t^2 \Pi_t = \Psi_T, \end{aligned} \quad (43)$$

where $\{\mathbf{W}_n^*\}$ denotes the optimal global s -adapter and edge p -adapters matrices.