

APPENDIX A
PROOF OF THEOREM 1

We bring in the following assumptions from [1], [2] for analytical tractability.

Assumption 1. The loss function is L -smooth as $\|\nabla F(\mathbf{w}_1) - \nabla F(\mathbf{w}_2)\| \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|$ for arbitrary given \mathbf{w}_1 and \mathbf{w}_2 .

Assumption 2. The expected squared norm of stochastic gradients for each vehicle k is upper-bounded by

$$\mathbb{E}\|\nabla f(\mathbf{w}_k^{r,m})\|^2 \leq G^2, \forall k, \forall r, \forall m.$$

Assumption 3. The variance of mini-batch gradients is upper-bounded by

$$\|g(\mathbf{w}) - \nabla f(\mathbf{w})\|^2 \leq \sigma^2.$$

Assumption 4. The divergence between local and global loss functions is bounded by

$$\frac{1}{K} \sum_{k=1}^K \|\nabla f(\mathbf{w}) - \nabla F(\mathbf{w})\|^2 \leq \epsilon_g^2, \forall \mathbf{w}.$$

Based on the L -smoothness assumption, we have the following equation:

$$F(\mathbf{w}^{r+1}) - F(\mathbf{w}^r) \leq \langle \nabla F(\mathbf{w}^r), \mathbf{w}^{r+1} - \mathbf{w}^r \rangle + L/2 \|\mathbf{w}^{r+1} - \mathbf{w}^r\|^2. \quad (\text{A.1})$$

We give the expression of the local update of each model $\mathbf{w}_k^{m,r}$ with each modality m for vehicle k as

$$\mathbf{w}_{k,e+1}^{m,r} = \mathbf{w}_{k,e}^{m,r} - \eta \mathbf{g}_{k,e}^{m,r}, \quad (\text{A.2})$$

here, η is the learning rate, and $\mathbf{g}_{k,e}^{m,r}$ is the gradient descent for modality m in round r epoch $e \in \{0, \dots, E-1\}$. We then substitute the local update with the gradient and model updates of all modalities from all vehicles and take the expectation on both sides as

$$\begin{aligned} \mathbb{E}[F(\mathbf{w}^{r+1}) - F(\mathbf{w}^r)] &\leq \mathbb{E}\langle \mathbf{w}^{r+1} - \mathbf{w}^r, \nabla F(\mathbf{w}^r) \rangle + \frac{L}{2} \mathbb{E}\|\mathbf{w}^{r+1} - \mathbf{w}^r\|^2 \\ &\leq \sum_{m=1}^M \mathbb{E}\langle \nabla F(\mathbf{w}^{m,r}), -\frac{1}{K} \sum_{k=1}^K \sum_{e=0}^{E-1} \eta \mathbf{g}_{k,e}^{m,r} \rangle + \sum_{m=1}^M \frac{L}{2} \mathbb{E}\left\| \frac{1}{K} \sum_{k=1}^K \sum_{e=0}^{E-1} \eta \mathbf{g}_{k,e}^{m,r} \right\|^2 \end{aligned} \quad (\text{A.3})$$

For the left side of the sum, we obtain the following expressions according to Assumptions 2 and 3:

$$\begin{aligned} \sum_{m=1}^M \mathbb{E}\langle \nabla F(\mathbf{w}^{m,r}), -\frac{1}{K} \sum_{k=1}^K \sum_{e=0}^{E-1} \eta \mathbf{g}_{k,e}^{m,r} \rangle &\leq -\eta \sum_{m=1}^M \frac{1}{K} \sum_{k=1}^K \sum_{e=0}^{E-1} \mathbb{E}\langle \nabla F(\mathbf{w}^{m,r}), \nabla f_k(\mathbf{w}_{k,e}^{m,r}) \rangle \\ &\leq -\eta \sum_{m=1}^M \frac{1}{K} \sum_{k=1}^K \sum_{e=0}^{E-1} \left(\mathbb{E}\|\nabla F(\mathbf{w}^{m,r})\|^2 + \mathbb{E}\|\nabla f_k(\mathbf{w}_{k,e}^{m,r})\|^2 \right. \\ &\quad \left. - \mathbb{E}\|\nabla F(\mathbf{w}^{m,r}) - \nabla F(\mathbf{w}_{k,e}^{m,r}) + \nabla F(\mathbf{w}_{k,e}^{m,r}) - \nabla f_k(\mathbf{w}_{k,e}^{m,r})\|^2 \right) \\ &\leq -\eta \sum_{m=1}^M \frac{1}{K} \sum_{k=1}^K \left(\mathbb{E}\|\nabla F(\mathbf{w}^{m,r})\|^2 E + \sum_{e=0}^{E-1} \mathbb{E}\|\nabla f_k(\mathbf{w}_{k,e}^{m,r})\|^2 \right. \\ &\quad \left. - 2 \sum_{e=0}^{E-1} \mathbb{E}[\|\nabla F(\mathbf{w}^{m,r}) - \nabla F(\mathbf{w}_{k,e}^{m,r})\|^2 + \|\nabla F(\mathbf{w}_{k,e}^{m,r}) - \nabla f_k(\mathbf{w}_{k,e}^{m,r})\|^2] \right) \\ &\leq -\eta \sum_{m=1}^M \frac{1}{K} \sum_{k=1}^K \left(\mathbb{E}\|\nabla F(\mathbf{w}^{m,r})\|^2 E - 2 \sum_{e=0}^{E-1} L^2 \|\mathbf{w}^{m,r} - \mathbf{w}_{k,e}^{m,r}\|^2 \right) + 2\eta EM\epsilon_g^2. \end{aligned} \quad (\text{A.4})$$

Considering the norm $\|\mathbf{w}^{m,r} - \mathbf{w}_{k,e}^{m,r}\|^2$, we get the following expressions based on the Lemma 3 in [2]:

$$\mathbb{E}\left[\sum_{k=1}^K \|\mathbf{w}^{m,r} - \mathbf{w}_{k,e}^{m,r}\|^2\right] \leq 4\eta^2(E-1)^2 G^2. \quad (\text{A.5})$$

For the right side of the sum, we obtain

$$\frac{L}{2} \mathbb{E} \left\| \frac{1}{K} \sum_{k=1}^K \sum_{e=0}^{E-1} \eta \mathbf{g}_{k,e}^{m,r} \right\|^2 \leq \frac{L\eta^2 E^2}{2} \delta^2. \quad (\text{A.6})$$

Combining the above results, we obtain the following expression:

$$\mathbb{E}[F(\mathbf{w}^{r+1}) - F(\mathbf{w}^r)] \leq -\eta \frac{1}{K} \sum_{k=1}^K \sum_{m=1}^M \left(\mathbb{E} \|\nabla F(\mathbf{w}^{m,r})\|^2 E - 8E(E-1)^2 \eta^2 G^2 \right) + 2EM\epsilon_g^2 \eta + \frac{L\eta^2 E^2 \delta^2 M}{2} \quad (\text{A.7})$$

Then we rearrange the above expression and add all the terms from $r = \{0, 1, \dots, R\}$ to get the below expression

$$\frac{1}{R} \sum_{r=0}^{R-1} E [\|\nabla F(\mathbf{w}^{m,r})\|^2] \leq \frac{2\chi^m (F(\mathbf{w}^0) - F(\mathbf{w}^R))}{\eta ER} + 2\epsilon_g^2 + 8(E-1)^2 \eta^2 G^2 + \frac{L\eta E \delta^2}{2}, \quad (\text{A.8})$$

where χ^m is the contribution ratio of modality m to the training optimization during model training. Hence, Theorem 1 is proved.

APPENDIX B PROOF OF LEMMA 1

According to (A.7), the loss divergence of the global model in each round is

$$\mathbb{E}[F(\mathbf{w}^{r+1}) - F(\mathbf{w}^r)] \leq -\eta \left(\mathbb{E} \|\nabla F(\mathbf{w}^r)\|^2 E - 8E(E-1)^2 \eta^2 G^2 \right) + 2EM\epsilon_g^2 \eta + \frac{L\eta^2 E^2 \delta^2 M}{2}. \quad (\text{A.9})$$

We add the norm F^* to both sides and rearrange the terms to get the following expressions:

$$\mathbb{E}[F(\mathbf{w}^{r+1}) - F^*] \leq \mathbb{E}[F(\mathbf{w}^r) - F^*] - \eta \left(\mathbb{E} \|\nabla F(\mathbf{w}^r)\|^2 E - 8E(E-1)^2 \eta^2 G^2 \right) + 2EM\epsilon_g^2 \eta + \frac{L\eta^2 E^2 \delta^2 M}{2}. \quad (\text{A.10})$$

According to Lemma 8 in [3], the upper bound is derived as

$$\mathcal{O} \left(\frac{dr_0}{R} + \frac{c_1^{\frac{1}{2}} r_0^{\frac{1}{2}}}{R^{\frac{1}{2}}} + \frac{c_2^{\frac{1}{3}} r_0^{\frac{2}{3}}}{R^{\frac{2}{3}}} \right). \quad (\text{A.11})$$

We set $c_1 = LE^2 \delta^2 M$, $c_2 = E(E-1)^2 G^2$, $d = 1$, $r_0 = F(\mathbf{w}^r) - F^*$ to derive the following result:

$$\min_{0 \leq r \leq R-1} \mathbb{E} \|\nabla F(\mathbf{w}^r)\|^2 \leq \mathcal{O} \left(\frac{\Delta}{RE} + \frac{E\delta\sqrt{LM}\Delta^{\frac{1}{2}}}{R^{\frac{1}{2}}} + \frac{\Delta^{\frac{2}{3}}(E(E-1)^2 G^2)^{\frac{1}{3}}}{R^{\frac{2}{3}}} \right) \quad (\text{A.12})$$

where we set $F(\mathbf{w}^r) - F^* = \Delta$ and $\min_{0 \leq r \leq R-1} \mathbb{E} \|\nabla F(\mathbf{w}^r)\|^2 \leq \frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E} \|\nabla F(\mathbf{w}^r)\|^2$ is utilized. Hence, Lemma 1 is proved.

REFERENCES

- [1] C. Feng, H. H. Yang, D. Hu, Z. Zhao, T. Q. S. Quek, and G. Min, "Mobility-aware cluster federated learning in hierarchical wireless networks," *IEEE Trans. Wireless Commun.*, vol. 21, no. 10, pp. 8441–8458, Oct. 2022.
- [2] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of FedAvg on non-IID data," 2020. [Online]. Available: <https://arxiv.org/abs/1907.02189>
- [3] Y. Li and X. Lyu, "Convergence analysis of sequential federated learning on heterogeneous data," in *Proceedings of the 37th International Conference on Neural Information Processing Systems*, ser. NIPS '23, 2023.