

Robust Federated Learning: The Case of Affine Distribution Shifts

Appendix A Additional Numerical Experiments

A.1 Experimental Setup

In the experiments, we simulated a federated learning scenario with $n = 10$ nodes where each node observes $m = 5000$ training samples. We also divided the extra 10,000 samples in each dataset to two validation and test sets containing 5000 samples each. For CIFAR-10 samples, we applied the standard normalization and scaled and linearly mapped the pixel intensity values to interval $[-1, 1]$. We applied batch normalization [44] in order to stabilize training and used the ADAM optimizer [45] with stepsize value 10^{-4} and default beta parameters $\beta_1 = 0.9$ and $\beta_2 = 0.99$ to optimize the neural net's parameters for $T = 100$ epochs (10000 iterations).

We did cross validation to choose $\lambda \in \{0.1, 0.5, 1, 5, 10, 50\}$ and chose the λ -value resulting in the closest additive penalty $\frac{1}{n} \sum_{i=1}^n [\|\Lambda^{i*} - I\|_2^2 + \|\delta^{i*}\|_2^2]$ to 10 percent of the average sample norm, i.e. $\frac{0.1}{m} \sum_{i=1}^m \|\mathbf{x}_i^{\text{val}}\|_2^2$, over the $m = 5000$ validation samples. To perform GDA optimization, we applied two ascent steps per descent step with stepsize $\frac{1}{2\lambda}$. In order to simulate an affine distribution shift, we manipulated each $\tilde{\mathbf{x}}_j^i$ in the original training dataset via an affine transformation chosen randomly at each node:

$$\mathbf{x}_j^i = (I_d + \tilde{\Lambda}^i) \tilde{\mathbf{x}}_j^i + \tilde{\delta}^i. \quad (9)$$

Here, each $\tilde{\Lambda}^i$ is a random matrix with i.i.d. Gaussian entries according to $\mathcal{N}(0, \frac{\sigma^2}{d})$, and $\tilde{\delta}^i$ is a random Gaussian vector according to $\mathcal{N}(0, \sigma^2 I_d)$ where we set $\sigma = 0.01$. In test time, we did not apply any random affine transformation to test samples and instead considered the following three scenarios: (1) no perturbation, (2) adversarial affine distribution shift obtained by optimizing the inner maximization in (1) using projected gradient descent, 3) adversarial perturbations designed by the projected gradient descent algorithm. We used 100 projected gradient steps with stepsize 0.1.

We considered three baselines in the experiments: (1) FedAvg where the server node averages the updated parameters of the local nodes after every gradient step, (2) Distributed FGM training where the nodes perform fast adversarial training [9] by optimizing an ℓ_2 -norm bounded perturbation δ_j^i using one gradient step followed by projection onto the ball $\{\delta_j^i : \|\delta_j^i\|_2 \leq \epsilon_{\text{fgm}}\}$, and (3) Distributed PGD training where each node performs PGD adversarial training [8] similar to distributed FGM but uses 10 projected gradient steps, each followed by projection onto $\{\delta_j^i : \|\delta_j^i\|_2 \leq \epsilon_{\text{pgd}}\}$. We used the value $\epsilon_{\text{fgm}} = \epsilon_{\text{pgd}} = 0.05 \mathbb{E}[\|\mathbf{x}_i\|_2]$ in the experiments. We observed training instability after achieving perfect training accuracy for the baseline FedAvg algorithm, and hence performed early stopping to avoid the instability in the FedAvg experiments. We did not encounter the instability issue in FedRobust experiments.

A.2 Numerical Results for MNIST data

We repeated the CIFAR experiments in Figures 1 and 2 for the MNIST dataset. Figure 5 shows the numerical results under affine distribution shifts. The figure's top row includes the plots for fixed maximum delta norm $\|\delta\|_2 \leq 1$ and different levels of maximum allowed $\|\Lambda - I\|_F$, while in the bottom row we fix the maximum allowed linear shift $\|\Lambda - I\|_F \leq 0.6$ and evaluate the test accuracy under different levels of $\|\delta\|_2$. As shown in the plots, FedRobust results in the best performance in most of the evaluations, which indicates the superior performance of FedRobust against affine distribution shifts. Figure 6 shows the test accuracy of the trained networks under different levels of adversarial PGD perturbations. The figure's experiments again shows that FedRobust can effectively shield against PGD adversarial attacks and achieve a comparable performance to PGD and FGM adversarial training.

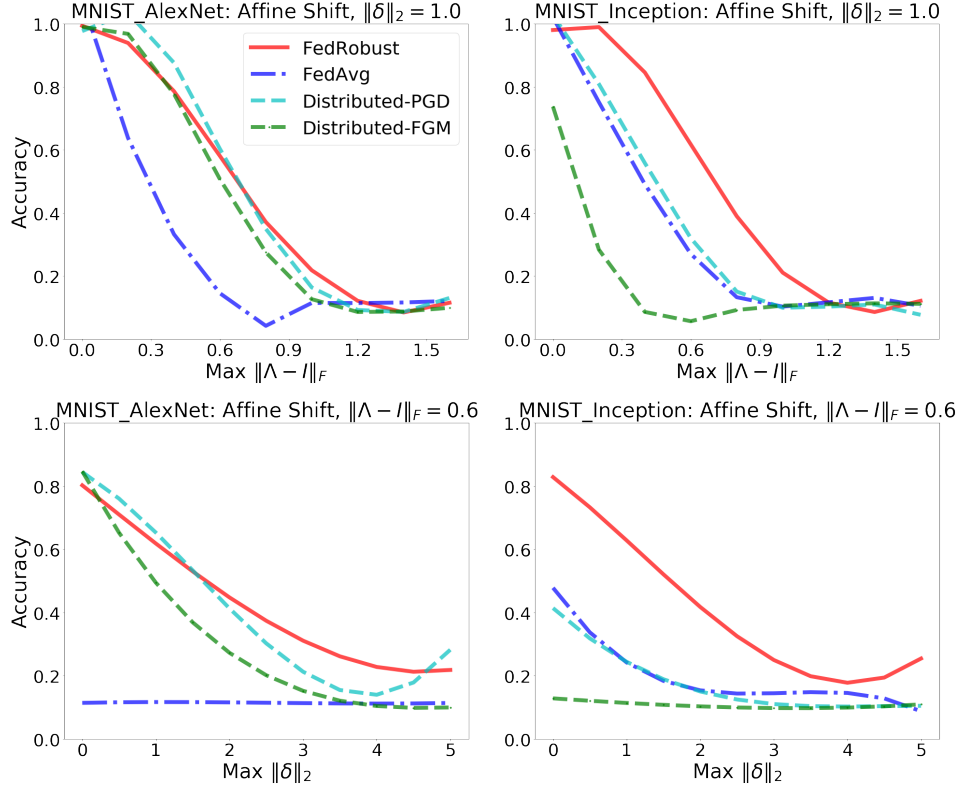


Figure 5: Trained networks' test accuracy under affine distribution shifts in the MNIST experiments. Top row: constraining $\|\delta\|_2 \leq 1$ and changing maximum allowed $\|\Lambda - I\|_F$, bottom row: constraining $\|\Lambda - I\|_F \leq 0.6$ and changing maximum allowed $\|\delta\|_2$.

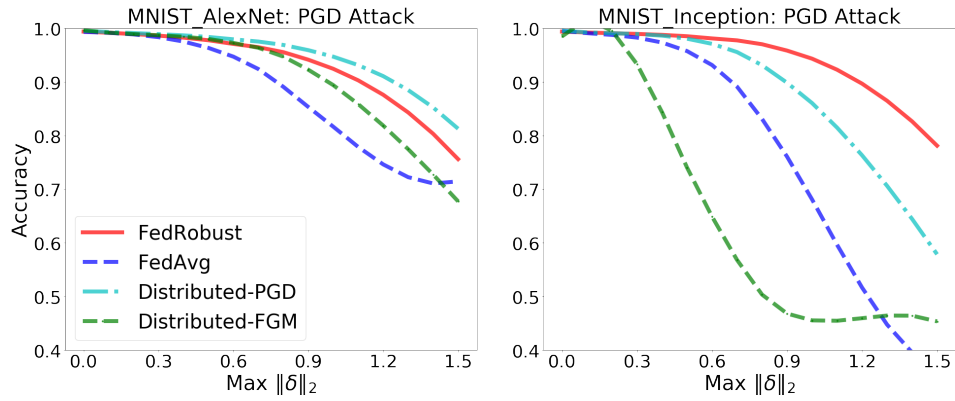


Figure 6: Trained networks' test accuracy under PGD perturbations in the MNIST experiments. X-axis shows the maximum allowed ℓ_2 -norm for PGD perturbations.

Appendix B Preliminaries and Useful Lemmas

In this section, we provide preliminary and useful results in order to prove Theorems 1 and 2. For notational convenience, we use the following short-hand notations:

Notation	Description
$\psi_t^i = (\Lambda_t^i, \delta_t^i)$	maximization variables of node i iteration t
$\Psi_t = (\psi_t^1; \dots; \psi_t^n)$	concatenation of all nodes' maximization models at iteration t
$\bar{\mathbf{w}}_t = \frac{1}{n} \sum_{i \in [n]} \mathbf{w}_t^i$	average model at iteration t
$a_t = \mathbb{E}[\Phi(\bar{\mathbf{w}}_t)] - \Phi^*$	optimality gap measure between $\Phi(\bar{\mathbf{w}}_t)$ and $\min_{\mathbf{w}} \Phi(\mathbf{w})$
$b_t = \mathbb{E}[\Phi(\bar{\mathbf{w}}_t) - f(\bar{\mathbf{w}}_t, \Psi_t)]$	optimality gap measure between $f(\bar{\mathbf{w}}_t, \Psi_t)$ and $\max_{\Psi} f(\bar{\mathbf{w}}_t, \Psi)$
$e_t = \frac{1}{n} \sum_{i \in [n]} \mathbb{E} \ \mathbf{w}_t^i - \bar{\mathbf{w}}_t\ ^2$	average deviation of the local models from the average model at iteration t
$g_t = \mathbb{E} \left\ \frac{1}{n} \sum_{i \in [n]} \nabla_{\mathbf{w}} f^i(\mathbf{w}_t^i, \psi_t^i) \right\ ^2$	norm squared of local gradients w.r.t \mathbf{w} at iteration t
$h_t = \mathbb{E} \left\ \nabla \Phi(\bar{\mathbf{w}}_t) - \frac{1}{n} \sum_{i \in [n]} \nabla_{\mathbf{w}} f^i(\mathbf{w}_t^i, \psi_t^i) \right\ ^2$	norm squared of deviation in gradients w.r.t \mathbf{w} of $\max_{\Psi} f(\bar{\mathbf{w}}_t, \Psi)$ and local functions $f^i(\mathbf{w}_t^i, \psi_t^i)$

Table 1: Table of notations.

Now, we present a set of useful lemmas and observations which we will invoke to prove the convergence results for both PL-PL and nonconvex-PL loss cases. The following lemma establishes the Lipschitz gradient parameter for the global function given those of the local objectives.

Lemma 1. *If the local functions f^i 's have Lipschitz gradients with parameters stated in Assumption 3, then the global function f has also Lipschitz gradients as follows: for any $\mathbf{w}, \mathbf{w}', \Psi, \Psi'$ it holds that*

$$\begin{aligned} \|\nabla_{\mathbf{w}} f(\mathbf{w}, \Psi) - \nabla_{\mathbf{w}} f(\mathbf{w}', \Psi)\| &\leq L_1 \|\mathbf{w} - \mathbf{w}'\|, \|\nabla_{\mathbf{w}} f(\mathbf{w}, \Psi) - \nabla_{\mathbf{w}} f(\mathbf{w}, \Psi')\| \leq \frac{L_{12}}{\sqrt{n}} \|\Psi - \Psi'\|_F, \\ \|\nabla_{\Psi} f(\mathbf{w}, \Psi) - \nabla_{\Psi} f(\mathbf{w}', \Psi)\|_F &\leq \frac{L_{21}}{\sqrt{n}} \|\mathbf{w} - \mathbf{w}'\|, \|\nabla_{\Psi} f(\mathbf{w}, \Psi) - \nabla_{\Psi} f(\mathbf{w}, \Psi')\|_F \leq \frac{L_2}{n} \|\Psi - \Psi'\|_F. \end{aligned} \quad (10)$$

Proof. We defer the proof to Section E.1. □

Recall the definition of the function $\Phi(\cdot)$, that is,

$$\Phi(\mathbf{w}) := \max_{\Psi} f(\mathbf{w}, \Psi) = \max_{\psi^1, \dots, \psi^n} \frac{1}{n} \sum_{i \in [n]} f^i(\mathbf{w}, \psi^i) = \max_{(\Lambda^1, \delta^1), \dots, (\Lambda^n, \delta^n)} \frac{1}{n} \sum_{i \in [n]} f^i(\mathbf{w}, \Lambda^i, \delta^i). \quad (11)$$

Next lemma shows that Φ has Lipschitz gradients and characterizes its parameter.

Lemma 2 ([32]). *If Assumptions 3 and 4 (ii) hold, that is, the local objectives have Lipschitz gradients and $-f(\mathbf{w}, \cdot)$ is μ_2 -PL, then we have*

$$\nabla \Phi(\mathbf{w}) = \nabla_{\mathbf{w}} f(\mathbf{w}, \Psi^*(\mathbf{w})), \quad (12)$$

where $\Psi^*(\mathbf{w}) \in \arg \max_{\Psi} f(\mathbf{w}, \Psi)$ for any \mathbf{w} . Moreover, Φ has Lipschitz gradients with parameter $L_{\Phi} = L_1 + \frac{L_{12}L_{21}}{2n\mu_2}$.

Proof. We defer the proof to Section E.2. \square

Next lemma shows the contraction of the sequence $\{\mathbb{E}[\Phi(\bar{\mathbf{w}}_t)]\}_{t \geq 0}$ when running the update rule of FedRobust method in Algorithm 1. Please refer to Table 1 to recall the definition of h_t and g_t .

Lemma 3. *If Assumptions 2 and 3 hold, then the iterates of FedRobust satisfy the following contraction inequality for any iteration $t \geq 0$*

$$\mathbb{E}[\Phi(\bar{\mathbf{w}}_{t+1})] - \mathbb{E}[\Phi(\bar{\mathbf{w}}_t)] \leq -\frac{\eta_1}{2} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 + \frac{\eta_1}{2} h_t - \frac{\eta_1}{2} (1 - \eta_1 L_\Phi) g_t + \eta_1^2 \frac{L_\Phi}{2} \frac{\sigma_{\mathbf{w}}^2}{n}. \quad (13)$$

Proof. We defer the proof to Section E.3. \square

Next lemma further bounds h_t w.r.t. the two sequences b_t and e_t .

Lemma 4. *If Assumptions 3 and 4 (ii) hold, that is, the local objectives have Lipschitz gradients and $-f(\mathbf{w}, \cdot)$ is μ_2 -PL, then we have*

$$h_t \leq \frac{4L_{12}^2}{\mu_2 n} b_t + 2L_1^2 e_t. \quad (14)$$

Proof. We defer the proof to Section E.4. \square

Next lemma establishes a contraction bound on the sequence b_t .

Lemma 5. *If Assumptions 2, 3 and 4 (ii) hold, then the sequence of $\{b_t\}_{t \geq 0}$ generated by the FedRobust iterations with $\eta_2 \leq 1/L_2$ satisfies the following contraction bound:*

$$\begin{aligned} b_{t+1} \leq & (1 - \mu_2 \eta_2 n) \left(1 + \eta_1 \frac{4L_{12}^2}{\mu_2 n} \right) b_t + \frac{\eta_1}{2} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 + \frac{\eta_1^2}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) g_t \\ & + (\eta_1 L_1^2 + \eta_2 L_{21}^2) e_t + \frac{\eta_1^2}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \frac{\sigma_{\mathbf{w}}^2}{n} + \frac{\eta_2^2}{2} L_2 \sigma_\psi^2, \end{aligned} \quad (15)$$

where L_Φ is the Lipschitz gradient parameter of the function $\Phi(\cdot)$ characterized in Lemma 2.

Proof. We defer the proof to Section E.5. \square

Next lemma bounds e_t , that is the average deviation of local parameter models from their average.

Lemma 6. *If Assumptions 1, 2 and 3 hold and the step-size η_1 satisfies $32\eta_1^2(\tau - 1)^2 L_1^2 \leq 1$, then the sequence $e_t = \frac{1}{n} \sum_{i \in [n]} \mathbb{E} \|\mathbf{w}_t^i - \bar{\mathbf{w}}_t\|^2$ is bounded as follows*

$$e_t \leq 16\eta_1^2(\tau - 1)^2 \rho^2 + 4\eta_1^2(\tau - 1)(n + 1) \frac{\sigma_{\mathbf{w}}^2}{n} + 20\eta_1^2(\tau - 1) \sum_{l=t_c+1}^{t-1} g_l, \quad (16)$$

where t_c denotes the index of the most recent server-worker communication, i.e. $t_c = \lfloor \frac{t}{\tau} \rfloor \tau$ and we also denote $\rho^2 := 3\rho_f^2 + 6L_{12}^2(\epsilon_1^2 + \epsilon_2^2)$.

Proof. We defer the proof to Section E.6. \square

Next generic lemma is adopted from [16].

Lemma 7. *Assume that two non-negative sequences $\{P_t\}_{t \geq 0}$ and $\{g_t\}_{t \geq 0}$ satisfy the following inequality for each iteration $t \geq 0$ and some constants $0 < \Upsilon < 1$, $L \geq 0$, $B \geq 0$, and $\Gamma \geq 0$:*

$$P_{t+1} \leq \Upsilon P_t - \frac{\eta_1}{2} (1 - \eta_1 L) g_t + \eta_1^2 B \sum_{l=t_c+1}^{t-1} g_l + \Gamma, \quad (17)$$

where $t_c = \lfloor \frac{t}{\tau} \rfloor \tau$. Then, for each $t \geq 0$ we have

$$P_t \leq \Upsilon^t P_0 + \frac{\Gamma}{1 - \Upsilon}, \quad (18)$$

if η_1 satisfies the following condition

$$\eta_1 \left(L + \frac{2B}{\Upsilon^{\tau-1}(1-\Upsilon)} \right) \leq 1. \quad (19)$$

Proof. We defer the proof to Section E.7. \square

Next lemma bounds the overall optimality gap b_t averaged over T iterations.

Lemma 8. *If Assumptions 2, 3 and 4 (ii) hold and the step-sizes satisfy the conditions $\eta_2 \leq 1/L_2$ and $\frac{\eta_2}{\eta_1} \geq \frac{8L_{12}^2}{\mu_2^2 n^2}$, then the average of the sequence $\{b_t\}_{t=0}^{T-1}$ generated from the FedRobust can be bounded as follows:*

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} b_t &\leq \frac{4L_2^2}{\mu_2^2 n^2} \frac{\epsilon_1^2 + \epsilon_2^2}{\eta_2 T} + \frac{\eta_1}{\eta_2} \frac{1}{\mu_2 n} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 \\ &\quad + \frac{\eta_1^2}{\eta_2} \frac{1}{\mu_2 n} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \frac{1}{T} \sum_{t=0}^{T-1} g_t + \frac{1}{\eta_2} \frac{2}{\mu_2 n} (\eta_1 L_1^2 + \eta_2 L_{21}^2) \frac{1}{T} \sum_{t=0}^{T-1} e_t \\ &\quad + \frac{\eta_1^2}{\eta_2} \frac{1}{\mu_2 n} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \frac{\sigma_{\mathbf{w}}^2}{n} + \eta_2 \frac{L_2}{\mu_2 n} \sigma_\psi^2, \end{aligned} \quad (20)$$

where L_Φ is the Lipschitz gradient parameter of the function $\Phi(\cdot)$ characterized in Lemma 2 and ϵ_1, ϵ_2 represent the radius of the affine perturbation balls, i.e. $\|\Lambda^i - I\| \leq \epsilon_1$ and $\|\delta^i\| \leq \epsilon_2$ for each node $i \in [n]$.

Proof. We defer the proof to Section E.8. \square

Next lemma bounds the averaged local model deviations e_t over T iterations.

Lemma 9. *If Assumptions 1, 2 and 3 hold and the step-size η_1 satisfies $32\eta_1^2(\tau-1)^2 L_1^2 \leq 1$, then the average of the sequence e_t over $t = 0, \dots, T-1$ is bounded as follows*

$$\frac{1}{T} \sum_{t=0}^{T-1} e_t \leq 20\eta_1^2(\tau-1)^2 \frac{1}{T} \sum_{t=0}^{T-1} g_t + 16\eta_1^2(\tau-1)^2 \rho^2 + 8\eta_1^2(\tau-1)(n+1) \frac{\sigma_{\mathbf{w}}^2}{n}. \quad (21)$$

Proof. We defer the proof to Section E.9. \square

Appendix C Proof of Theorem 1

Having established the key lemmas, now we proceed to prove Theorem 1 for any $\beta \leq 1/2$. To show the convergence of the sequence $P_t = a_t + \beta b_t$, we firstly need to establish a contraction inequality on P_{t+1} with respect to P_t . We begin by the following bound on the sequence $a_t = \mathbb{E}[\Phi(\bar{\mathbf{w}}_t)] - \Phi^*$ which is directly implied from Lemma 3:

$$a_{t+1} \leq a_t - \frac{\eta_1}{2} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 + \frac{\eta_1}{2} h_t - \frac{\eta_1}{2} (1 - \eta_1 L_\Phi) g_t + \eta_1^2 \frac{L_\Phi}{2} \frac{\sigma_{\mathbf{w}}^2}{n}. \quad (22)$$

Using Lemma 4 that shows $h_t \leq 4L_{12}^2 b_t / (\mu_2 n) + 2L_1^2 e_t$, the bound in (22) yields that

$$a_{t+1} \leq a_t - \frac{\eta_1}{2} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 + \eta_1 \frac{2L_{12}^2}{\mu_2 n} b_t + \eta_1 L_1^2 e_t - \frac{\eta_1}{2} (1 - \eta_1 L_\Phi) g_t + \eta_1^2 \frac{L_\Phi}{2} \frac{\sigma_{\mathbf{w}}^2}{n}. \quad (23)$$

Next, we employ the result of Lemma 5 which establishes a contraction bound on the b_t sequence. Putting together with (23) implies that

$$\begin{aligned}
P_{t+1} &= a_{t+1} + \beta b_{t+1} \\
&\leq a_t - \frac{\eta_1}{2} (1 - \beta) \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 \\
&\quad + \beta \left(\eta_1 \frac{2L_{12}^2}{\beta \mu_2 n} + (1 - \mu_2 \eta_2 n) \left(1 + \eta_1 \frac{4L_{12}^2}{\mu_2 n} \right) \right) b_t \\
&\quad - \left(\frac{\eta_1}{2} (1 - \eta_1 L_\Phi) - \eta_1^2 \frac{\beta}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \right) g_t \\
&\quad + \left(\eta_1 L_1^2 + \beta (\eta_1 L_1^2 + \eta_2 L_{21}^2) \right) e_t \\
&\quad + \frac{\eta_1^2}{2} \left(L_\Phi + \beta (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \right) \frac{\sigma_{\mathbf{w}}^2}{n} + \eta_2^2 L_2 \frac{\beta}{2} \sigma_\psi^2. \tag{24}
\end{aligned}$$

We begin simplifying the above bound by first considering the first two terms in RHS of (24). We can show that the function $\Phi(\cdot)$ is μ_1 -PL [31], which implies that

$$\mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 \geq 2\mu_1 \mathbb{E} [\Phi(\bar{\mathbf{w}}_t)] - \Phi^* = 2\mu_1 a_t. \tag{25}$$

Therefore, for any $\beta \leq 1/2$ we have

$$a_t - \frac{\eta_1}{2} (1 - \beta) \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 \leq \left(1 - \frac{1}{2} \mu_1 \eta_1 \right) a_t, \tag{26}$$

which implies the coefficient of a_t in (24) is bounded by $1 - \frac{1}{2} \mu_1 \eta_1$. Next, the coefficient of βb_t in (24) can be bounded as follows:

$$\begin{aligned}
\eta_1 \frac{2L_{12}^2}{\beta \mu_2 n} + (1 - \mu_2 \eta_2 n) \left(1 + \eta_1 \frac{4L_{12}^2}{\mu_2 n} \right) &= 1 - \eta_1 \frac{L_1 L_2}{\mu_2 n} \left(\frac{\mu_2^2 \eta_2 n}{\eta_1 L_1 L_2} - \frac{2L_{21}^2}{\beta L_1 L_2} - 4(1 - \mu_2 \eta_2 n) \frac{L_{21}^2}{L_1 L_2} \right) \\
&\stackrel{(a)}{\leq} 1 - \eta_1 \frac{L_1 L_2}{\mu_2 n} \\
&\stackrel{(b)}{\leq} 1 - \frac{1}{2} \mu_1 \eta_1, \tag{27}
\end{aligned}$$

where (a) holds for our choice of β and assuming $\frac{\mu_2^2 \eta_2 n}{\eta_1 L_1 L_2} \geq 1 + (4 + \frac{2}{\beta}) \frac{L_{12}^2}{L_1 L_2}$ and (b) is implies from the fact that

$$\frac{\eta_1 \frac{L_1 L_2}{\mu_2 n}}{\frac{1}{2} \mu_1 \eta_1} = 2 \left(\frac{L_1}{\mu_1} \right) \left(\frac{L_2}{\mu_2 n} \right) \geq 1. \tag{28}$$

Now that we have bounded the coefficients of a_t and βb_t in (24), rearranging the terms and using the assumption $\eta_2 \leq 1/L_2$ simplifies the contraction on P_t as follows

$$P_{t+1} \leq \left(1 - \frac{1}{2} \mu_1 \eta_1 \right) P_t - \frac{\eta_1}{2} (1 - \eta_1 \hat{L}_\beta) g_t + \tilde{L}_\beta e_t + \eta_1^2 \frac{\hat{L}_\beta}{2} \frac{\sigma_{\mathbf{w}}^2}{n} + \eta_2^2 \frac{L_2}{2} \beta \sigma_\psi^2, \tag{29}$$

where we picked the following notations for convenient of the exposition

$$\tilde{L}_\beta = (1 + \beta) \eta_1 L_1^2 + \beta \eta_2 L_{21}^2, \quad \hat{L}_\beta = (1 + \beta) L_\Phi + \beta L_1 + 2\beta \frac{L_{21}^2}{L_2}. \tag{30}$$

Next, we use Lemma 6 which for $32\eta_1^2(\tau - 1)^2 L_1^2 \leq 1$ provides an upper bound on e_t with respect to g_t . We can write

$$\begin{aligned}
P_{t+1} &\leq \left(1 - \frac{1}{2} \mu_1 \eta_1 \right) P_t - \frac{\eta_1}{2} (1 - \eta_1 \hat{L}_\beta) g_t + 20\eta_1^2 \tilde{L}_\beta (\tau - 1) \sum_{l=t_c+1}^{t-1} g_l \\
&\quad + 16\eta_1^2 \tilde{L}_\beta (\tau - 1)^2 \rho^2 + 4\eta_1^2 \tilde{L}_\beta (\tau - 1) (n + 1) \frac{\sigma_{\mathbf{w}}^2}{n} + \eta_1^2 \frac{\hat{L}_\beta}{2} \frac{\sigma_{\mathbf{w}}^2}{n} + \eta_2^2 \frac{L_2}{2} \beta \sigma_\psi^2. \tag{31}
\end{aligned}$$

We have shown in Lemma 7 that how a such contraction sequence converges. In particular, let us pick the following notations and apply the result of Lemma 7 to contraction in (31)

$$\begin{aligned}
L &= \hat{L}_\beta, \\
\Upsilon &= 1 - \frac{1}{2}\mu_1\eta_1, \\
B &= 20\tilde{L}_\beta(\tau - 1), \\
\Gamma &= 16\eta_1^2\tilde{L}_\beta(\tau - 1)^2\rho^2 + 4\eta_1^2\tilde{L}_\beta(\tau - 1)(n + 1)\frac{\sigma_{\mathbf{w}}^2}{n} + \eta_1^2\frac{\hat{L}_\beta}{2}\frac{\sigma_{\mathbf{w}}^2}{n} + \eta_2^2\frac{L_2}{2}\beta\sigma_\psi^2.
\end{aligned} \tag{32}$$

It implies that if the step-sizes satisfy the following condition

$$\eta_1 \left(\hat{L}_\beta + \frac{80\tilde{L}_\beta(\tau - 1)}{\eta_1\mu_1 \left(1 - \frac{1}{2}\mu_1\eta_1\right)^{\tau-1}} \right) \leq 1, \tag{33}$$

then we have

$$P_t \leq \left(1 - \frac{1}{2}\mu_1\eta_1\right)^t P_0 + 32\eta_1\frac{\tilde{L}_\beta}{\mu_1}(\tau - 1)^2\rho^2 + 8\eta_1\frac{\tilde{L}_\beta}{\mu_1}(\tau - 1)(n + 1)\frac{\sigma_{\mathbf{w}}^2}{n} + \eta_1\frac{\hat{L}_\beta}{\mu_1}\frac{\sigma_{\mathbf{w}}^2}{n} + \frac{\eta_2^2}{\eta_1}\frac{L_2}{\mu_1}\beta\sigma_\psi^2, \tag{34}$$

which concludes the proof of Theorem 1. Note to hold this result, in addition to condition (33), we have assumed the following constraints on the step-sizes as well

$$\eta_2 L_2 \leq 1, \quad 32\eta_1^2(\tau - 1)^2 L_1^2 \leq 1, \quad \frac{\mu_2^2 \eta_2 n}{\eta_1 L_1 L_2} \geq 1 + \left(4 + \frac{2}{\beta}\right) \frac{L_{12}^2}{L_1 L_2}. \tag{35}$$

Appendix D Proof of Theorem 2

We begin the proof by combining the results of Lemmas 3 and 4 which yields that for every iteration $t = 0, \dots, T - 1$ we have

$$\mathbb{E}[\Phi(\bar{\mathbf{w}}_{t+1})] - \mathbb{E}[\Phi(\bar{\mathbf{w}}_t)] \leq -\frac{\eta_1}{2}\mathbb{E}\|\nabla\Phi(\bar{\mathbf{w}}_t)\|^2 - \frac{\eta_1}{2}(1 - \eta_1 L_\Phi)g_t + \eta_1\frac{2L_{12}^2}{\mu_2 n}b_t + \eta_1 L_1^2 e_t + \eta_1^2\frac{L_\Phi}{2}\frac{\sigma_{\mathbf{w}}^2}{n}. \tag{36}$$

Summing up all the T inequalities in (36) for $t = 0, \dots, T - 1$ and dividing by T yields the following

$$\begin{aligned}
\frac{1}{T}(\mathbb{E}[\Phi(\bar{\mathbf{w}}_T)] - \Phi(\bar{\mathbf{w}}_0)) &\leq -\frac{\eta_1}{2}\frac{1}{T}\sum_{t=0}^{T-1}\mathbb{E}\|\nabla\Phi(\bar{\mathbf{w}}_t)\|^2 \\
&\quad - \frac{\eta_1}{2}(1 - \eta_1 L_\Phi)\frac{1}{T}\sum_{t=0}^{T-1}g_t \\
&\quad + \eta_1\frac{2L_{12}^2}{\mu_2 n}\frac{1}{T}\sum_{t=0}^{T-1}b_t \\
&\quad + \eta_1 L_1^2\frac{1}{T}\sum_{t=0}^{T-1}e_t \\
&\quad + \eta_1^2\frac{L_\Phi}{2}\frac{\sigma_{\mathbf{w}}^2}{n}.
\end{aligned} \tag{37}$$

Next we use Lemmas 8 and then Lemma 9 to replace the terms $\frac{1}{T} \sum_{t=0}^{T-1} b_t$ and $\frac{1}{T} \sum_{t=0}^{T-1} e_t$ and rewrite the above bound in terms of $\frac{1}{T} \sum_{t=0}^{T-1} g_t$. It yields that

$$\begin{aligned} \frac{1}{T} (\mathbb{E}[\Phi(\bar{\mathbf{w}}_T)] - \Phi(\bar{\mathbf{w}}_0)) &\leq -\frac{\eta_1}{2} \left(1 - \eta_1 \frac{4L_{12}^2 L_2}{\mu_2^2 n^2} \right) \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 \\ &\quad - \frac{\eta_1}{2} \left(1 - \eta_1 (\hat{L} + 40\tilde{L}(\tau-1)^2) \right) \frac{1}{T} \sum_{t=0}^{T-1} g_t \\ &\quad + \frac{\eta_1}{\eta_2} \frac{8L_{12}^2 L_2^2}{\mu_2^3 n^3} \frac{\epsilon_1^2 + \epsilon_2^2}{T} + 16\eta_1^2 \tilde{L}(\tau-1)^2 \rho^2 + \frac{\eta_1^2}{2} \hat{L} \frac{\sigma_{\mathbf{w}}^2}{n} + \eta_1 \eta_2 \frac{4L_{12}^2}{\mu_2^2 n^2} \hat{L} \sigma_{\psi}^2, \end{aligned} \quad (38)$$

where we adopt the following short-hand notations

$$\tilde{L} = \frac{3}{2} \eta_1 L_1^2 + \frac{1}{2} \eta_2 L_{21}^2, \quad \hat{L} = \frac{3}{2} L_{\Phi} + \frac{1}{2} L_1 + \frac{L_{21}^2}{L_2}. \quad (39)$$

Finally, we use the assumption $\eta_1 (\hat{L} + 40\tilde{L}(\tau-1)^2) \leq 1$ to remove the term $\frac{1}{T} \sum_{t=0}^{T-1} g_t$ and apply $\frac{\eta_1}{\eta_2} \leq \frac{\mu_2^2 n^2}{8L_{12}^2}$ to simply the bound and conclude the proof:

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 &\leq \frac{4\Delta_{\Phi}}{\eta_1 T} + \frac{4L_2^2}{\mu_2^2 n^2} \frac{\epsilon_1^2 + \epsilon_2^2}{\eta_1 T} + 64\eta_1 \tilde{L}(\tau-1)^2 \rho^2 \\ &\quad + 16\eta_1 \tilde{L}(\tau-1)(n+1) \frac{\sigma_{\mathbf{w}}^2}{n} + 2\eta_1 \hat{L} \frac{\sigma_{\mathbf{w}}^2}{n} + \frac{\eta_2^2}{\eta_1} L_2 \sigma_{\psi}^2. \end{aligned} \quad (40)$$

Appendix E Proof of Useful Lemmas

E.1 Proof of Lemma 1

Proof of all four cases in the claim is simple. We derive the proof for the fourth one as an instance. Recall definition of the global function f , that is

$$f(\mathbf{w}, \Psi) = \frac{1}{n} \sum_{i \in [n]} f^i(\mathbf{w}, \psi^i). \quad (41)$$

Therefore, the gradient of f with respect to Ψ is

$$\nabla_{\Psi} f(\mathbf{w}, \Psi) = \begin{pmatrix} \frac{\partial}{\partial \psi^1} f(\mathbf{w}, \Psi) \\ \vdots \\ \frac{\partial}{\partial \psi^n} f(\mathbf{w}, \Psi) \end{pmatrix} = \frac{1}{n} \begin{pmatrix} \nabla_{\psi} f^1(\mathbf{w}, \psi^1) \\ \vdots \\ \nabla_{\psi} f^n(\mathbf{w}, \psi^n) \end{pmatrix}. \quad (42)$$

We can then write for any $\mathbf{w}, \Psi = (\psi^1; \dots; \psi^n)$, $\Psi' = (\psi'^1; \dots; \psi'^n)$ and using Assumption 3 that

$$\begin{aligned} \|\nabla_{\Psi} f(\mathbf{w}, \Psi) - \nabla_{\Psi} f(\mathbf{w}, \Psi')\|_F^2 &= \frac{1}{n^2} \sum_{i \in [n]} \|\nabla_{\psi} f^i(\mathbf{w}, \psi^i) - \nabla_{\psi} f^i(\mathbf{w}, \psi'^i)\|_F^2 \\ &\leq \frac{L_2^2}{n^2} \sum_{i \in [n]} \|\psi^i - \psi'^i\|_F^2 \\ &= \frac{L_2^2}{n^2} \|\Psi - \Psi'\|_F^2. \end{aligned} \quad (43)$$

E.2 Proof of Lemma 2

The detailed proof can be found in [32], Lemma A.5. Note that in our case, according to Lemma 1 the function f has Lipschitz gradients with constants $L_1, L_{12}/\sqrt{n}, L_{21}/\sqrt{n}, L_2/n$; implying the Lipschitz gradient parameter of the function Φ to be

$$L_{\Phi} = L_1 + \frac{(L_{12}/\sqrt{n})(L_{21}/\sqrt{n})}{2\mu_2} = L_1 + \frac{L_{12}L_{21}}{2n\mu_2}. \quad (44)$$

E.3 Proof of Lemma 3

We invoke Lemma 2 which shows that the gradient of the function $\Phi(\cdot)$ is L_Φ -Lipschitz. We can write

$$\begin{aligned} \Phi(\bar{\mathbf{w}}_{t+1}) - \Phi(\bar{\mathbf{w}}_t) &\leq \langle \nabla \Phi(\bar{\mathbf{w}}_t), \bar{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t \rangle + \frac{L_\Phi}{2} \|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t\|^2 \\ &= -\eta_1 \left\langle \nabla \Phi(\bar{\mathbf{w}}_t), \frac{1}{n} \sum_{i \in [n]} \tilde{\nabla}_{\mathbf{w}} f^i(\mathbf{w}_t^i, \boldsymbol{\psi}_t^i) \right\rangle + \eta_1^2 \frac{L_\Phi}{2} \left\| \frac{1}{n} \sum_{i \in [n]} \tilde{\nabla}_{\mathbf{w}} f^i(\mathbf{w}_t^i, \boldsymbol{\psi}_t^i) \right\|^2, \end{aligned} \quad (45)$$

where we use the update rule of FedRobust and note that the difference of averaged models can be written as $\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t = -\eta_1 \frac{1}{n} \sum_{i \in [n]} \tilde{\nabla}_{\mathbf{w}} f^i(\mathbf{w}_t^i, \boldsymbol{\psi}_t^i)$. Moreover, since the stochastic gradients $\tilde{\nabla}_{\mathbf{w}} f^i$ are unbiased and variance-bounded by $\sigma_{\mathbf{w}}^2$, we can take expectation from both sides of (45) and further simplify it as follows

$$\mathbb{E}[\Phi(\bar{\mathbf{w}}_{t+1}) - \Phi(\bar{\mathbf{w}}_t)] \leq -\frac{\eta_1}{2} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 + \frac{\eta_1}{2} h_t - \frac{\eta_1}{2} (1 - \eta_1 L_\Phi) g_t + \eta_1^2 \frac{L_\Phi}{2} \frac{\sigma_{\mathbf{w}}^2}{n}. \quad (46)$$

In above, we used the inequality $2\langle \mathbf{a}, \mathbf{b} \rangle = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 - \|\mathbf{a} - \mathbf{b}\|^2$ as well as the notations for g_t and h_t as defined in Table 1.

E.4 Proof of Lemma 4

We begin bounding h_t by adding/subtracting the term $\nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_t, \Psi_t)$ and use the inequality $\|\mathbf{a} + \mathbf{b}\|^2 \leq 2\|\mathbf{a}\|^2 + 2\|\mathbf{b}\|^2$ to write

$$\begin{aligned} h_t &= \mathbb{E} \left\| \nabla \Phi(\bar{\mathbf{w}}_t) - \frac{1}{n} \sum_{i \in [n]} \nabla_{\mathbf{w}} f^i(\mathbf{w}_t^i, \boldsymbol{\psi}_t^i) \right\|^2 \\ &\leq 2\mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t) - \nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_t, \Psi_t)\|^2 + 2\mathbb{E} \left\| \nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_t, \Psi_t) - \frac{1}{n} \sum_{i \in [n]} \nabla_{\mathbf{w}} f^i(\mathbf{w}_t^i, \boldsymbol{\psi}_t^i) \right\|^2. \end{aligned} \quad (47)$$

The first term in RHS of (47) can be bounded as follows:

$$\begin{aligned} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t) - \nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_t, \Psi_t)\|^2 &= \mathbb{E} \|\nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_t, \Psi^*(\bar{\mathbf{w}}_t)) - \nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_t, \Psi_t)\|^2 \\ &\stackrel{(a)}{\leq} \frac{L_{12}^2}{n} \mathbb{E} \|\Psi^*(\bar{\mathbf{w}}_t) - \Psi_t\|_F^2 \\ &\stackrel{(b)}{\leq} \frac{2L_{12}^2}{\mu_2 n} \mathbb{E} [\Phi(\bar{\mathbf{w}}_t) - f(\bar{\mathbf{w}}_t, \Psi_t)] \\ &\stackrel{(c)}{=} \frac{2L_{12}^2}{\mu_2 n} b_t. \end{aligned} \quad (48)$$

In above and to derive (a), we employ the result of Lemma 1 which shows that given Assumption 3, the gradient function $\nabla_{\mathbf{w}} f(\mathbf{w}, \cdot)$ is L_{12}/\sqrt{n} Lipschitz. To derive (b), we use Assumption 4 (ii) and lastly, (c) is implied from the definition of b_t . The second term in RHS of (47) can be bounded by noting that the local gradients $\nabla_{\mathbf{w}} f^i(\cdot, \boldsymbol{\psi}^i)$ are L_1 -Lipschitz, which we can write

$$\begin{aligned} \mathbb{E} \left\| \nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_t, \Psi_t) - \frac{1}{n} \sum_{i \in [n]} \nabla_{\mathbf{w}} f^i(\mathbf{w}_t^i, \boldsymbol{\psi}_t^i) \right\|^2 &= \mathbb{E} \left\| \frac{1}{n} \sum_{i \in [n]} \nabla_{\mathbf{w}} f^i(\bar{\mathbf{w}}_t, \boldsymbol{\psi}_t^i) - \frac{1}{n} \sum_{i \in [n]} \nabla_{\mathbf{w}} f^i(\mathbf{w}_t^i, \boldsymbol{\psi}_t^i) \right\|^2 \\ &\leq \frac{L_1^2}{n} \sum_{i \in [n]} \mathbb{E} \|\mathbf{w}_t^i - \bar{\mathbf{w}}_t\|^2 \\ &= L_1^2 e_t. \end{aligned} \quad (49)$$

Finally, plugging (48) and (49) back in (47) implies the claim of the lemma, that is

$$h_t \leq \frac{4L_{12}^2}{\mu_2 n} b_t + 2L_1^2 e_t. \quad (50)$$

E.5 Proof of Lemma 5

We begin the proof by noting the definition of b_t and use the fact that the gradients $\nabla_\Psi f(\mathbf{w}, \cdot)$ are $\frac{L_2}{n}$ -Lipschitz (Refer to Lemma 1). We can accordingly write

$$\begin{aligned} \Phi(\bar{\mathbf{w}}_{t+1}) - f(\bar{\mathbf{w}}_{t+1}, \Psi_{t+1}) &\leq \Phi(\bar{\mathbf{w}}_{t+1}) - f(\bar{\mathbf{w}}_{t+1}, \Psi_t) - \langle \nabla_\Psi f(\bar{\mathbf{w}}_{t+1}, \Psi_t), \Psi_{t+1} - \Psi_t \rangle \\ &\quad + \frac{L_2}{2n} \|\Psi_{t+1} - \Psi_t\|_F^2. \end{aligned} \quad (51)$$

In this work, we define the inner product for any two matrices A, B as follows

$$\langle A, B \rangle := \text{Tr}(A^\top B). \quad (52)$$

Note that according to the ascent update rule of FedRobust in Algorithm 1, we can write

$$\Psi_{t+1} - \Psi_t = \eta_2 \tilde{\partial}_t f, \quad (53)$$

where we adopt the following short-hand notation for the stochastic gradients at iteration t with respect to the maximization variables $\psi_t^i = (\Lambda_t^i, \delta_t^i)$

$$\tilde{\partial}_t f = \begin{pmatrix} \tilde{\nabla}_\psi f^1(\mathbf{w}_t^1, \psi_t^1) \\ \vdots \\ \tilde{\nabla}_\psi f^n(\mathbf{w}_t^n, \psi_t^n) \end{pmatrix} = \begin{pmatrix} \tilde{\nabla}_\Lambda f^1(\mathbf{w}_t^1, \Lambda_t^1, \delta_t^1) & \tilde{\nabla}_\delta f^1(\mathbf{w}_t^1, \Lambda_t^1, \delta_t^1) \\ \vdots & \vdots \\ \tilde{\nabla}_\Lambda f^n(\mathbf{w}_t^n, \Lambda_t^n, \delta_t^n) & \tilde{\nabla}_\delta f^n(\mathbf{w}_t^n, \Lambda_t^n, \delta_t^n) \end{pmatrix}. \quad (54)$$

We also denote the gradients by $\partial_t f = \mathbb{E}[\tilde{\partial}_t f]$ where the expectation is with respect to the randomness in stochastic gradients $\tilde{\nabla}_\psi f^i$. According to Assumption 2, each of the local stochastic gradients $\tilde{\nabla}_\psi f^i(\mathbf{w}_t^i, \psi_t^i)$ are variance-bounded by σ_ψ^2 . Therefore, we can bound the variance of $\tilde{\partial}_t f$ as $\mathbb{E}\|\tilde{\partial}_t f - \partial_t f\|_F^2 \leq n\sigma_\psi^2$. Now, we can plug these back in (51) which implies

$$\begin{aligned} \Phi(\bar{\mathbf{w}}_{t+1}) - \mathbb{E}f(\bar{\mathbf{w}}_{t+1}, \Psi_{t+1}) &\leq \Phi(\bar{\mathbf{w}}_{t+1}) - f(\bar{\mathbf{w}}_{t+1}, \Psi_t) - \eta_2 \frac{n}{2} \|\nabla_\Psi f(\bar{\mathbf{w}}_{t+1}, \Psi_t)\|_F^2 + \eta_2^2 \frac{L_2}{2} \sigma_\psi^2 \\ &\quad + \eta_2 \frac{n}{2} \left\| \nabla_\Psi f(\bar{\mathbf{w}}_{t+1}, \Psi_t) - \frac{1}{n} \partial_t f \right\|_F^2 - \frac{\eta_2}{2n} (1 - \eta_2 L_2) \|\partial_t f\|_F^2, \end{aligned} \quad (55)$$

where the expectation is with respect to the randomness of the stochastic gradients $\tilde{\partial}_t f$ while conditioning on all the randomness history. Now recall from Assumption 4 (ii) that $-f(\bar{\mathbf{w}}_{t+1}, \cdot)$ is μ_2 -PL implying that $\|\nabla_\Psi f(\bar{\mathbf{w}}_{t+1}, \Psi_t)\|_F^2 \geq 2\mu_2(\Phi(\bar{\mathbf{w}}_{t+1}) - f(\bar{\mathbf{w}}_{t+1}, \Psi_t))$. Moreover, assume that $\eta_2 \leq 1/L_2$ to remove the last term in (55). Putting altogether implies that

$$\begin{aligned} \Phi(\bar{\mathbf{w}}_{t+1}) - \mathbb{E}f(\bar{\mathbf{w}}_{t+1}, \Psi_{t+1}) &\leq (1 - \mu_2 \eta_2 n) (\Phi(\bar{\mathbf{w}}_{t+1}) - f(\bar{\mathbf{w}}_{t+1}, \Psi_t)) + \eta_2^2 \frac{L_2}{2} \sigma_\psi^2 \\ &\quad + \eta_2 \frac{n}{2} \left\| \nabla_\Psi f(\bar{\mathbf{w}}_{t+1}, \Psi_t) - \frac{1}{n} \partial_t f \right\|_F^2. \end{aligned} \quad (56)$$

Next, we continue to bound the last term in RHS of (56). We can write

$$\begin{aligned} \left\| \nabla_\Psi f(\bar{\mathbf{w}}_{t+1}, \Psi_t) - \frac{1}{n} \partial_t f \right\|_F^2 &= \frac{1}{n^2} \sum_{i \in [n]} \left\| \nabla_\psi f^i(\bar{\mathbf{w}}_{t+1}, \psi_t^i) - \nabla_\psi f^i(\mathbf{w}_t^i, \psi_t^i) \right\|_F^2 \\ &\leq \frac{L_{21}^2}{n^2} \sum_{i \in [n]} \left\| \bar{\mathbf{w}}_{t+1} - \mathbf{w}_t^i \right\|^2 \\ &\leq \frac{2L_{21}^2}{n^2} \sum_{i \in [n]} \left\| \mathbf{w}_t^i - \bar{\mathbf{w}}_t \right\|^2 + \frac{2L_{21}^2}{n} \left\| \bar{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t \right\|^2, \end{aligned} \quad (57)$$

where the first inequality above uses Assumption 3 on Lipschitz continuity of local gradients and the second inequality simply uses the inequality $\|\mathbf{a} + \mathbf{b}\|^2 \leq 2\|\mathbf{a}\|^2 + 2\|\mathbf{b}\|^2$. Next, let us bound the term $\|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t\|^2$ in expectation as follows. Using the descent update rule in Algorithm 1 and considering Assumption 2 on variance of the stochastic gradients $\tilde{\nabla}_{\mathbf{w}} f^i$ we can write

$$\begin{aligned} \mathbb{E}\|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t\|^2 &= \eta_1^2 \mathbb{E} \left\| \frac{1}{n} \sum_{i \in [n]} \tilde{\nabla}_{\mathbf{w}} f^i(\mathbf{w}_t^i, \psi_t^i) \right\|^2 \\ &\leq \eta_1^2 \mathbb{E} \left\| \frac{1}{n} \sum_{i \in [n]} \nabla_{\mathbf{w}} f^i(\mathbf{w}_t^i, \psi_t^i) \right\|^2 + \eta_1^2 \frac{\sigma_{\mathbf{w}}^2}{n} \\ &= \eta_1^2 g_t + \eta_1^2 \frac{\sigma_{\mathbf{w}}^2}{n}, \end{aligned} \quad (58)$$

where we use the short-hand notation of g_t also listed in Table 1. Plugging (58) back in (57) and noting the notation $e_t = \frac{1}{n} \sum_{i \in [n]} \mathbb{E}\|\mathbf{w}_t^i - \bar{\mathbf{w}}_t\|^2$ implies that

$$\mathbb{E} \left\| \nabla_{\Psi} f(\bar{\mathbf{w}}_{t+1}, \Psi_t) - \frac{1}{n} \partial_t f \right\|_F^2 \leq \frac{2L_{21}^2}{n} e_t + \eta_1^2 \frac{2L_{21}^2}{n} g_t + \eta_1^2 \frac{2L_{21}^2}{n} \frac{\sigma_{\mathbf{w}}^2}{n}. \quad (59)$$

Before proceeding to bound more terms, let us recall what we have shown till this point. We plug (59) back in (56), take the expectation with respect to all the sources of randomness and use the notation $b_t = \mathbb{E}[\Phi(\bar{\mathbf{w}}_t) - f(\bar{\mathbf{w}}_t, \Psi_t)]$ to conclude

$$\begin{aligned} b_{t+1} &\leq (1 - \mu_2 \eta_2 n) \mathbb{E} [\Phi(\bar{\mathbf{w}}_{t+1}) - f(\bar{\mathbf{w}}_{t+1}, \Psi_t)] \\ &\quad + \eta_2 L_{21}^2 e_t + \eta_1^2 \eta_2 L_{21}^2 g_t + \eta_1^2 \eta_2 L_{21}^2 \frac{\sigma_{\mathbf{w}}^2}{n} + \eta_2^2 \frac{L_2}{2} \sigma_{\Psi}^2. \end{aligned} \quad (60)$$

To bound the term $\mathbb{E} [\Phi(\bar{\mathbf{w}}_{t+1}) - f(\bar{\mathbf{w}}_{t+1}, \Psi_t)]$, we can decompose it to the following three terms:

$$\Phi(\bar{\mathbf{w}}_{t+1}) - f(\bar{\mathbf{w}}_{t+1}, \Psi_t) = \Phi(\bar{\mathbf{w}}_t) - f(\bar{\mathbf{w}}_t, \Psi_t) + f(\bar{\mathbf{w}}_t, \Psi_t) - f(\bar{\mathbf{w}}_{t+1}, \Psi_t) + \Phi(\bar{\mathbf{w}}_{t+1}) - \Phi(\bar{\mathbf{w}}_t). \quad (61)$$

Given the Lipschitz gradient assumption for the local functions in Assumption 3 and using Lemma 1 on Lipschitz gradient for the global function, we can write

$$f(\bar{\mathbf{w}}_t, \Psi_t) - f(\bar{\mathbf{w}}_{t+1}, \Psi_t) \leq -\langle \nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_t, \Psi_t), \bar{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t \rangle + \frac{L_1}{2} \|\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t\|^2, \quad (62)$$

where $\bar{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t = -\eta_1 \frac{1}{n} \sum_{i \in [n]} \tilde{\nabla}_{\mathbf{w}} f^i(\mathbf{w}_t^i, \psi_t^i)$. Taking expectation from both sides of (62) implies that

$$\begin{aligned} \mathbb{E} [f(\bar{\mathbf{w}}_t, \Psi_t) - f(\bar{\mathbf{w}}_{t+1}, \Psi_t)] &\stackrel{(a)}{\leq} \eta_1 \mathbb{E} \|\nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_t, \Psi_t) - \nabla \Phi(\bar{\mathbf{w}}_t)\|^2 + \eta_1 \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 \\ &\quad + \left(\frac{\eta_1}{2} + \eta_1^2 \frac{L_1}{2} \right) g_t + \eta_1^2 \frac{L_1}{2} \frac{\sigma_{\mathbf{w}}^2}{n} \\ &\stackrel{(b)}{\leq} \eta_1 \frac{2L_{12}^2}{\mu_2 n} b_t + \eta_1 \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 + \left(\frac{\eta_1}{2} + \eta_1^2 \frac{L_1}{2} \right) g_t + \eta_1^2 \frac{L_1}{2} \frac{\sigma_{\mathbf{w}}^2}{n}, \end{aligned} \quad (63)$$

where in inequality (a) we use the inequality $2\langle \mathbf{a}, \mathbf{b} \rangle \leq \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2$ and also the result in (58). To derive (b), we use Assumptions 3 and 4 (ii), result of Lemma 1 and the notation $b_t = \mathbb{E}[\Phi(\bar{\mathbf{w}}_t) - f(\bar{\mathbf{w}}_t, \Psi_t)]$ to write

$$\begin{aligned} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t) - \nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_t, \Psi_t)\|^2 &= \mathbb{E} \|\nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_t, \Psi^*(\bar{\mathbf{w}}_t)) - \nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_t, \Psi_t)\|^2 \\ &\leq \frac{L_{12}^2}{n} \mathbb{E} \|\Psi^*(\bar{\mathbf{w}}_t) - \Psi_t\|_F^2 \\ &\leq \frac{2L_{12}^2}{\mu_2 n} \mathbb{E} [\Phi(\bar{\mathbf{w}}_t) - f(\bar{\mathbf{w}}_t, \Psi_t)] \\ &= \frac{2L_{12}^2}{\mu_2 n} b_t. \end{aligned} \quad (64)$$

We now have all the ingredients to conclude the claim of Lemma 5. To do so, we combine the result of Lemma 3 which bounds the term $\mathbb{E}[\Phi(\bar{\mathbf{w}}_{t+1})] - \mathbb{E}[\Phi(\bar{\mathbf{w}}_t)]$, Lemma 4 that shows $h_t \leq 4L_{12}^2 b_t / (\mu_2 n) + 2L_1^2 e_t$, and the bound (63); plug back in (61) and then in (60) and conclude the claim of the lemma, that is

$$\begin{aligned} b_{t+1} &\leq (1 - \mu_2 \eta_2 n) \left(1 + \eta_1 \frac{4L_{12}^2}{\mu_2 n} \right) b_t + \frac{\eta_1}{2} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 + \frac{\eta_1^2}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) g_t \\ &\quad + (\eta_1 L_1^2 + \eta_2 L_{21}^2) e_t + \frac{\eta_1^2}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \frac{\sigma_{\mathbf{w}}^2}{n} + \frac{\eta_2^2}{2} L_2 \sigma_\psi^2, \end{aligned} \quad (65)$$

E.6 Proof of Lemma 6

To prove this lemma, we first need to establish an intermediate step, which is stated in the following.

Proposition 1. *If Assumptions 1, 2 and 3 hold, then*

$$e_t \leq 16\eta_1^2(\tau - 1)L_1^2 \sum_{l=t_c+1}^{t-1} e_l + 10\eta_1^2(\tau - 1) \sum_{l=t_c+1}^{t-1} g_l + 8\eta_1^2(\tau - 1)^2 \rho^2 + 4\eta_1^2(\tau - 1)(n + 1) \frac{\sigma_{\mathbf{w}}^2}{n}. \quad (66)$$

Proof of Proposition 1. Consider an iteration $t \geq 1$ and let t_c denote the index of the most recent communication between the workers and the server, i.e. $t_c = \lfloor \frac{t}{\tau} \rfloor \tau$. Therefore, all the workers share the same local minimization model at iteration $t_c + 1$, i.e. $\mathbf{w}_{t_c+1}^1 = \dots = \mathbf{w}_{t_c+1}^n = \bar{\mathbf{w}}_{t_c+1}$. According to the update rule of FedRobust, we can write for each node i that

$$\begin{aligned} \mathbf{w}_{t_c+2}^i &= \mathbf{w}_{t_c+1}^i - \eta_1 \tilde{\nabla}_{\mathbf{w}} f^i(\mathbf{w}_{t_c+1}^i, \psi_{t_c+1}^i), \\ &\vdots \\ \mathbf{w}_t^i &= \mathbf{w}_{t-1}^i - \eta_1 \tilde{\nabla}_{\mathbf{w}} f^i(\mathbf{w}_{t-1}^i, \psi_{t-1}^i). \end{aligned} \quad (67)$$

Summing up all the equalities in (67) yields that

$$\mathbf{w}_t^i = \mathbf{w}_{t_c+1}^i - \eta_1 \sum_{l=t_c+1}^{t-1} \tilde{\nabla}_{\mathbf{w}} f^i(\mathbf{w}_l^i, \psi_l^i). \quad (68)$$

Therefore, the difference of the local models \mathbf{w}_t^i and their average $\bar{\mathbf{w}}_t$ can be written as

$$\begin{aligned} \mathbf{w}_t^i - \bar{\mathbf{w}}_t &= \mathbf{w}_{t_c+1}^i - \eta_1 \sum_{l=t_c+1}^{t-1} \tilde{\nabla}_{\mathbf{w}} f^i(\mathbf{w}_l^i, \psi_l^i) - \left(\bar{\mathbf{w}}_{t_c+1} - \eta_1 \frac{1}{n} \sum_{j \in [n]} \sum_{l=t_c+1}^{t-1} \tilde{\nabla}_{\mathbf{w}} f^j(\mathbf{w}_l^j, \psi_l^j) \right) \\ &= -\eta_1 \left(\sum_{l=t_c+1}^{t-1} \tilde{\nabla}_{\mathbf{w}} f^i(\mathbf{w}_l^i, \psi_l^i) - \frac{1}{n} \sum_{j \in [n]} \sum_{l=t_c+1}^{t-1} \tilde{\nabla}_{\mathbf{w}} f^j(\mathbf{w}_l^j, \psi_l^j) \right). \end{aligned} \quad (69)$$

This yields the following bound on each local deviation from the average $\mathbb{E} \|\mathbf{w}_t^i - \bar{\mathbf{w}}_t\|^2$:

$$\begin{aligned} \mathbb{E} \|\mathbf{w}_t^i - \bar{\mathbf{w}}_t\|^2 &= \eta_1^2 \mathbb{E} \left\| \sum_{l=t_c+1}^{t-1} \tilde{\nabla}_{\mathbf{w}} f^i(\mathbf{w}_l^i, \psi_l^i) - \frac{1}{n} \sum_{j \in [n]} \sum_{l=t_c+1}^{t-1} \tilde{\nabla}_{\mathbf{w}} f^j(\mathbf{w}_l^j, \psi_l^j) \right\|^2 \\ &\leq 2\eta_1^2 \mathbb{E} \left\| \sum_{l=t_c+1}^{t-1} \tilde{\nabla}_{\mathbf{w}} f^i(\mathbf{w}_l^i, \psi_l^i) \right\|^2 + 2\eta_1^2 \mathbb{E} \left\| \frac{1}{n} \sum_{j \in [n]} \sum_{l=t_c+1}^{t-1} \tilde{\nabla}_{\mathbf{w}} f^j(\mathbf{w}_l^j, \psi_l^j) \right\|^2 \\ &\stackrel{(a)}{\leq} \underbrace{2\eta_1^2 \mathbb{E} \left\| \sum_{l=t_c+1}^{t-1} \nabla_{\mathbf{w}} f^i(\mathbf{w}_l^i, \psi_l^i) \right\|^2}_{T_3} + \underbrace{2\eta_1^2 \mathbb{E} \left\| \frac{1}{n} \sum_{j \in [n]} \sum_{l=t_c+1}^{t-1} \nabla_{\mathbf{w}} f^j(\mathbf{w}_l^j, \psi_l^j) \right\|^2}_{T_4} \\ &\quad + 2\eta_1^2(t - t_c - 1)(n + 1) \frac{\sigma_{\mathbf{w}}^2}{n}, \end{aligned} \quad (70)$$

where we used Assumption 2 to bound the variance of the stochastic gradients and derive (a). The term T_4 in (70) can simply be bounded as

$$T_4 \leq \mathbb{E} \left\| \frac{1}{n} \sum_{j \in [n]} \sum_{l=t_c+1}^{t-1} \nabla_{\mathbf{w}} f^j(\mathbf{w}_l^j, \boldsymbol{\psi}_l^j) \right\|^2 \leq (t - t_c - 1) \sum_{l=t_c+1}^{t-1} \mathbb{E} \left\| \frac{1}{n} \sum_{j \in [n]} \nabla_{\mathbf{w}} f^j(\mathbf{w}_l^j, \boldsymbol{\psi}_l^j) \right\|^2 \quad (71)$$

Note that t_c denotes the latest server-worker communication before iteration t , hence $t - t_c \leq \tau$ where τ is the duration of local updates in each round. Therefore, we have

$$T_4 \leq (\tau - 1) \sum_{l=t_c+1}^{t-1} \mathbb{E} \left\| \frac{1}{n} \sum_{j \in [n]} \nabla_{\mathbf{w}} f^j(\mathbf{w}_l^j, \boldsymbol{\psi}_l^j) \right\|^2 \leq (\tau - 1) \sum_{l=t_c+1}^{t-1} g_l \quad (72)$$

Now we proceed to bound the term T_3 in (70) as follows:

$$\begin{aligned} T_3 &= \mathbb{E} \left\| \sum_{l=t_c+1}^{t-1} \nabla_{\mathbf{w}} f^i(\mathbf{w}_l^i, \boldsymbol{\psi}_l^i) \right\|^2 \\ &\leq (\tau - 1) \sum_{l=t_c+1}^{t-1} \mathbb{E} \left\| \nabla_{\mathbf{w}} f^i(\mathbf{w}_l^i, \boldsymbol{\psi}_l^i) \right\|^2 \\ &\leq 4(\tau - 1) \sum_{l=t_c+1}^{t-1} \mathbb{E} \left\| \nabla_{\mathbf{w}} f^i(\mathbf{w}_l^i, \boldsymbol{\psi}_l^i) - \nabla_{\mathbf{w}} f^i(\bar{\mathbf{w}}_l, \boldsymbol{\psi}_l^i) \right\|^2 \\ &\quad + 4(\tau - 1) \sum_{l=t_c+1}^{t-1} \mathbb{E} \left\| \nabla_{\mathbf{w}} f^i(\bar{\mathbf{w}}_l, \boldsymbol{\psi}_l^i) - \frac{1}{n} \sum_{j \in [n]} \nabla_{\mathbf{w}} f^j(\bar{\mathbf{w}}_l, \boldsymbol{\psi}_l^j) \right\|^2 \\ &\quad + 4(\tau - 1) \sum_{l=t_c+1}^{t-1} \mathbb{E} \left\| \frac{1}{n} \sum_{j \in [n]} \nabla_{\mathbf{w}} f^j(\bar{\mathbf{w}}_l, \boldsymbol{\psi}_l^j) - \frac{1}{n} \sum_{j \in [n]} \nabla_{\mathbf{w}} f^j(\mathbf{w}_l^j, \boldsymbol{\psi}_l^j) \right\|^2 \\ &\quad + 4(\tau - 1) \sum_{l=t_c+1}^{t-1} \mathbb{E} \left\| \frac{1}{n} \sum_{j \in [n]} \nabla_{\mathbf{w}} f^j(\mathbf{w}_l^j, \boldsymbol{\psi}_l^j) \right\|^2 \end{aligned} \quad (73)$$

We can simply this bound by using Assumption 3 on Lipschitz gradients for the local objectives f^i s and applying the notations for e_l and g_l to derive

$$\begin{aligned} T_3 &\leq 4(\tau - 1) L_1^2 \sum_{l=t_c+1}^{t-1} \mathbb{E} \left\| \mathbf{w}_l^i - \bar{\mathbf{w}}_l \right\|^2 + 4(\tau - 1) \sum_{l=t_c+1}^{t-1} \mathbb{E} \left\| \nabla_{\mathbf{w}} f^i(\bar{\mathbf{w}}_l, \boldsymbol{\psi}_l^i) - \nabla_{\mathbf{w}} f(\bar{\mathbf{w}}_l, \boldsymbol{\Psi}_l) \right\|^2 \\ &\quad + 4(\tau - 1) L_1^2 \sum_{l=t_c+1}^{t-1} e_l + 4(\tau - 1) \sum_{l=t_c+1}^{t-1} g_l \end{aligned} \quad (74)$$

We can plug (72) and (74) into (70) and take the average of the both sides over $i = 1, \dots, n$. This implies that

$$e_t \leq 16\eta_1^2(\tau - 1) L_1^2 \sum_{l=t_c+1}^{t-1} e_l + 10\eta_1^2(\tau - 1) \sum_{l=t_c+1}^{t-1} g_l + 8\eta_1^2(\tau - 1)^2 \rho^2 + 4\eta_1^2(\tau - 1)(n + 1) \frac{\sigma_{\mathbf{w}}^2}{n}. \quad (75)$$

In above, we used the result of Proposition 2 that given Assumption 1, bounds the gradient diversity $\frac{1}{n} \sum_{i \in [n]} \left\| \nabla_{\mathbf{w}} f^i(\mathbf{w}, \boldsymbol{\psi}^i) - \nabla_{\mathbf{w}} f(\mathbf{w}, \boldsymbol{\Psi}) \right\|^2 \leq \rho^2$, where $\rho^2 = 3\rho_f^2 + 6L_{12}^2(\epsilon_1^2 + \epsilon_2^2)$. We defer the proof of this proposition to the end of this section. This concludes the proof of Proposition 1. \square

Having set the required intermediate steps, we resume the proof of Lemma 6. According to Proposition 1, we can write the term e_t as follows

$$e_t \leq C_1 \sum_{l=t_c+1}^{t-1} e_l + C_2 \sum_{l=t_c+1}^{t-1} g_l + C_3 \quad (76)$$

where we use the following short-hand coefficients

$$\begin{aligned} C_1 &:= 16\eta_1^2(\tau-1)L_1^2 \\ C_2 &:= 10\eta_1^2(\tau-1) \\ C_3 &:= 8\eta_1^2(\tau-1)^2\rho^2 + 4\eta_1^2(\tau-1)(n+1)\frac{\sigma_w^2}{n}. \end{aligned} \quad (77)$$

We can then write this bound for every iteration in $[t_c + 1 : t]$, that is

$$\begin{aligned} e_{t_c+1} &= 0 \\ e_{t_c+2} &\leq C_1 e_{t_c+1} + C_2 g_{t_c+1} + C_3 \\ &\vdots \\ e_t &\leq C_1 (e_{t_c+1} + \dots + e_{t-1}) + C_2 (g_{t_c+1} + \dots + g_{t-1}) + C_3. \end{aligned} \quad (78)$$

Summing all of the inequalities results in the following

$$\sum_{l=t_c+1}^{t-1} e_l \leq C_1(\tau-1) \sum_{l=t_c+1}^{t-1} e_l + C_2(\tau-1) \sum_{l=t_c+1}^{t-1} g_l + C_3(\tau-1). \quad (79)$$

We can further rearrange the terms above and write

$$\sum_{l=t_c+1}^{t-1} e_l \leq \frac{C_2(\tau-1)}{1-C_1(\tau-1)} \sum_{l=t_c+1}^{t-1} g_l + \frac{C_3(\tau-1)}{1-C_1(\tau-1)}. \quad (80)$$

Now, if we assume that $C_1(\tau-1) \leq 1/2$, then we get the following bound on $\sum_{l=t_c+1}^{t-1} e_l$

$$\sum_{l=t_c+1}^{t-1} e_l \leq 2C_2(\tau-1) \sum_{l=t_c+1}^{t-1} g_l + 2C_3(\tau-1) \quad (81)$$

Plugging back in (99) and using the assumption $C_1(\tau-1) \leq 1/2$ yields that

$$\begin{aligned} e_t &\leq C_1 \left(2C_2(\tau-1) \sum_{l=t_c+1}^{t-1} g_l + 2C_3(\tau-1) \right) + C_2 \sum_{l=t_c+1}^{t-1} g_l + C_3 \\ &\leq 2C_2 \sum_{l=t_c+1}^{t-1} g_l + 2C_3, \end{aligned} \quad (82)$$

which concludes the proof of Lemma 6. Lastly, we present the following proposition along with its proof which we used this result to prove Proposition 1.

Proposition 2. *An immediate implication of Assumptions 1 and 3 is that for any w, Ψ , the diversity of the local gradients is bounded in the following sense*

$$\frac{1}{n} \sum_{i \in [n]} \left\| \nabla_w f^i(w, \psi^i) - \nabla_w f(w, \Psi) \right\|^2 \leq \rho^2, \quad (83)$$

where we denote $\rho^2 = 3\rho_f^2 + 6L_{12}^2(\epsilon_1^2 + \epsilon_2^2)$.

Proof of Proposition 2. The proof is simply implied from Assumptions 1 and 3 by writing

$$\begin{aligned} \frac{1}{n} \sum_{i \in [n]} \left\| \nabla_w f^i(w, \psi^i) - \nabla_w f(w, \Psi) \right\|^2 &\leq 3 \frac{1}{n} \sum_{i \in [n]} \left\| \nabla_w f^i(w, \Lambda^i, \delta^i) - \nabla_w f^i(w, I, 0) \right\|^2 \\ &\quad + 3 \frac{1}{n} \sum_{i \in [n]} \left\| \nabla_w f^i(w) - \nabla_w f(w) \right\|^2 \\ &\quad + 3 \frac{1}{n} \sum_{i \in [n]} \left\| \nabla_w f(w, I, 0) - \nabla_w f(w, \Psi) \right\|^2 \\ &\leq 3\rho_f^2 + 6L_{12}^2(\epsilon_1^2 + \epsilon_2^2). \end{aligned} \quad (84)$$

□

E.7 Proof of Lemma 7

[16] proves a similar claim for $\Gamma = 0$. For completeness, we provide the proof for general case when $\Gamma \neq 0$. Let t_c denote the index of the most recent communication round, i.e. $t_c = \lfloor \frac{t}{\tau} \rfloor \tau$. We can write $t = t_c + r$ where $1 \leq r \leq \tau$. Starting from $r = 1$, we can write

$$\begin{aligned} P_{t_c+2} &\leq \Upsilon P_{t_c+1} - \frac{\eta_1}{2} (1 - \eta_1 L) g_{t_c+1} + \Gamma \\ &\leq \Upsilon P_{t_c+1} + \Gamma, \end{aligned} \quad (85)$$

where the last inequality holds if

$$\eta_1 L \leq 1. \quad (86)$$

We can continue for $r = 2$ as follows

$$\begin{aligned} P_{t_c+3} &\leq \Upsilon P_{t_c+2} - \frac{\eta_1}{2} (1 - \eta_1 L) g_{t_c+2} + \eta_1^2 B g_{t_c+1} + \Gamma \\ &\stackrel{(a)}{\leq} \Upsilon^2 P_{t_c+1} - \frac{\eta_1}{2} \Upsilon \left(1 - \eta_1 L - \eta_1 \frac{2B}{\Upsilon} \right) g_{t_c+1} + \Gamma(1 + \Upsilon) \\ &\stackrel{(b)}{\leq} \Upsilon^2 P_{t_c+1} + \Gamma(1 + \Upsilon) \end{aligned} \quad (87)$$

where (a) is due to the inequality $P_{t_c+2} \leq \Upsilon P_{t_c+1} - \frac{\eta_1}{2} (1 - \eta_1 L) g_{t_c+1} + \Gamma$ and (b) holds if

$$1 - \eta_1 L - \eta_1 \frac{2B}{\Upsilon} \geq 0, \quad (88)$$

or equivalently

$$\eta_1 \left(L + \frac{2B}{\Upsilon} \right) \leq 1. \quad (89)$$

We can continue the same argument up to $r + 1$ and write

$$P_{t_c+r+1} \leq \Upsilon^r P_{t_c+1} + \Gamma(1 + \Upsilon + \dots + \Upsilon^{r-1}), \quad (90)$$

if the step-size is as small as follows

$$\eta_1 \left(L + \frac{2B}{\Upsilon^{r-1}} (1 + \Upsilon + \dots + \Upsilon^{r-2}) \right) \leq 1. \quad (91)$$

Since $1 + \Upsilon + \dots + \Upsilon^{r-2} \leq \frac{1}{1-\Upsilon}$, then the following condition implies all the previous ones on η

$$\eta_1 \left(L + \frac{2B}{\Upsilon^{r-1}(1-\Upsilon)} \right). \quad (92)$$

Moreover, since $\Upsilon < 1$, then the strongest condition on η is (92) when we put the largest possible value for r which is τ , yielding

$$\eta_1 \left(L + \frac{2B}{\Upsilon^{\tau-1}(1-\Upsilon)} \right). \quad (93)$$

Lastly, we note that $1 + \Upsilon + \dots + \Upsilon^{r-1} \leq \frac{1}{1-\Upsilon}$ in (90), and the claim is concluded.

E.8 Proof of Lemma 8

Recall the result of Lemma 5 in which we showed that if $\eta_2 \leq 1/L_2$, then the following contraction bound on the sequence $\{b_t\}_{t \geq 0}$ holds:

$$\begin{aligned} b_{t+1} &\leq (1 - \mu_2 \eta_2 n) \left(1 + \eta_1 \frac{4L_{12}^2}{\mu_2 n} \right) b_t + \frac{\eta_1}{2} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 + \frac{\eta_1^2}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) g_t \\ &\quad + (\eta_1 L_1^2 + \eta_2 L_{21}^2) e_t + \frac{\eta_1^2}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \frac{\sigma_{\mathbf{w}}^2}{n} + \frac{\eta_2^2}{2} L_2 \sigma_\psi^2, \end{aligned} \quad (94)$$

and consider the coefficient of b_t in above. A simple calculation yields that if the step-sizes satisfy the condition $\frac{\eta_2}{\eta_1} \geq \frac{8L_{12}^2}{\mu_2^2 n^2}$, then we have

$$(1 - \mu_2 \eta_2 n) \left(1 + \eta_1 \frac{4L_{12}^2}{\mu_2 n} \right) \leq 1 - \frac{1}{2} \mu_2 \eta_2 n. \quad (95)$$

Now, we denote $\gamma = 1 - \frac{1}{2} \mu_2 \eta_2 n$ and apply (94) to all iterations $t = 0, \dots, T-1$, which yields that

$$\begin{aligned} b_0 &\leq \frac{2L_2^2}{\mu_2 n} (\epsilon_1^2 + \epsilon_2^2), \\ b_1 &\leq \gamma b_0 + \frac{\eta_1}{2} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 + \frac{\eta_1^2}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) g_0 + (\eta_1 L_1^2 + \eta_2 L_{21}^2) e_0 \\ &\quad + \frac{\eta_1^2}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \frac{\sigma_{\mathbf{w}}^2}{n} + \frac{\eta_2^2}{2} L_2 \sigma_\psi^2, \\ &\quad \vdots \\ b_{T-1} &\leq \gamma b_{T-2} + \frac{\eta_1}{2} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 + \frac{\eta_1^2}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) g_{T-2} + (\eta_1 L_1^2 + \eta_2 L_{21}^2) e_{T-2} \\ &\quad + \frac{\eta_1^2}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \frac{\sigma_{\mathbf{w}}^2}{n} + \frac{\eta_2^2}{2} L_2 \sigma_\psi^2. \end{aligned} \quad (96)$$

Taking the average of the T inequalities above yields that

$$\begin{aligned} (1 - \gamma) \frac{1}{T} \sum_{t=0}^{T-1} b_t &\leq \frac{2L_2^2}{\mu_2 n} \frac{\epsilon_1^2 + \epsilon_2^2}{T} + \frac{\eta_1}{2} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 \\ &\quad + \frac{\eta_1^2}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \frac{1}{T} \sum_{t=0}^{T-1} g_t + (\eta_1 L_1^2 + \eta_2 L_{21}^2) \frac{1}{T} \sum_{t=0}^{T-1} e_t \\ &\quad + \frac{\eta_1^2}{2} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \frac{\sigma_{\mathbf{w}}^2}{n} + \frac{\eta_2^2}{2} L_2 \sigma_\psi^2. \end{aligned} \quad (97)$$

We can further divide both sides of (97) by $1 - \gamma$ and conclude

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} b_t &\leq \frac{4L_2^2}{\mu_2^2 n^2} \frac{\epsilon_1^2 + \epsilon_2^2}{\eta_2 T} + \frac{\eta_1}{\eta_2} \frac{1}{\mu_2 n} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla \Phi(\bar{\mathbf{w}}_t)\|^2 \\ &\quad + \frac{\eta_1^2}{\eta_2} \frac{1}{\mu_2 n} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \frac{1}{T} \sum_{t=0}^{T-1} g_t + \frac{1}{\eta_2} \frac{2}{\mu_2 n} (\eta_1 L_1^2 + \eta_2 L_{21}^2) \frac{1}{T} \sum_{t=0}^{T-1} e_t \\ &\quad + \frac{\eta_1^2}{\eta_2} \frac{1}{\mu_2 n} (L_1 + L_\Phi + 2\eta_2 L_{21}^2) \frac{\sigma_{\mathbf{w}}^2}{n} + \eta_2 \frac{L_2}{\mu_2 n} \sigma_\psi^2. \end{aligned} \quad (98)$$

E.9 Proof of Lemma 9

We begin by noting the result of Proposition 1 in which we showed the following bound on e_t

$$e_t \leq C_1 \sum_{l=t_c+1}^{t-1} e_l + C_2 \sum_{l=t_c+1}^{t-1} g_l + C_3, \quad (99)$$

where we defined the coefficients C_1, C_2, C_3 in (77) and recall here for more convenient:

$$\begin{aligned} C_1 &:= 16\eta_1^2(\tau - 1)L_1^2 \\ C_2 &:= 10\eta_1^2(\tau - 1) \\ C_3 &:= 8\eta_1^2(\tau - 1)^2 \rho^2 + 4\eta_1^2(\tau - 1)(n + 1) \frac{\sigma_{\mathbf{w}}^2}{n}. \end{aligned} \quad (100)$$

Next, we apply this bound to each iteration $t = 0, \dots, T - 1$ as follows

$$\begin{aligned}
e_0 &= 0 \\
\begin{cases} e_1 &= 0 \\ e_2 &\leq C_1 e_1 + C_2 g_1 + C_3 \\ \vdots & \\ e_\tau &\leq C_1 (e_1 + \dots + e_{\tau-1}) + C_2 (g_1 + \dots + g_{\tau-1}) + C_3 \end{cases} \\
\begin{cases} e_{\tau+1} &= 0 \\ e_{\tau+2} &\leq C_1 e_{\tau+1} + C_2 g_{\tau+1} + C_3 \\ \vdots & \\ e_{2\tau} &\leq C_1 (e_{\tau+1} + \dots + e_{2\tau-1}) + C_2 (g_{\tau+1} + \dots + g_{2\tau-1}) + C_3 \\ \vdots & \end{cases} \\
\begin{cases} e_{T_c+1} &= 0 \\ e_{T_c+2} &\leq C_1 e_{T_c+1} + C_2 g_{T_c+1} + C_3 \\ \vdots & \\ e_{T-1} &\leq C_1 (e_{T_c+1} + \dots + e_{T-2}) + C_2 (g_{T_c+1} + \dots + g_{T-2}) + C_3, \end{cases}
\end{aligned} \tag{101}$$

where $T_c = \lfloor \frac{T}{\tau} \rfloor \tau$ denote the index of the most recent communication between the workers and the server before iteration T . Summing the above inequalities yields that

$$\sum_{t=0}^{T-1} e_t \leq C_1 (\tau - 1) \sum_{t=0}^{T-1} e_t + C_2 (\tau - 1) \sum_{t=0}^{T-1} g_t + C_3 T. \tag{102}$$

Now if we assume that $C_1 (\tau - 1) = 16\eta_1^2 (\tau - 1)^2 L_1^2 \leq \frac{1}{2}$, the the claim is concluded by rearranging the terms in (102):

$$\frac{1}{T} \sum_{t=0}^{T-1} e_t \leq 2C_2 (\tau - 1) \frac{1}{T} \sum_{t=0}^{T-1} g_t + 2C_3. \tag{103}$$

Appendix F Proof of Theorem 3

Fix a distribution \tilde{P} and consider

$$\max_{\Lambda, \delta} \mathbb{E}_{\tilde{P}} [\ell(f_{\mathbf{w}}(\Lambda \mathbf{x} + \delta))] - \lambda \|\delta\|_2^2 - \lambda \|\Lambda - I\|_F^2 \tag{104}$$

Assuming a 1-Lipschitz loss ℓ with 1-Lipschitz gradient, based on [36]’s Lemma 7 the above function’s gradient with respect to δ has a Lipschitz constant bounded by

$$\text{Lip}(\nabla f_{\mathbf{w}}) := \left(\prod_{i=1}^L \|\mathbf{w}_i\|_{\sigma} \right) \sum_{i=1}^l \prod_{j=1}^i \|\mathbf{w}_j\|_{\sigma}.$$

Similarly, the expected loss’s derivative with respect to Λ will also be Lipschitz in the spectral norm with a Lipschitz constant upper-bounded by

$$B \text{Lip}(\nabla f_{\mathbf{w}}) = B \left(\prod_{i=1}^L \|\mathbf{w}_i\|_{\sigma} \right) \sum_{i=1}^l \prod_{j=1}^i \|\mathbf{w}_j\|_{\sigma}.$$

Given weights in \mathbf{w} , we denote the optimal solution for δ and Λ by $\delta_{\mathbf{w}}$ and $\Lambda_{\mathbf{w}}$, respectively. To apply the Pac-Bayes generalization analysis, we need to bound the change in $\delta_{\mathbf{w}}, \Lambda_{\mathbf{w}}$ caused by perturbing \mathbf{w} to $\mathbf{w} + \mathbf{u}$. Note that since $\lambda > (1 + B) \text{Lip}(\nabla f_{\mathbf{w}})$, the maximization problem for optimizing $\Lambda_{\mathbf{w}}, \delta_{\mathbf{w}}$ is maximizing a strongly-concave objective whose solutions will satisfy:

$$\begin{aligned}
\delta_{\mathbf{w}} &= \frac{1}{\lambda} \mathbb{E} [\nabla \ell \circ f_{\mathbf{w}}(\Lambda_{\mathbf{w}} \mathbf{x} + \delta_{\mathbf{w}})], \\
\Lambda_{\mathbf{w}} - I &= \frac{1}{\lambda} \mathbb{E} [(\nabla \ell \circ f_{\mathbf{w}}(\Lambda_{\mathbf{w}} \mathbf{x} + \delta_{\mathbf{w}})) \mathbf{X}^{\top}]
\end{aligned}$$

which are norm-bounded by $\frac{\text{Lip}(\ell \circ f_w)}{\lambda} \leq \frac{\prod_{i=1}^d \|w_i\|_\sigma}{\lambda}$ and $B \frac{\text{Lip}(\ell \circ f_w)}{\lambda} \leq B \frac{\prod_{i=1}^d \|w_i\|_\sigma}{\lambda}$, respectively. Therefore, for a norm-bounded perturbation u where $\|u_i\|_\sigma \leq \frac{1}{L} \|w_i\|_\sigma$ we can write

$$\begin{aligned}
& \|\delta_{w+u} - \delta_w\|_2 + \|\Lambda_{w+u} - \Lambda_w\|_\sigma \\
&= \left\| \frac{1}{\lambda} \mathbb{E}[\nabla \ell(f_{w+u}(\Lambda_{w+u} \mathbf{X} + \delta_{w+u}))] - \frac{1}{\lambda} \mathbb{E}[\nabla \ell(f_w(\Lambda_w \mathbf{X} + \delta_w))] \right\|_2 \\
&\quad + \left\| \frac{1}{\lambda} \mathbb{E}[\nabla \ell(f_{w+u}(\Lambda_{w+u} \mathbf{X} + \delta_{w+u})) \mathbf{X}^\top] - \frac{1}{\lambda} \mathbb{E}[\nabla \ell(f_w(\Lambda_w \mathbf{X} + \delta_w)) \mathbf{X}^\top] \right\|_\sigma \\
&= \left\| \frac{1}{\lambda} \mathbb{E}[\nabla \ell(f_{w+u}(\Lambda_{w+u} \mathbf{X} + \delta_{w+u})) - \nabla \ell(f_w(\Lambda_w \mathbf{X} + \delta_w))] \right\|_2 \\
&\quad + \left\| \frac{1}{\lambda} \mathbb{E}[(\nabla \ell(f_{w+u}(\Lambda_{w+u} \mathbf{X} + \delta_{w+u})) - \nabla \ell(f_w(\Lambda_w \mathbf{X} + \delta_w))) \mathbf{X}^\top] \right\|_\sigma \\
&\leq \left\| \frac{1}{\lambda} \mathbb{E}[\nabla \ell(f_{w+u}(\Lambda_{w+u} \mathbf{X} + \delta_{w+u})) - \nabla \ell(f_w(\Lambda_{w+u} \mathbf{X} + \delta_{w+u}))] \right\|_2 \\
&\quad + \left\| \frac{1}{\lambda} \mathbb{E}[\nabla \ell(f_w(\Lambda_{w+u} \mathbf{X} + \delta_{w+u})) - \nabla \ell(f_w(\Lambda_w \mathbf{X} + \delta_{w+u}))] \right\|_2 \\
&\quad + \left\| \frac{1}{\lambda} \mathbb{E}[\nabla \ell(f_w(\Lambda_w \mathbf{X} + \delta_{w+u})) - \nabla \ell(f_w(\Lambda_w \mathbf{X} + \delta_w))] \right\|_2 \\
&\quad + \left\| \frac{1}{\lambda} \mathbb{E}[(\nabla \ell(f_{w+u}(\Lambda_{w+u} \mathbf{X} + \delta_{w+u})) - \nabla \ell(f_w(\Lambda_{w+u} \mathbf{X} + \delta_{w+u}))) \mathbf{X}^\top] \right\|_\sigma \\
&\quad + \left\| \frac{1}{\lambda} \mathbb{E}[(\nabla \ell(f_w(\Lambda_{w+u} \mathbf{X} + \delta_{w+u})) - \nabla \ell(f_w(\Lambda_w \mathbf{X} + \delta_{w+u}))) \mathbf{X}^\top] \right\|_\sigma \\
&\quad + \left\| \frac{1}{\lambda} \mathbb{E}[(\nabla \ell(f_w(\Lambda_w \mathbf{X} + \delta_{w+u})) - \nabla \ell(f_w(\Lambda_w \mathbf{X} + \delta_w))) \mathbf{X}^\top] \right\|_\sigma \\
&\leq \frac{(B+1) \text{Lip}(\ell \circ f_w)}{\lambda} (\|\delta_{w+u} - \delta_w\|_2 + \|\Lambda_{w+u} - \Lambda_w\|_\sigma) \\
&\quad + (B+1) e^2 \left(\prod_{i=1}^L \|w_i\|_\sigma \right) \sum_{i=1}^d \left[\frac{\|u_i\|_\sigma}{\|w_i\|_\sigma} + B \left(\prod_{j=1}^i \|w_j\|_\sigma \right) \sum_{j=1}^i \frac{\|u_j\|_\sigma}{\|w_j\|_\sigma} \right],
\end{aligned}$$

where the last inequality follows from Lemma 3 in [36]. As a result,

$$\begin{aligned}
& \|\delta_{w+u} - \delta_w\|_2 + \|\Lambda_{w+u} - \Lambda_w\|_\sigma \\
&\leq \frac{\lambda}{\lambda - (B+1) \text{Lip}(\ell \circ f_w)} \left[(B+1) e^2 \left(\prod_{i=1}^L \|w_i\|_\sigma \right) \sum_{i=1}^d \left[\frac{\|u_i\|_\sigma}{\|w_i\|_\sigma} + B \left(\prod_{j=1}^i \|w_j\|_\sigma \right) \sum_{j=1}^i \frac{\|u_j\|_\sigma}{\|w_j\|_\sigma} \right] \right].
\end{aligned}$$

Then, we can bound the change in the loss function caused by perturbing w at any $\|x\|_2 \leq B$ with any norm-bounded $\|u_i\|_\sigma \leq \frac{1}{L} \|w_i\|_\sigma$:

$$\begin{aligned}
& \|f_{w+u}(\Lambda_{w+u} \mathbf{X} + \delta_{w+u}) - f_w(\Lambda_w \mathbf{X} + \delta_w)\|_2 \\
&\leq \|f_{w+u}(\Lambda_{w+u} \mathbf{X} + \delta_{w+u}) - f_w(\Lambda_{w+u} \mathbf{X} + \delta_{w+u})\|_2 \\
&\quad + \|f_w(\Lambda_{w+u} \mathbf{X} + \delta_{w+u}) - f_w(\Lambda_w \mathbf{X} + \delta_{w+u})\|_2 \\
&\quad + \|f_w(\Lambda_w \mathbf{X} + \delta_{w+u}) - f_w(\Lambda_w \mathbf{X} + \delta_w)\|_2 \\
&\leq eB \left(\prod_{i=1}^L \|w_i\|_\sigma \right) \sum_{i=1}^L \frac{\|u_i\|_2}{\|w_i\|_2} + (1+B) \left(\prod_{i=1}^d \|w_i\|_\sigma \right) \\
&\quad \frac{e^2}{\lambda - (B+1) \text{Lip}(\nabla f_w)} \sum_{i=1}^L \left[\frac{\|u_i\|_\sigma}{\|w_i\|_\sigma} + B \left(\prod_{j=1}^i \|w_j\|_\sigma \right) \sum_{j=1}^i \frac{\|u_j\|_\sigma}{\|w_j\|_\sigma} \right].
\end{aligned}$$

Now, for a fixed weight vector \tilde{w} we consider a multivariate Gaussian distribution Q with zero-mean and diagonal covaraince matrix for perturbation u where each entry u_i has standard deviation

$\kappa_i = \frac{\|\tilde{\mathbf{w}}_i\|_\sigma}{\sqrt[k]{\prod_{i=1}^L \|\tilde{\mathbf{w}}_i\|_\sigma}} \kappa$ with κ chosen as

$$\kappa = \frac{\gamma}{8e^5 L \sqrt{2d \log(4dL)} B \left(\prod_{i=1}^L \|\tilde{\mathbf{w}}_i\|_\sigma \right) \left(1 + \frac{\lambda}{\lambda - (1+B) \text{Lip}(\nabla f_{\mathbf{w}})} \sum_{i=1}^L \prod_{j=1}^i \|\tilde{\mathbf{w}}_j\|_\sigma \right)}. \quad (105)$$

Also, for any \mathbf{w} which satisfies $\|\mathbf{w}_i\|_\sigma - \|\tilde{\mathbf{w}}_i\|_\sigma \leq \frac{\eta}{4L} \|\tilde{\mathbf{w}}_i\|_\sigma$, we have $\overline{\text{Lip}}(\ell \circ f_{\mathbf{w}}) \leq e^{\eta/2} \lambda (1 - \eta) \leq (1 - \eta/2) \lambda$. Therefore,

$$\begin{aligned} & \text{KL}(P_{\mathbf{w}+\mathbf{u}} \| Q) \\ & \leq \sum_{i=1}^d \frac{\|\mathbf{w}_i\|_F^2}{2\kappa_i^2} \\ & \leq O \left(L^2 B^2 d \log(dL) \frac{(\prod_{i=1}^L \|\tilde{\mathbf{w}}_i\|_\sigma^2) \left(1 + \frac{1}{\lambda - (1+B) \text{Lip}(\nabla f_{\mathbf{w}})} \sum_{i=1}^L \prod_{j=1}^i \|\tilde{\mathbf{w}}_j\|_\sigma \right)^2}{\gamma^2} \sum_{i=1}^d \frac{\|\mathbf{w}_i\|_F^2}{\|\tilde{\mathbf{w}}_i\|_\sigma^2} \right) \\ & \leq O \left(L^2 B^2 d \log(dL) \frac{(\prod_{i=1}^L \|\mathbf{w}_i\|_\sigma^2) \left(1 + \frac{1}{\lambda - (1+B) \text{Lip}(\nabla f_{\mathbf{w}})} \sum_{i=1}^L \prod_{j=1}^i \|\mathbf{w}_j\|_\sigma \right)^2}{\gamma^2} \sum_{i=1}^d \frac{\|\mathbf{w}_i\|_F^2}{\|\mathbf{w}_i\|_\sigma^2} \right) \end{aligned}$$

Now we plug the above result into [36]’s Lemma 1, implying that given a fixed underlying distribution P and any $\xi > 0$ with probability at least $1 - \xi$ for any \mathbf{w} satisfying $\|\mathbf{w}_i\|_\sigma - \|\tilde{\mathbf{w}}_i\|_\sigma \leq \frac{\eta}{4L} \|\tilde{\mathbf{w}}_i\|_\sigma$ we have

$$\mathcal{L}_{0-1}^{\text{adv}}(\mathbf{w}) - \hat{\mathcal{L}}_\gamma^{\text{adv}}(\mathbf{w}) \leq O \left(\sqrt{\frac{B^2 L^2 d \log(Ld) \lambda^2 \left(\prod_{i=1}^L \|\mathbf{w}_i\|_\sigma \sum_{i=1}^L \frac{\|\mathbf{w}_i\|_F^2}{\|\mathbf{w}_i\|_\sigma^2} \right)^2 + \log \frac{m}{\xi}}{m \gamma^2 (\lambda - (1+B) \text{Lip}(\nabla f_{\mathbf{w}}))^2}} \right). \quad (106)$$

Now we use a cover of size $O(\frac{L}{\eta} \log M)$ points where for any feasible $\|\mathbf{w}_i\|_\sigma$ we can find a point a_i in the cover such that $|\|\mathbf{w}_i\|_\sigma - a_i| \leq \frac{\eta}{4L} a_i$. As a result, we can cover the space of feasible \mathbf{w}_i ’s with $O((\frac{L}{\eta} \log M)^L L)$ number of points. This proves that for a fixed underlying distribution for every $\xi > 0$, with probability at least $\xi > 0$ for any feasible norm-bounded \mathbf{w} we have

$$\mathcal{L}_{0-1}^{\text{adv}}(\mathbf{w}) - \hat{\mathcal{L}}_\gamma^{\text{adv}}(\mathbf{w}) \leq O \left(\sqrt{\frac{B^2 L^2 d \log(Ld) \lambda^2 \left(\prod_{i=1}^L \|\mathbf{w}_i\|_\sigma \sum_{i=1}^L \frac{\|\mathbf{w}_i\|_F^2}{\|\mathbf{w}_i\|_\sigma^2} \right)^2 + L \log \frac{mL \log(M)}{\eta \xi}}{m \gamma^2 (\lambda - (1+B) \text{Lip}(\nabla \ell \circ f_{\mathbf{w}}))^2}} \right). \quad (107)$$

To apply the result to the network of n nodes, we apply a union bound to have the bound hold simultaneously for the distribution of every node, which proves for every $\xi > 0$ with probability at least $1 - \xi$ the average worst-case loss of the nodes satisfies the following margin-based bound:

$$\mathcal{L}_{0-1}^{\text{adv}}(\mathbf{w}) - \hat{\mathcal{L}}_\gamma^{\text{adv}}(\mathbf{w}) \leq O \left(\sqrt{\frac{B^2 L^2 d \log(Ld) \lambda^2 \left(\prod_{i=1}^L \|\mathbf{w}_i\|_\sigma \sum_{i=1}^L \frac{\|\mathbf{w}_i\|_F^2}{\|\mathbf{w}_i\|_\sigma^2} \right)^2 + L \log \frac{nmL \log(M)}{\eta \xi}}{m \gamma^2 (\lambda - (1+B) \text{Lip}(\nabla f_{\mathbf{w}}))^2}} \right). \quad (108)$$

Therefore, the proof is complete.

Appendix G Proof of Theorem 4

Define random vector $\mathbf{U} = \Lambda \mathbf{X} + \delta$. According to the definition of optimal transport cost $W_c(P_{\mathbf{X}}, P_{\mathbf{U}})$ for quadratic $c(\mathbf{x}, \mathbf{u}) = \frac{1}{2} \|\mathbf{x} - \mathbf{u}\|_2^2$,

$$W_c(P_{\mathbf{X}}, P_{\mathbf{U}}) := \min_{P_{\mathbf{X}}, \mathbf{U} \in \Pi(P_{\mathbf{X}}, P_{\mathbf{U}})} \mathbb{E} \left[\frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_2^2 \right] \quad (109)$$

where $\Pi(P_{\mathbf{X}}, P_{\mathbf{U}})$ contains any joint distribution $P_{\mathbf{X}, \mathbf{U}}$ with marginals $P_{\mathbf{X}}, P_{\mathbf{U}}$. One distribution in $\Pi(P_{\mathbf{X}}, P_{\mathbf{U}})$ is the joint distribution of $(\mathbf{X}, \Lambda \mathbf{X} + \delta)$ implying that

$$\begin{aligned}
W_c(P_{\mathbf{X}}, P_{\mathbf{U}}) &\leq \frac{1}{2} \mathbb{E}[\|\mathbf{X} - \Lambda \mathbf{X} - \delta\|_2^2] \\
&= \frac{1}{2} \mathbb{E}[\|(I - \Lambda)\mathbf{X} - \delta\|_2^2] \\
&\stackrel{(a)}{\leq} \mathbb{E}[\|(I - \Lambda)\mathbf{X}\|_2^2] + \|\delta\|_2^2 \\
&\stackrel{(b)}{\leq} \text{Tr}((I - \Lambda)(I - \Lambda)^\top \mathbb{E}[\mathbf{X}\mathbf{X}^\top]) + \|\delta\|_2^2 \\
&\stackrel{(c)}{\leq} \lambda \text{Tr}((I - \Lambda)(I - \Lambda)^\top) + \|\delta\|_2^2 \\
&\stackrel{(d)}{\leq} \lambda \|I - \Lambda\|_F^2 + \|\delta\|_2^2 \\
&\leq \max\{\lambda, 1\} (\|I - \Lambda\|_F^2 + \|\delta\|_2^2).
\end{aligned}$$

In the above, (a) holds since for every two vectors $\mathbf{u}_1, \mathbf{u}_2$ we have $\|\mathbf{u}_1 + \mathbf{u}_2\|_2^2 = \|\mathbf{u}_1\|_2^2 + \|\mathbf{u}_2\|_2^2 + 2\mathbf{u}_1^\top \mathbf{u}_2 \leq 2(\|\mathbf{u}_1\|_2^2 + \|\mathbf{u}_2\|_2^2)$. (b) follows from the fact that $\mathbb{E}[\|(I - \Lambda)\mathbf{X}\|_2^2] = \mathbb{E}[\text{Tr}((I - \Lambda)\mathbf{X}\mathbf{X}^\top(I - \Lambda)^\top)] = \text{Tr}((I - \Lambda)(I - \Lambda)^\top \mathbb{E}[\mathbf{X}\mathbf{X}^\top])$. (c) holds because of the theorem's assumption implying that $\mathbb{E}[\mathbf{X}\mathbf{X}^\top] \leq \lambda I$. Last, (d) holds because we have $\text{Tr}(AA^\top) = \|A\|_F^2$ for every A . Therefore, the proof is complete.