# Convergence Analysis Details of FedATMV

This document provides the detailed proof for the convergence theorem of the FedATMV presented in the main paper. We first restate the assumptions and the main theorem for clarity.

## 1 Assumptions

**A1.** (*L-smoothness*) The global loss function $F(w)$ and all local loss functions $f_i(w)$ for $i \in \{0, 1, \ldots, N\}$ are continuously differentiable and $L$-smooth. That is, for any $w, v \in \mathbb{R}^d$, there exists a constant $L > 0$ such that:

$$\|\nabla F(w) - \nabla F(v)\| \le L\|w - v\|. \tag{1}$$

**A2.** (*Unbiased and Bounded Variance Gradients*) The stochastic gradients computed on clients and the server are unbiased estimators of the true local gradients, and their variances are bounded. For any client $i \in \{1, \ldots, N\}$ and server ($i = 0$), there exists a constant $\sigma^2 \ge 0$ such that:

$$\mathbb{E}[g_i(w)] = \nabla f_i(w) \quad \text{and} \quad \mathbb{E}[\|g_i(w) - \nabla f_i(w)\|^2] \le \sigma^2. \tag{2}$$

**A3.** (*Bounded Gradient Divergence*) The dissimilarity between local client data distributions is bounded. We assume that the expected squared norm of the difference between local and global gradients is bounded by a constant $\zeta^2 \ge 0$:

$$\frac{1}{M} \sum_{i \in \mathcal{S}_t} \|\nabla f_i(w) - \nabla F(w)\|^2 \le \zeta^2. \tag{3}$$

**A4.** (*Bounded Model and Update Norms*) The adaptive parameters $\lambda_t$ and $\rho_t$ are bounded during training by constants $\lambda_{\max}$ and $\rho_{\max}$. Additionally, we assume the norm of the global model's gradient and the server's local gradient are bounded by a constant $G^2$: $\|\nabla F(w)\|^2 \le G^2$ and $\|\nabla f_0(w)\|^2 \le G^2$.

## 2 Main Theorem and Proof

**Theorem 1** (Convergence of FedATMV). *Let Assumptions A1-A4 hold. Consider the FedATMV algorithm with client learning rate $\eta$, server learning rate $\eta_0$, and total rounds $T$. If we set $\eta = \eta_0 = \mathcal{O}(1/\sqrt{T})$, then for a sufficiently large $T$, the convergence of FedATMV is bounded. The bound is influenced by our adaptive parameters $\lambda_{\max}$ and $\rho_{\max}$:*

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}[\|\nabla F(w_{t-1})\|^2] \le \frac{F(w_0) - F^*}{\eta KT} + 2L\eta K\zeta^2 + \frac{2L\eta K\sigma^2}{M}$$

$$+ \frac{1}{\eta K}\left(\lambda_{\max}\eta_0 EG^2 + 2L\rho_{\max}^2 G^2 + L\lambda_{\max}^2 \eta_0^2 E^2 G^2\right), \tag{4}$$

*where $F^*$ is the optimal value of $F(w)$. This implies that FedATMV achieves a convergence rate of $\mathcal{O}(1/\sqrt{T})$.*

*Proof.* The proof analyzes the expected one-round progress, which is the difference $\mathbb{E}[F(w_t)] - \mathbb{E}[F(w_{t-1})]$. From the $L$-smoothness of the global loss function $F(w)$ (Assumption A1), we have the descent lemma:

$$\mathbb{E}[F(w_t)] \leq \mathbb{E}[F(w_{t-1})] + \mathbb{E}[\langle \nabla F(w_{t-1}), w_t - w_{t-1}\rangle]$$
$$+ \frac{L}{2}\mathbb{E}[\|w_t - w_{t-1}\|^2]. \tag{5}$$

The proof proceeds by bounding the two expectation terms on the right-hand side: the inner product term and the squared norm term.

## 2.1 Part 1: Bounding the Inner Product Term $\mathbb{E}[\langle \nabla F(w_{t-1}), w_t - w_{t-1}\rangle]$

The one-round update $w_t - w_{t-1}$ can be decomposed into the client aggregation contribution and the server update contribution:

$$w_t - w_{t-1} = (\bar{w}_t - w_{t-1}) + (w_t - \bar{w}_t). \tag{6}$$

Thus, the inner product can be split:

$$\mathbb{E}[\langle \nabla F(w_{t-1}), w_t - w_{t-1}\rangle] = \mathbb{E}[\langle \nabla F(w_{t-1}), \bar{w}_t - w_{t-1}\rangle]$$
$$+ \mathbb{E}[\langle \nabla F(w_{t-1}), w_t - \bar{w}_t\rangle]. \tag{7}$$

Let's analyze each part. For the client aggregation part:

$$\mathbb{E}[\bar{w}_t - w_{t-1}] = \mathbb{E}\left[\frac{1}{M}\sum_{i\in\mathcal{S}_t}(w_t^{(i,K)} - w_{t-1})\right]$$
$$= \mathbb{E}\left[\frac{1}{M}\sum_{i\in\mathcal{S}_t}(w_{t-1,var}^i - w_{t-1} - \eta\sum_{k=0}^{K-1}g_t^{(i,k)})\right]$$
$$\overset{(a)}{=} \mathbb{E}\left[\frac{1}{M}\sum_{i\in\mathcal{S}_t}\left(-\eta\sum_{k=0}^{K-1}\nabla f_i(w_t^{(i,k)})\right)\right]$$
$$= -\eta\sum_{k=0}^{K-1}\mathbb{E}\left[\frac{1}{M}\sum_{i\in\mathcal{S}_t}\nabla f_i(w_t^{(i,k)})\right]$$
$$= -\eta\sum_{k=0}^{K-1}\mathbb{E}[\nabla F(w_t^{(\cdot,k)})], \tag{8}$$

where in (a) we used $\mathbb{E}[c_{t-1,i}] = 0$ over the random shuffling of coefficients, which makes the variation term $\mathbb{E}[w_{t-1,var}^i - w_{t-1}]$ equal to zero. We also used the unbiased gradient assumption (A2). The inner product for the client part is then:

$$\mathbb{E}[\langle \nabla F(w_{t-1}), \bar{w}_t - w_{t-1}\rangle]$$
$$= -\eta K\mathbb{E}[\langle \nabla F(w_{t-1}), \nabla F(w_{t-1})\rangle] + \text{Drift}$$
$$= -\eta K\|\nabla F(w_{t-1})\|^2 + \mathcal{O}(\eta^2 K^2 L\zeta^2). \tag{9}$$

2

The drift term arises from the deviation of local models from the global model, and its bound is standard in FL analysis.

For the server update part of the inner product:

$$\mathbb{E}[\langle \nabla F(w_{t-1}), w_t - \bar{w}_t \rangle] = \mathbb{E}[\langle \nabla F(w_{t-1}), \lambda_t(w_t^{(0,E)} - \bar{w}_t) \rangle]$$

$$= \mathbb{E}[\langle \nabla F(w_{t-1}), \lambda_t(-\eta_0 \sum_{e=0}^{E-1} g_t^{(0,e)}) \rangle]$$

$$\leq \mathbb{E}[|\langle \nabla F(w_{t-1}), -\lambda_t \eta_0 E \nabla f_0(\bar{w}_t) \rangle|]$$

$$\leq \lambda_{\max} \eta_0 E \cdot \mathbb{E}[\|\nabla F(w_{t-1})\| \|\nabla f_0(\bar{w}_t)\|]$$

$$\leq \lambda_{\max} \eta_0 E G^2, \tag{10}$$

where we used the Cauchy-Schwarz inequality and the bounded gradient assumption (A4). Combining these, the total inner product is bounded by:

$$\mathbb{E}[\langle \ldots \rangle] \leq -\eta K \|\nabla F(w_{t-1})\|^2 + \mathcal{O}(\eta^2) + \lambda_{\max} \eta_0 E G^2. \tag{11}$$

## 2.2 Part 2: Bounding the Squared Norm Term $\mathbb{E}[\|w_t - w_{t-1}\|^2]$

We use the inequality $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$:

$$\mathbb{E}[\|w_t - w_{t-1}\|^2] \leq 2\mathbb{E}[\|\bar{w}_t - w_{t-1}\|^2] + 2\mathbb{E}[\|w_t - \bar{w}_t\|^2]. \tag{12}$$

Let's bound each term. For the client aggregation part:

$$\mathbb{E}[\|\bar{w}_t - w_{t-1}\|^2] \leq \frac{1}{M} \sum_{i \in \mathcal{S}_t} \mathbb{E}[\|w_t^{(i,K)} - w_{t-1}\|^2]$$

$$\leq \frac{2}{M} \sum_i \mathbb{E}[\|w_{t-1,var}^i - w_{t-1}\|^2] + \frac{2}{M} \sum_i \mathbb{E}[\|w_t^{(i,K)} - w_{t-1,var}^i\|^2]$$

$$\overset{(b)}{\leq} 2\rho_{\max}^2 G^2 + \frac{2}{M} \sum_i \mathbb{E}[\| - \eta \sum_{k=0}^{K-1} g_t^{(i,k)}\|^2]$$

$$\leq 2\rho_{\max}^2 G^2 + 2\eta^2 K \sum_{k=0}^{K-1} \frac{1}{M} \sum_i \mathbb{E}[\|g_t^{(i,k)}\|^2]$$

$$\leq 2\rho_{\max}^2 G^2 + 2\eta^2 K^2(\sigma^2 + G^2 + \zeta^2). \tag{13}$$

In (b), we substituted $w_{t-1,var}^i = w_{t-1} + c_{t-1,i} \cdot \rho_{t-1} \cdot \Delta w_{t-1}$ (i.e., equation (18) from the original paper), and then bounded the variation term $\|c_{t-1,i}\rho_{t-1}\Delta w_{t-1}\|^2$ using Assumption A4. The second term is a standard bound for $K$ steps of local SGD, accounting for gradient variance ($\sigma^2$), bounded true gradient ($G^2$), and non-IID divergence ($\zeta^2$).

For the server update part:

$$\mathbb{E}[\|w_t - \bar{w}_t\|^2] = \mathbb{E}[\|\lambda_t(w_t^{(0,E)} - \bar{w}_t)\|^2]$$

$$\leq \lambda_{\max}^2 \mathbb{E}[\| - \eta_0 \sum_{e=0}^{E-1} g_t^{(0,e)}\|^2]$$

$$\leq \lambda_{\max}^2 \eta_0^2 E \sum_{e=0}^{E-1} \mathbb{E}[\|g_t^{(0,e)}\|^2] \leq \lambda_{\max}^2 \eta_0^2 E^2(\sigma^2 + G^2). \tag{14}$$

3

Combining the bounds for the total squared norm:

$$\mathbb{E}[\|w_t - w_{t-1}\|^2] \leq 4\rho_{\max}^2 G^2 + 4\eta^2 K^2(\sigma^2 + G^2 + \zeta^2)$$
$$+ 2\lambda_{\max}^2 \eta_0^2 E^2(\sigma^2 + G^2). \tag{15}$$

For simplicity in the final theorem, we can absorb smaller terms into larger ones.

## 2.3  Part 3: Combining the Bounds and Finalizing the Proof

Substitute the bounds from Eq. (11) and Eq. (15) back into the descent lemma Eq. (5):

$$\mathbb{E}[F(w_t)] \leq \mathbb{E}[F(w_{t-1})] - \eta K\|\nabla F(w_{t-1})\|^2 + \lambda_{\max}\eta_0 E G^2$$
$$+ \frac{L}{2}[4\rho_{\max}^2 G^2 + 4\eta^2 K^2(\sigma^2/M + \zeta^2 + G^2)$$
$$+ 2\lambda_{\max}^2 \eta_0^2 E^2(\sigma^2 + G^2)]. \tag{16}$$

Rearranging to isolate the gradient term and simplifying higher-order $\eta$ terms:

$$\eta K\|\nabla F(w_{t-1})\|^2 \leq \mathbb{E}[F(w_{t-1}) - F(w_t)] + \lambda_{\max}\eta_0 E G^2$$
$$+ 2L\eta^2 K^2\zeta^2 + \frac{2L\eta^2 K^2\sigma^2}{M}$$
$$+ 2L\rho_{\max}^2 G^2 + L\lambda_{\max}^2 \eta_0^2 E^2 G^2 + \dots \tag{17}$$

Now, we sum this inequality from $t = 1$ to $T$:

$$\eta K \sum_{t=1}^{T} \mathbb{E}[\|\nabla F(w_{t-1})\|^2] \leq \sum_{t=1}^{T} \mathbb{E}[F(w_{t-1}) - F(w_t)]$$
$$+ T \cdot (\text{Error Terms per round}). \tag{18}$$

The first term on the right is a telescoping sum: $\sum_{t=1}^{T} \mathbb{E}[F(w_{t-1}) - F(w_t)] = \mathbb{E}[F(w_0) - F(w_T)] \leq F(w_0) - F^*$, where $F^*$ is the minimum value of $F(w)$. Dividing both sides by $\eta KT$:

$$\frac{1}{T}\sum_{t=0}^{T-1} \mathbb{E}[\|\nabla F(w_t)\|^2] \leq \frac{F(w_0) - F^*}{\eta KT} + 2L\eta K\zeta^2 + \frac{2L\eta K\sigma^2}{M}$$
$$+ \frac{\lambda_{\max}\eta_0 E G^2}{\eta K} + \frac{2L\rho_{\max}^2 G^2}{\eta K} + \frac{L\lambda_{\max}^2 \eta_0^2 E^2 G^2}{\eta K}. \tag{19}$$

By setting $\eta = \eta_0 = c/\sqrt{T}$ for a small constant $c$, the first term becomes $\mathcal{O}(1/\sqrt{T})$. The other error terms are either constants (multiplied by $\eta$ or $\eta_0$, making them $\mathcal{O}(1/\sqrt{T})$) or $\mathcal{O}(1/T)$. Thus, the dominant term dictating the convergence rate is $\mathcal{O}(1/\sqrt{T})$. This completes the proof. $\square$