# Layerwise and Dimensionwise Locally Adaptive Optimization Method (Supplementary Material)

**Plan of the supplementary material:** The supplementary material of this paper is composed of two main parts. Section A contains detailed proofs of our results and Section B where additional runs are provided. In particular, Theorem 1 is proved in subsection A.2.

## A Theoretical Analysis

We first recall in Table 2 some important notations that will be used in our following analysis.

| | | |
|---|---|---|
| $R, T$ | := | Number of communications rounds and local iterations (resp.) |
| $n, D, i$ | := | Total number of clients, portion sampled uniformly and client index |
| $\mathsf{h}, \ell$ | := | Total number of layers in the DNN and its index |
| $\phi(\cdot)$ | := | Scaling factor in FED-LAMBupdate |
| $\overline{\theta}$ | := | Global model (after periodic averaging) |
| $p_{r,i}^t$ | := | ratio computed at round $r$, local iteration $t$ and for device $i$. $p_{r,i}^{\ell,t}$ denotes its component at layer $\ell$ |

Table 2: Summary of notations used in the paper.

We now provide the proofs for the theoretical results of the main paper, including the intermediary Lemmas and the main convergence result, Theorem 1.

### A.1 Intermediary Lemmas

**Lemma.** *Consider $\{\overline{\theta_r}\}_{r>0}$, the sequence of parameters obtained running Algorithm 1. Then for $i \in [\![n]\!]$:*

$$\|\overline{\theta_r} - \theta_{r,i}\|^2 \leq \alpha^2 M^2 \phi_M^2 \frac{(1-\beta_2)p}{v_0} \ ,$$

*where $\phi_M$ is defined in H4 and $p$ is the total number of dimensions $p = \sum_{\ell=1}^{\mathsf{h}} p_\ell$.*

*Proof.* Assuming the simplest case when $T = 1$, i.e. one local iteration, then by construction of Algorithm 1, we have for all $\ell \in [\![\mathsf{h}]\!]$, $i \in [\![n]\!]$ and $r > 0$:

$$\theta_{r,i}^\ell = \overline{\theta_r}^\ell - \alpha \phi(\|\theta_{r,i}^{\ell,t-1}\|) p_{r,i}^j / \|p_{r,i}^\ell\| = \overline{\theta_r}^\ell - \alpha \phi(\|\theta_{r,i}^{\ell,t-1}\|) \frac{m_{r,i}^t}{\sqrt{v_r^t}} \frac{1}{\|p_{r,i}^\ell\|}$$

leading to

$$\|\overline{\theta_r} - \theta_{r,i}\|^2 = \sum_{\ell=1}^{\mathsf{h}} \left\langle \overline{\theta_r}^\ell - \theta_{r,i}^\ell \,|\, \overline{\theta_r}^\ell - \theta_{r,i}^\ell \right\rangle$$

$$\leq \alpha^2 M^2 \phi_M^2 \frac{(1-\beta_2)p}{v_0} \ ,$$

which concludes the proof. $\qquad\square$

**Lemma.** *Consider $\{\overline{\theta_r}\}_{r>0}$, the sequence of parameters obtained running Algorithm 1. Then for $r > 0$:*

$$\left\| \frac{\overline{\nabla f(\theta_r)}}{\sqrt{v_r^t}} \right\|^2 \geq \frac{1}{2} \left\| \frac{\nabla f(\overline{\theta_r})}{\sqrt{v_r^t}} \right\|^2 - \overline{L} \alpha^2 M^2 \phi_M^2 \frac{(1-\beta_2)p}{v_0}$$

*where $M$ is defined in H2, $p$ is the total number of dimensions $p = \sum_{\ell=1}^{\mathsf{h}} p_\ell$ and $\phi_M$ is defined in H4.*

*Proof.* Consider the following sequence:

$$\left\|\frac{\overline{\nabla}f(\theta_r)}{\sqrt{v_r^t}}\right\|^2 \geq \frac{1}{2}\left\|\frac{\nabla f(\overline{\theta_r})}{\sqrt{v_r^t}}\right\|^2 - \left\|\frac{\overline{\nabla}f(\theta_r) - \nabla f(\overline{\theta_r})}{\sqrt{v_r^t}}\right\|^2 ,$$

where the inequality is due to the Cauchy-Schwartz inequality.

Under the smoothness assumption H1 and using Lemma 1, we have

$$\left\|\frac{\overline{\nabla}f(\theta_r)}{\sqrt{v_r^t}}\right\|^2 \geq \frac{1}{2}\left\|\frac{\nabla f(\overline{\theta_r})}{\sqrt{v_r^t}}\right\| - \left\|\frac{\overline{\nabla}f(\theta_r) - \nabla f(\overline{\theta_r})}{\sqrt{v_r^t}}\right\|^2$$

$$\geq \frac{1}{2}\left\|\frac{\nabla f(\overline{\theta_r})}{\sqrt{v_r^t}}\right\|^2 - \overline{L}\alpha^2 M^2 \phi_M^2 \frac{(1-\beta_2)p}{v_0} ,$$

which concludes the proof. □

## A.2 Proof of Theorem 1

We now develop a proof for the two intermediary lemmas, Lemma 1 and Lemma 2, in the case when each local model is obtained after more than one local update. Then the two quantities, either the gap between the periodically averaged parameter and each local update, *i.e.*, $\|\overline{\theta_r} - \theta_{r,i}\|^2$, and the ratio of the average gradient, more particularly its relation to the gradient of the average global model (*i.e.*, $\left\|\frac{\overline{\nabla}f(\theta_r)}{\sqrt{v_r^t}}\right\|$ and $\left\|\frac{\nabla f(\overline{\theta_r})}{\sqrt{v_r^t}}\right\|$), are impacted.

**Theorem.** *Assume **H1-H4**. Consider $\{\overline{\theta_r}\}_{r>0}$, the sequence of parameters obtained running Algorithm 1 with a decreasing learning rate $\alpha$. Let the number of local epochs be $T \geq 1$ and $\lambda = 0$. Then, at iteration $\tau$, we have:*

$$\frac{1}{\tau}\sum_{t=1}^{\tau}\mathbb{E}\left[\left\|\frac{\nabla f(\overline{\theta}_t)}{\hat{v}_t^{1/4}}\right\|^2\right] \leq \sqrt{\frac{M^2 p}{n}}\frac{\mathbb{E}[f(\overline{\theta}_1)] - \min_{\theta\in\Theta}f(\theta)}{\mathsf{h}\alpha_r\tau} + \frac{\phi_M\sigma^2}{\tau n}\sqrt{\frac{1-\beta_2}{M^2 p}}$$

$$+ 4\alpha\left[\frac{\alpha^2 L_\ell}{\sqrt{v_0}}M^2(T-1)^2\phi_M^2(1-\beta_2)p + \frac{M^2}{\sqrt{v_0}} + \phi_M^2\sqrt{M^2 + p\sigma^2} + \phi_M\frac{\mathsf{h}\sigma^2}{\sqrt{n}}\right] + cst.$$

*Proof.* Using H1, we have:

$$f(\overline{\vartheta}_{r+1}) \leq f(\overline{\vartheta}_r) + \left\langle \nabla f(\overline{\vartheta}_r) \,|\, \overline{\vartheta}_{r+1} - \overline{\vartheta}_r \right\rangle + \sum_{\ell=1}^{L}\frac{L_\ell}{2}\|\overline{\vartheta}_{r+1}^\ell - \overline{\vartheta}_r^\ell\|^2$$

$$\leq f(\overline{\vartheta}_r) + \sum_{\ell=1}^{\mathsf{h}}\sum_{j=1}^{p_\ell}\nabla_\ell f(\overline{\vartheta}_r)^j(\overline{\vartheta}_{r+1}^{\ell,j} - \overline{\vartheta}_r^{\ell,j}) + \sum_{\ell=1}^{L}\frac{L_\ell}{2}\|\overline{\vartheta}_{r+1}^\ell - \overline{\vartheta}_r^\ell\|^2 .$$

Taking expectations on both sides leads to:

$$-\mathbb{E}[\left\langle \nabla f(\overline{\vartheta}_r) \,|\, \overline{\vartheta}_{r+1} - \overline{\vartheta}_r \right\rangle] \leq \mathbb{E}[f(\overline{\vartheta}_r) - f(\overline{\vartheta}_{r+1})] + \sum_{\ell=1}^{L}\frac{L_\ell}{2}\mathbb{E}[\|\overline{\vartheta}_{r+1}^\ell - \overline{\vartheta}_r^\ell\|^2] . \tag{5}$$

Yet, we observe that, using the classical intermediate quantity, used for proving convergence results of adaptive optimization methods, see for instance [29], we have:

$$\overline{\vartheta}_r = \overline{\theta}_r + \frac{\beta_1}{1-\beta_1}(\overline{\theta}_r - \overline{\theta}_{r-1}) , \tag{6}$$

where $\bar{\theta}_r$ denotes the average of the local models at round $r$. Then for each layer $\ell$,

$$\bar{\vartheta}_{r+1}^\ell - \bar{\vartheta}_r^\ell = \frac{1}{1-\beta_1}(\bar{\theta}_{r+1}^\ell - \bar{\theta}_r^\ell) - \frac{\beta_1}{1-\beta_1}(\bar{\theta}_r^\ell - \bar{\theta}_{r-1}^\ell) \tag{7}$$

$$= \frac{\alpha_r}{1-\beta_1}\frac{1}{n}\sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\|p_{r,i}^\ell\|}p_{r,i}^\ell - \frac{\alpha_{r-1}}{1-\beta_1}\frac{1}{n}\sum_{i=1}^n \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\|p_{r-1,i}^\ell\|}p_{r-1,i}^\ell \tag{8}$$

$$= \frac{\alpha\beta_1}{1-\beta_1}\frac{1}{n}\sum_{i=1}^n \left( \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t}\|p_{r-1,i}^\ell\|} \right)m_{r-1}^t + \frac{\alpha}{n}\sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|}g_{r,i}^t \, , \tag{9}$$

where we have assumed a constant learning rate $\alpha$.

We note for all $\theta \in \Theta$, the majorant $G > 0$ such that $\phi(\|\theta\|) \leq G$. Then, following (5), we obtain:

$$-\mathbb{E}[\langle \nabla f(\bar{\vartheta}_r) \,|\, \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle] \leq \mathbb{E}[f(\bar{\vartheta}_r) - f(\bar{\vartheta}_{r+1})] + \sum_{\ell=1}^L \frac{L_\ell}{2}\mathbb{E}[\|\bar{\vartheta}_{r+1} - \bar{\vartheta}_r\|^2] \, . \tag{10}$$

Developing the LHS of (10) using (7) leads to

$$\langle \nabla f(\bar{\vartheta}_r) \,|\, \bar{\vartheta}_{r+1} - \bar{\vartheta}_r \rangle = \sum_{\ell=1}^h \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j(\bar{\vartheta}_{r+1}^{\ell,j} - \bar{\vartheta}_r^{\ell,j}) \tag{11}$$

$$= \frac{\alpha\beta_1}{1-\beta_1}\frac{1}{n}\sum_{\ell=1}^h \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j \left[ \sum_{i=1}^n \left( \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|} - \frac{\phi(\|\theta_{r-1,i}^\ell\|)}{\sqrt{v_{r-1}^t}\|p_{r-1,i}^\ell\|} \right)m_{r-1}^t \right] \tag{12}$$

$$\underbrace{-\frac{\alpha}{n}\sum_{\ell=1}^h \sum_{j=1}^{p_\ell} \nabla_\ell f(\bar{\vartheta}_r)^j \sum_{i=1}^n \frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_r^t}\|p_{r,i}^\ell\|}g_{r,i}^{t,l,j}}_{=A_1} \, . \tag{13}$$

We change all index $r$ to iteration $t$. Suppose $T$ is the number of local iterations. We can write (13) as

$$A_1 = -\alpha_t \langle \nabla f(\bar{\vartheta}_t), \frac{\bar{g}_t}{\sqrt{\hat{v}_t}} \rangle,$$

where $\bar{g}_t = \frac{1}{n}\sum_{i=1}^n \bar{g}_{t,i}$, with $\bar{g}_{t,i} = \left[ \frac{\phi(\|\theta_{t,i}^1\|)}{\|p_{t,i}^1\|}g_{t,i}^1, ..., \frac{\phi(\|\theta_{t,i}^L\|)}{\|p_{t,i}^L\|}g_{t,i}^L \right]$ representing the normalized gradient (concatenated by layers) of the $i$-th device. It holds that

$$\langle \nabla f(\bar{\vartheta}_t), \frac{\bar{g}_t}{\sqrt{\hat{v}_t}} \rangle = \frac{1}{2}\|\frac{\nabla f(\bar{\vartheta}_t)}{\hat{v}_t^{1/4}}\|^2 + \frac{1}{2}\|\frac{\bar{g}_t}{\hat{v}_t^{1/4}}\|^2 - \|\frac{\nabla f(\bar{\vartheta}_t) - \bar{g}_t}{\hat{v}_t^{1/4}}\|^2. \tag{14}$$

To bound the last term on the RHS, we have

$$\|\frac{\nabla f(\bar{\vartheta}_t) - \bar{g}_t}{\hat{v}_t^{1/4}}\|^2 = \|\frac{\frac{1}{n}\sum_{i=1}^n(\nabla f(\bar{\vartheta}_t) - \bar{g}_{t,i})}{\hat{v}_t^{1/4}}\|^2$$

$$\leq \frac{1}{n}\sum_{i=1}^n \|\frac{\nabla f(\bar{\vartheta}_t) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}}\|^2$$

$$\leq \frac{2}{n}\sum_{i=1}^n \left( \|\frac{\nabla f(\bar{\vartheta}_t) - \nabla f(\bar{\theta}_t)}{\hat{v}_t^{1/4}}\|^2 + \|\frac{\nabla f(\bar{\theta}_t) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}}\|^2 \right).$$

By Lipschitz smoothness of the loss function, the first term admits

$$
\frac{2}{n}\sum_{i=1}^{n}\|\frac{\nabla f_i(\bar{\vartheta}_t) - \nabla f_i(\bar{\theta}_t)}{\hat{v}_t^{1/4}}\|^2 \leq \frac{2}{n\sqrt{v_0}}\sum_{i=1}^{n}L_\ell\|\bar{\vartheta}_t - \bar{\theta}_t\|^2
$$

$$
= \frac{2L_\ell}{n\sqrt{v_0}}\frac{\beta_1^2}{(1-\beta_1)^2}\sum_{i=1}^{n}\|\bar{\theta}_t - \bar{\theta}_{t-1}\|^2
$$

$$
\leq \frac{2\alpha_r^2 L_\ell}{n\sqrt{v_0}}\frac{\beta_1^2}{(1-\beta_1)^2}\sum_{l=1}^{L}\sum_{i=1}^{n}\|\frac{\phi(\|\theta_{t,i}^l\|)}{\|p_{t,i}^l\|}p_{t,i}^l\|^2
$$

$$
\leq \frac{2\alpha_r^2 L_\ell p\phi_M^2}{\sqrt{v_0}}\frac{\beta_1^2}{(1-\beta_1)^2}.
$$

For the second term,

$$
\frac{2}{n}\sum_{i=1}^{n}\|\frac{\nabla f(\bar{\theta}_t) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}}\|^2 \leq \frac{4}{n}\Big(\underbrace{\sum_{i=1}^{n}\|\frac{\nabla f(\bar{\theta}_t) - \nabla f(\theta_{t,i})}{\hat{v}_t^{1/4}}\|^2}_{B_1} + \underbrace{\sum_{i=1}^{n}\|\frac{\nabla f(\theta_{t,i}) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}}\|^2}_{B_2}\Big). \tag{15}
$$

Using the smoothness of $f_i$ we can transform $B_1$ into consensus error by

$$
B_1 \leq \frac{L}{\sqrt{v_0}}\sum_{i=1}^{n}\|\bar{\theta}_t - \theta_{t,i}\|^2
$$

$$
= \frac{\alpha_r^2 L}{\sqrt{v_0}}\sum_{i=1}^{n}\sum_{l=1}^{L}\|\sum_{j=\lfloor t\rfloor_T+1}^{t}\Big(\frac{\phi(\|\theta_{j,i}^l\|)}{\|p_{j,i}^l\|}p_{j,i}^l - \frac{1}{n}\sum_{k=1}^{n}\frac{\phi(\|\theta_{j,k}^l\|)}{\|p_{j,k}^l\|}p_{j,k}^l\Big)\|^2 \tag{16}
$$

$$
\leq n\frac{\alpha_t^2 L}{\sqrt{v_0}}M^2(T-1)^2\phi_M^2(1-\beta_2)p
$$

where the last inequality stems from Lemma 1 in the particular case where $\theta_{t,i}$ are averaged every $ct+1$ local iterations for any integer $c$, since $(t-1) - (\lfloor t\rfloor_T + 1) + 1 \leq T - 1$.

We now develop the expectation of $B_2$ under the simplification that $\beta_1 = 0$:

$$
\mathbb{E}[B_2] = \mathbb{E}[\sum_{i=1}^{n}\|\frac{\nabla f(\theta_{t,i}) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}}\|^2]
$$

$$
\leq \frac{nM^2}{\sqrt{v_0}} + n\phi_M^2\sqrt{M^2 + p\sigma^2} - 2\sum_{i=1}^{n}\mathbb{E}[\langle\nabla f(\theta_{t,i}), \bar{g}_{t,i}\rangle/\sqrt{\hat{v}_t}]
$$

$$
= \frac{nM^2}{\sqrt{v_0}} + n\phi_M^2\sqrt{M^2 + p\sigma^2} - 2\sum_{i=1}^{n}\sum_{\ell=1}^{L}\mathbb{E}[\langle\nabla_\ell f(\theta_{t,i}), \frac{\phi(\|\theta_{t,i}^l\|)}{\|p_{t,i}^l\|}g_{t,i}^l\rangle/\sqrt{\hat{v}_t^l}]
$$

$$
= \frac{nM^2}{\sqrt{v_0}} + n\phi_M^2\sqrt{M^2 + p\sigma^2} - 2\sum_{i=1}^{n}\sum_{l=1}^{L}\sum_{i=1}^{p_l}\mathbb{E}[\nabla_l f(\theta_{t,i})^j\frac{\phi(\|\theta_{t,i}^{l,j}\|)}{\sqrt{\hat{v}_t^{l,j}}\|p_{t,i}^{l,j}\|}g_{t,i}^{l,j}]
$$

$$
\leq \frac{nM^2}{\sqrt{v_0}} + n\phi_M^2\sqrt{M^2 + p\sigma^2} - 2\sum_{i=1}^{n}\sum_{l=1}^{L}\sum_{i=1}^{p_l}\mathbb{E}\left[\sqrt{\frac{1-\beta_2}{M^2 p_\ell}}\phi(\|\theta_{r,i}^{l,j}\|)\nabla_l f(\theta_{t,i})^j g_{t,i}^{l,j}\right]
$$

$$
- 2\sum_{i=1}^{n}\sum_{l=1}^{L}\sum_{j=1}^{p_l}E\left[\Big(\phi(\|\theta_{r,i}^{l,j}\|)\nabla_l f(\theta_{t,i})^j\frac{g_{r,i}^{t,l,j}}{\|p_{r,i}^{l,j}\|}\Big)\mathbf{1}\Big(\text{sign}(\nabla_l f(\theta_{t,i})^j \neq \text{sign}(g_{r,i}^{t,l,j}))\Big)\right]
$$

where we use assumption H2, H3 and H4. Yet,

$$
-\mathbb{E}\left[\Big(\phi(\|\theta_{r,i}^{l,j}\|)\nabla_l f(\theta_{t,i})^j\frac{g_{r,i}^{t,l,j}}{\|p_{r,i}^{l,j}\|}\Big)\mathbf{1}\Big(\text{sign}(\nabla_l f(\theta_{t,i})^j \neq \text{sign}(g_{r,i}^{t,l,j}))\Big)\right]
$$

$$
\leq \phi_M\nabla_l f(\theta_{t,i})^j\mathbb{P}\left[\text{sign}(\nabla_l f(\theta_{t,i})^j \neq \text{sign}(g_{r,i}^{t,l,j}))\right]
$$

Then we have:

$$\mathbb{E}[B_2] \leq \frac{nM^2}{\sqrt{v_0}} + n\phi_M^2 \sqrt{M^2 + p\sigma^2} - 2\phi_m \sqrt{\frac{1-\beta_2}{M^2 p}} \sum_{i=1}^{n} \mathbb{E}[\|[\nabla f(\theta_{t,i})\|^2] + \phi_M \frac{h\sigma^2}{\sqrt{n}}$$

Thus, (15) becomes:

$$\frac{2}{n} \sum_{i=1}^{n} \|\frac{\nabla f_i(\bar{\theta}_t) - \bar{g}_{t,i}}{\hat{v}_t^{1/4}}\|^2 \leq 4 \left[ \frac{\alpha_t^2 Ll}{\sqrt{v_0}} \alpha_r^2 M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p + \frac{M^2}{\sqrt{v_0}} + \phi_M^2 \sqrt{M^2 + p\sigma^2} + \phi_M \frac{h\sigma^2}{\sqrt{n}} \right]$$

Substituting all ingredients into (14), we obtain

$$-\alpha_t \mathbb{E}[\langle \nabla f(\bar{\vartheta}_t), \frac{\bar{g}_t}{\sqrt{\hat{v}_t}} \rangle] \leq -\frac{\alpha_t}{2} \mathbb{E}\left[\|\frac{\nabla f(\bar{\vartheta}_t)}{\hat{v}_t^{1/4}}\|^2\right] - \frac{\alpha_t}{2} \mathbb{E}\left[\|\frac{\bar{g}_t}{\hat{v}_t^{1/4}}\|^2\right] + \frac{2\alpha_t^3 L_\ell p \phi_M^2}{\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2}$$
$$+ 4 \left[ \frac{\alpha_t^2 L}{\sqrt{v_0}} M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p + \frac{M^2}{\sqrt{v_0}} + \phi_M^2 \sqrt{M^2 + p\sigma^2} + \phi_M \frac{h\sigma^2}{\sqrt{n}} \right].$$

At the same time, we have

$$\mathbb{E}\left[\|\frac{\bar{g}_t}{\hat{v}_t^{1/4}}\|^2\right] = \frac{1}{n^2} \mathbb{E}\left[\|\frac{\sum_{i=1}^{n} \bar{g}_{t,i}}{\hat{v}_t^{1/4}}\|^2\right]$$
$$= \frac{1}{n^2} \mathbb{E}\left[ \sum_{l=1}^{L} \sum_{i=1}^{n} \|\frac{\phi(\|\theta_{t,i}^l\|)}{\hat{v}^{1/4}\|p_{t,i}^l\|} g_{t,i}^l\|^2 \right]$$
$$\geq \phi_m^2 (1-\beta_2) \mathbb{E}\left[\|\frac{1}{n} \sum_{i=1}^{n} \frac{\nabla f(\theta_{t,i})}{\hat{v}^{1/4}}\|^2\right]$$
$$= \phi_m^2 (1-\beta_2) \mathbb{E}\left[\|\frac{\overline{\nabla f}(\theta_t)}{\hat{v}^{1/4}}\|^2\right]$$

Regarding $\left\|\frac{\overline{\nabla f}(\theta_t)}{\hat{v}_t^{1/4}}\right\|^2$, we have

$$\left\|\frac{\overline{\nabla f}(\theta_t)}{\hat{v}_t^{1/4}}\right\|^2 \geq \frac{1}{2} \left\|\frac{\nabla f(\overline{\theta}_t)}{\hat{v}_t^{1/4}}\right\|^2 - \left\|\frac{\overline{\nabla f}(\theta_t) - \nabla f(\overline{\theta}_t)}{\hat{v}_t^{1/4}}\right\|^2$$
$$\geq \frac{1}{2} \left\|\frac{\nabla f(\overline{\theta}_t)}{\hat{v}_t^{1/4}}\right\|^2 - \left\|\frac{\frac{1}{n}\sum_{i=1}^{n}(\nabla f(\theta_{t,i}) - \nabla f(\bar{\theta}_i))}{\hat{v}_t^{1/4}}\right\|^2$$
$$\geq \frac{1}{2} \left\|\frac{\nabla f(\overline{\theta}_t)}{\hat{v}_t^{1/4}}\right\|^2 - \frac{\alpha_t^2 L_\ell}{\sqrt{v_0}} M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p,$$

where the last line is due to (16). Therefore, we have obtained

$$A_1 \leq -\frac{\phi_m^2 (1-\beta_2)}{2} \left\|\frac{\nabla f(\overline{\theta}_t)}{\hat{v}_t^{1/4}}\right\|^2 + \frac{\alpha_r^2 L_\ell}{\sqrt{v_0}} M^2 (T-1)^2 \phi_m^2 \phi_M^2 (1-\beta_2)^2 p + \frac{2\alpha^3 L_\ell p \phi_M^2}{\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2}$$
$$+ 4\alpha_t \left[ \frac{\alpha_t^2 L}{\sqrt{v_0}} M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p + \frac{M^2}{\sqrt{v_0}} + \phi_M^2 \sqrt{M^2 + p\sigma^2} + \phi_M \frac{h\sigma^2}{\sqrt{n}} \right].$$

Substitute back into (13), and leave other derivations unchanged. Assuming $M \leq 1$, we have the following:

$$\frac{1}{\tau} \sum_{t=1}^{\tau} \mathbb{E}\left[\left\|\frac{\nabla f(\bar{\theta}_t)}{\hat{v}_t^{1/4}}\right\|^2\right]$$

$$\lesssim \sqrt{\frac{M^2 p}{n}} \frac{f(\bar{\vartheta}_1) - \mathbb{E}[f(\bar{\vartheta}_{\tau+1})]}{\mathsf{h}\alpha_t \tau} + \frac{\alpha_t}{n^2} \sum_{r=1}^{\tau} \sum_{i=1}^{n} \sigma_i^2 \mathbb{E}\left[\left\|\frac{\phi(\|\theta_{r,i}^\ell\|)}{\sqrt{v_t}\|p_{r,i}^\ell\|}\right\|^2\right] + \frac{2\alpha^3 L_\ell p \phi_M^2}{\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2}$$

$$+ 4\alpha_t \left[\frac{\alpha_t^2 L_\ell}{\sqrt{v_0}} M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p + \frac{M^2}{\sqrt{v_0}} + \phi_M^2 \sqrt{M^2 + p\sigma^2} + \phi_M \frac{\mathsf{h}\sigma^2}{\sqrt{n}}\right] + \frac{\overline{L}\beta_1^2 \mathsf{h}(1-\beta_2) M^2 \phi_M^2 n}{2(1-\beta_1)^2 v_0}$$

$$+ \frac{\alpha_t \beta_1}{1-\beta_1} \sqrt{(1-\beta_2) p} \frac{\mathsf{h} M^2}{\sqrt{v_0}} + \overline{L}\alpha_t^2 M^2 \phi_M^2 \frac{(1-\beta_2)p}{Tv_0}$$

$$\leq \sqrt{\frac{M^2 p}{n}} \frac{\mathbb{E}[f(\bar{\theta}_1)] - \min_{\theta \in \Theta} f(\theta)}{\mathsf{h}\alpha_t \tau} + \frac{\phi_M \sigma^2}{\tau n} \sqrt{\frac{1-\beta_2}{M^2 p}}$$

$$+ 4\alpha_t \left[\frac{\alpha_t^2 L_\ell}{\sqrt{v_0}} M^2 (T-1)^2 \phi_M^2 (1-\beta_2) p + \frac{M^2}{\sqrt{v_0}} + \phi_M^2 \sqrt{M^2 + p\sigma^2} + \phi_M \frac{\mathsf{h}\sigma^2}{\sqrt{n}}\right]$$

$$+ \frac{\alpha_t \beta_1}{1-\beta_1} \sqrt{(1-\beta_2) p} \frac{\mathsf{h} M^2}{\sqrt{v_0}} + \overline{L}\alpha_t^2 M^2 \phi_M^2 \frac{(1-\beta_2)p}{Tv_0} + \frac{\overline{L}\beta_1^2 \mathsf{h}(1-\beta_2) M^2 \phi_M^2 n}{2(1-\beta_1)^2 v_0} + \frac{2\alpha^3 L_\ell p \phi_M^2}{\sqrt{v_0}} \frac{\beta_1^2}{(1-\beta_1)^2}.$$

This concludes the proof.

$\square$

## B  Additional Numerical Experiments

In below we provide two more results on CIFAR dataset trained on a Linux server with four Nvidia Tesla V100 cards, which is the hardware setting for all the experiments conducted in this paper. In Figure 5, we report the test accuracies of a ResNet-9 [10] trained on CIFAR-10 dataset, where the data is iid allocated among clients. We run 1 and 3 local epochs for 10 clients. From the figures, we observe similar advantage as the set of experiments presented in the main paper: faster convergence than local AMS and local SGD. In both cases, FED-LAMB also generalizes better than local SGD.
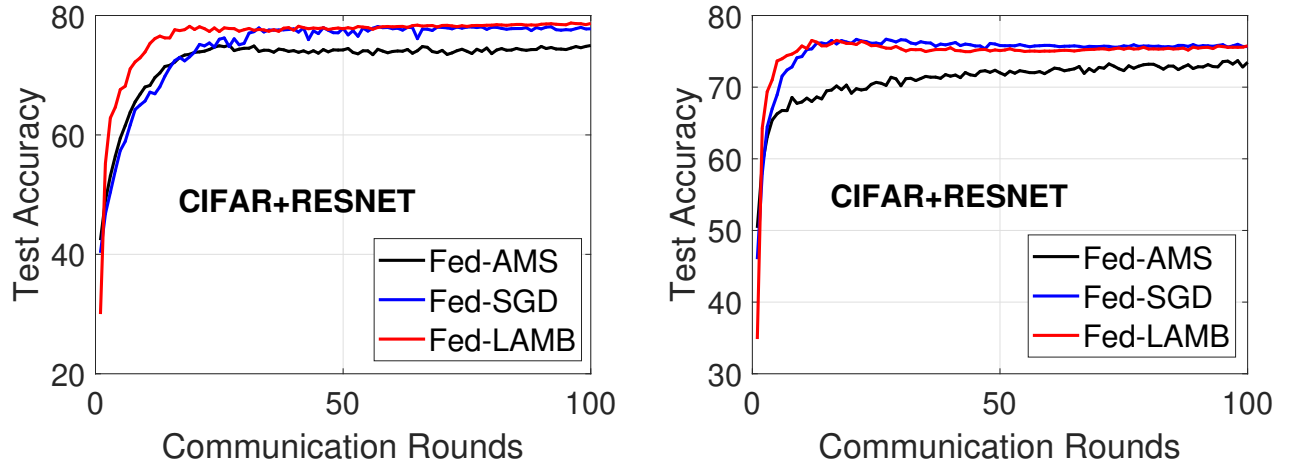


Figure 5: **From Left to Right**: Test accuracy on CIFAR+ResNet, with iid data distribution. 10 clients and (Left) 1 local epoch, (Right) 3 local epochs.