

## A PSEUDO-CODE OF MOMENTUM TRACKING

The pseudo-code for Momentum Tracking is given in Alg. 1, where  $\mathbf{Transmit}_{i \rightarrow j}(\cdot)$  denotes that node  $i$  transmits parameters to node  $j$  and  $\mathbf{Receive}_{i \leftarrow j}(\cdot)$  denotes that node  $i$  receives parameters from node  $j$ .

---

**Algorithm 1:** Update rules of Momentum Tracking at node  $i$ .

---

```

1: Input: Step size  $\eta > 0$ ,  $\beta \in (0, 1]$ , mixing matrix  $\mathbf{W}$ . Initialize  $\mathbf{c}_i$  and  $\mathbf{u}_i$  to
    $\frac{1}{1-\beta}(\nabla F_i(\mathbf{x}_i^{(0)}; \xi_i^{(0)}) - \frac{1}{N} \sum_j \nabla F_j(\mathbf{x}_j^{(0)}; \xi_j^{(0)}))$  for all  $i \in V$  and  $\mathbf{x}_i$  with the same parameter.
2: for  $r = 0, \dots, R$  do
3:    $\mathbf{u}_i^{(r+1)} \leftarrow \beta \mathbf{u}_i^{(r)} + \nabla F_i(\mathbf{x}_i^{(r)}; \xi_i^{(r)})$ .
4:   for  $j \in \mathcal{N}_i$  do
5:      $\mathbf{Transmit}_{i \rightarrow j}(\mathbf{x}_i^{(r)})$  and  $\mathbf{Receive}_{i \leftarrow j}(\mathbf{x}_j^{(r)})$ .
6:      $\mathbf{Transmit}_{i \rightarrow j}(\mathbf{c}_i^{(r)} - \mathbf{u}_i^{(r+1)})$  and  $\mathbf{Receive}_{i \leftarrow j}(\mathbf{c}_j^{(r)} - \mathbf{u}_j^{(r+1)})$ .
7:   end for
8:    $\mathbf{x}_i^{(r+1)} \leftarrow \sum_{j \in \mathcal{N}_i^+} W_{ij} \mathbf{x}_j^{(r)} - \eta (\mathbf{u}_i^{(r+1)} - \mathbf{c}_i^{(r)})$ .
9:    $\mathbf{c}_i^{(r+1)} \leftarrow \sum_{j \in \mathcal{N}_i^+} W_{ij} (\mathbf{c}_j^{(r)} - \mathbf{u}_j^{(r+1)}) + \mathbf{u}_i^{(r+1)}$ .
10: end for

```

---

## B ADDITIONAL DISCUSSION ABOUT CONVERGENCE RATE

Because Momentum Tracking is equivalent to Gradient Tracking when  $\beta = 0$ , Theorem 1 also provides the convergence rate of Gradient Tracking. In this section, we compare the convergence rate of Gradient Tracking provided in Theorem 1 to that provided by Koloskova et al. (2021).

From Theorem 1, we get the following statement.

**Corollary 1.** *Suppose that  $\beta = 0$  and the assumptions of Theorem 1 hold. Then, for any  $R \geq 1$ , there exists a step size  $\eta$  such that the average parameter  $\bar{\mathbf{x}} := \frac{1}{N} \sum_i \mathbf{x}_i$  generated by Eqs. (7-9) satisfies*

$$\frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \leq \mathcal{O} \left( \sqrt{\frac{r_0 \sigma^2 L}{NR}} + \left( \frac{r_0 \sigma L}{p^2 R} \right)^{\frac{2}{3}} + \frac{L r_0}{p^2 R} \right), \quad (11)$$

where  $r_0 := f(\bar{\mathbf{x}}^{(0)}) - f^*$ .

Then, under Assumptions 1, 2, 3, and 4, Koloskova et al. (2021) provided the convergence rate of Gradient Tracking as follows.

**Theorem 3** ((Koloskova et al., 2021)). *Suppose that Assumptions 1, 2, 3, and 4 hold. Then, for any round  $R > \frac{2}{p} \log(\frac{50}{p}(1 + \log \frac{1}{p}))$ , there exists a step size  $\eta$  that satisfies that the average parameter  $\bar{\mathbf{x}} := \frac{1}{N} \sum_i \mathbf{x}_i$  generated by Gradient Tracking satisfies*

$$\frac{1}{R} \sum_{r=0}^{R-1} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \leq \tilde{\mathcal{O}} \left( \sqrt{\frac{r_0 \sigma^2 L}{NR}} + \left( \frac{r_0 \sigma L}{(\sqrt{pc} + p\sqrt{N})R} \right)^{\frac{2}{3}} + \frac{L r_0}{pcR} \right), \quad (12)$$

where  $\tilde{\mathcal{O}}(\cdot)$  hides the polylogarithmic factors,  $c := 1 - \min\{\lambda_{\min}, 0\}^2$ , and  $\lambda_{\min}$  is the minimum eigenvalue of  $\mathbf{W}$ .

Comparing the convergence rates in Eqs. (11) and (12), the convergence rate in Eq. (12) is tighter than that in Eq. (11) because  $c \geq p$  for any mixing matrix  $\mathbf{W}$ . However, because the convergence rate in Eq. (12) holds only when the number of round  $R$  is larger than  $\frac{2}{p} \log(\frac{50}{p}(1 + \log \frac{1}{p}))$ , Theorem 3 can not describe the behavior of the convergence rate at the beginning of the training. In contrast, Corollary 1 provides the convergence rate for Gradient Tracking that holds for any  $R \geq 1$ .

## C ADDITIONAL EXPERIMENTS

### C.1 RESULTS WITH VARIOUS NETWORK TOPOLOGIES

We evaluated Momentum Tracking in more detail by changing the underlying network topology. Table 4 lists the test accuracy of all comparison methods when we set the underlying network topology to be a hypercube or a semantic exponential graph.

Table 4 indicates that when the data distributions held by each node are statistically homogeneous (i.e., 10-class), DSGDm, QG-DSGDm, DecentLaM, and Momentum Tracking are comparable and outperform DSGD and Gradient Tracking for all network topologies. When the data distributions are heterogeneous (i.e., 4-class), the results show that Momentum Tracking is more robust to data heterogeneity than DSGDm, QG-DSGDm, and DecentLaM and outperforms all comparison methods for all network topologies. Therefore, the results indicate that Momentum Tracking is robust to data heterogeneity regardless of the underlying network topology.

Table 4: Test accuracy on CIFAR-10 with different underlying network topologies.

	CIFAR-10			
	Hypercube		Semantic Exponential Graph	
	10-class	4-class	10-class	4-class
DSGD	63.3 $\pm$ 0.65	55.9 $\pm$ 4.11	64.0 $\pm$ 0.26	60.7 $\pm$ 1.82
Gradient Tracking	61.0 $\pm$ 1.34	60.2 $\pm$ 1.13	62.4 $\pm$ 0.53	62.4 $\pm$ 0.89
DSGDm	<b>73.2 <math>\pm</math> 0.09</b>	45.0 $\pm$ 5.90	<b>73.4 <math>\pm</math> 0.13</b>	51.5 $\pm$ 7.80
QG-DSGDm	73.0 $\pm$ 0.31	62.9 $\pm$ 3.68	<b>73.4 <math>\pm</math> 0.58</b>	70.2 $\pm$ 1.09
DecentLaM	72.9 $\pm$ 0.24	69.1 $\pm$ 4.05	72.9 $\pm$ 0.73	71.2 $\pm$ 1.72
Momentum Tracking	72.8 $\pm$ 0.15	<b>72.7 <math>\pm</math> 0.28</b>	72.7 $\pm$ 0.33	<b>72.9 <math>\pm</math> 0.07</b>

### C.2 INITIAL VALUE ANALYSIS

In this section, we discuss the initial values of  $\mathbf{c}_i$  and  $\mathbf{u}_i$ . Table 5 lists the test accuracy for Momentum Tracking when we initialize  $\mathbf{c}_i$  and  $\mathbf{u}_i$  to zero and when we initialize  $\mathbf{c}_i$  and  $\mathbf{u}_i$  as in Theorem 1. The results indicate that the test accuracy are almost equivalent on both settings. Hence, Theorem 1 requires  $\mathbf{c}_i$  and  $\mathbf{u}_i$  to be initialized to  $\frac{1}{1-\beta}(\nabla F_i(\mathbf{x}_i^{(0)}; \xi_i^{(0)}) - \frac{1}{N} \sum_{j=1}^N \nabla F_j(\mathbf{x}_j^{(0)}; \xi_j^{(0)}))$ . However, in practice,  $\mathbf{c}_i$  and  $\mathbf{u}_i$  can be initialized to zeros without any impact on accuracy.

Table 5: Test accuracy on FashionMNIST, SVHN, and CIFAR-10 with LeNet. “ $k$ -class” indicates that each node has only the data of randomly selected  $k$  classes.

	FashionMNIST			
	10-class	8-class	6-class	4-class
Momentum Tracking	89.5 $\pm$ 0.36	89.4 $\pm$ 0.05	88.9 $\pm$ 0.47	86.8 $\pm$ 1.56
Momentum Tracking ( $\mathbf{c}_i^{(0)} = \mathbf{u}_i^{(0)} = \mathbf{0}$ )	89.5 $\pm$ 0.38	89.4 $\pm$ 0.04	88.7 $\pm$ 0.63	85.8 $\pm$ 1.53
	SVHN			
	10-class	8-class	6-class	4-class
Momentum Tracking	92.6 $\pm$ 0.32	92.4 $\pm$ 0.40	92.3 $\pm$ 0.23	91.7 $\pm$ 0.53
Momentum Tracking ( $\mathbf{c}_i^{(0)} = \mathbf{u}_i^{(0)} = \mathbf{0}$ )	92.5 $\pm$ 0.34	92.3 $\pm$ 0.50	92.2 $\pm$ 0.29	92.0 $\pm$ 0.81
	CIFAR-10			
	10-class	8-class	6-class	4-class
Momentum Tracking	72.9 $\pm$ 0.59	73.0 $\pm$ 0.49	72.6 $\pm$ 0.41	70.7 $\pm$ 1.38
Momentum Tracking ( $\mathbf{c}_i^{(0)} = \mathbf{u}_i^{(0)} = \mathbf{0}$ )	72.8 $\pm$ 0.35	72.9 $\pm$ 0.32	73.0 $\pm$ 0.41	70.7 $\pm$ 2.00

### C.3 LEARNING CURVES

In this section, we present the learning curves for the results whose final accuracy are presented in Tables 2 and 3. Figs. 3, 4, and 5 show the learning curves for FashionMNSIT, SVHN, and CIFAR-10, respectively, with LeNet. Figs. 6 and 7 show the learning curves for CIFAR-10 with VGG-11 and ResNet-34, respectively.

When the data distributions are statistically homogeneous (i.e., 10-class), the results indicate that DSGDm, QG-DSGDm, DecentLaM, and Momentum Tracking are comparable and can achieve high accuracy faster than DSGD and Gradient Tracking. When the data distributions are statistically heterogeneous (e.g., 2-class and 4-class), the results indicate that the learning curves for Momentum Tracking are more stable than those for DSGDm, QG-DSGDm, and DecentLaM, and Momentum Tracking outperforms all comparison methods. In particular, in Figs. 6 and 7, the accuracy of DSGD, DSGDm, QG-DSGDm, and DecentLaM continue to oscillate in the final training phase in the 2-class setting, whereas the accuracy of Momentum Tracking and Gradient Tracking converge in the 2-class setting as well as in the 10-class setting. Therefore, Momentum Tracking is more robust to data heterogeneity than DSGDm, QG-DSGDm, and DecentLaM.

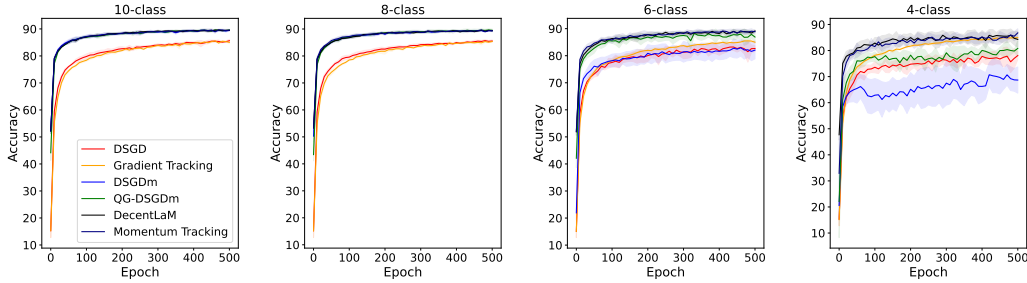


Figure 3: Learning curves on FashionMNIST. The accuracy is evaluated per 10 epochs.

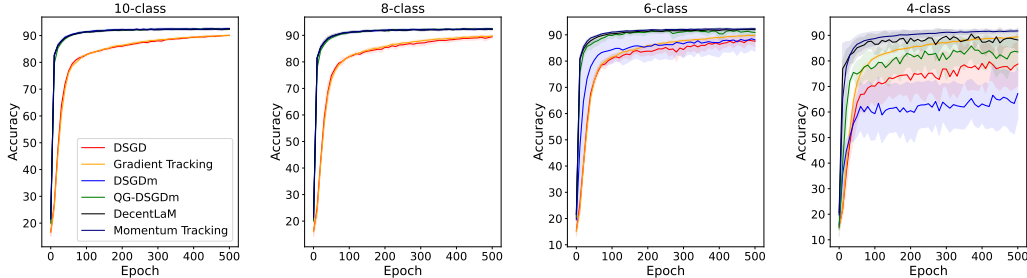


Figure 4: Learning curves on SVHN. The accuracy is evaluated per 10 epochs.

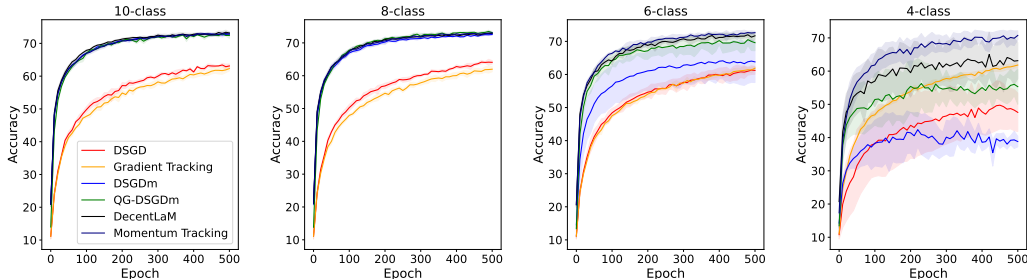


Figure 5: Learning curves on CIFAR-10. The accuracy is evaluated per 10 epochs.

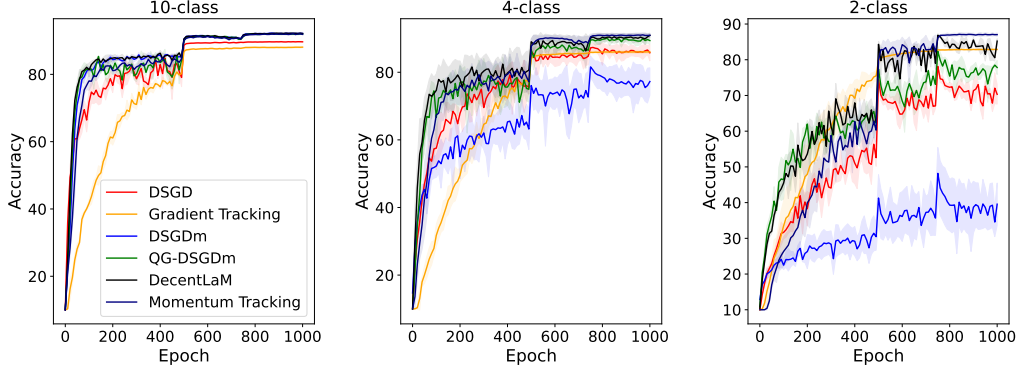


Figure 6: Learning curves on CIFAR-10 with VGG-11. The accuracy is evaluated per 10 epochs.

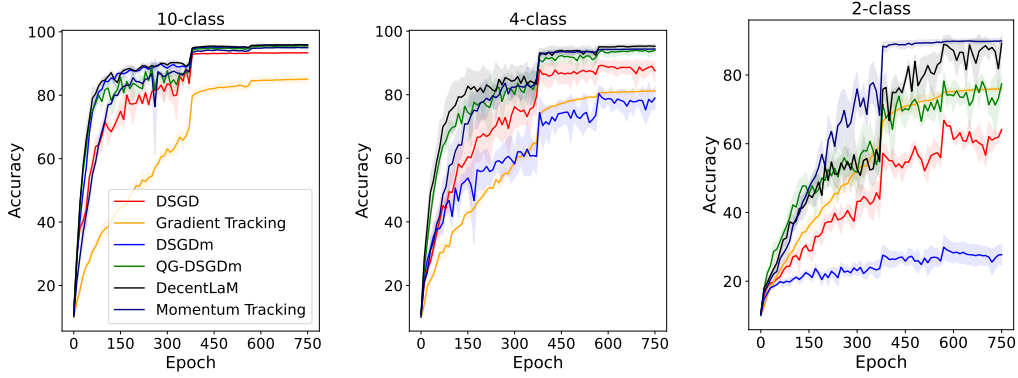


Figure 7: Learning curves on CIFAR-10 with ResNet-34. The accuracy is evaluated per 10 epochs.

#### C.4 SYNTHETIC EXPERIMENT

In this section, we evaluate the convergence rate in more detail using a synthetic dataset. Following the previous work (Koloskova et al., 2020), we set the dimension of parameter  $d$  to 50, the number of nodes  $N$  to 25, and the network topology to a ring consisting of  $N$  nodes. We then defined the local objective function  $f_i(\mathbf{x})$  to be  $\frac{1}{2}\|\mathbf{A}_i\mathbf{x} - \mathbf{b}_i\|^2$  where  $\mathbf{A}_i := i/\sqrt{N}$  and  $\mathbf{b}_i$  are sampled from  $\mathcal{N}(\mathbf{0}, \zeta^2/i^2\mathbf{1})$ , and we defined the stochastic gradient  $\nabla F_i(\mathbf{x}; \xi_i)$  to be  $\nabla f_i(\mathbf{x}) + \epsilon$  where  $\epsilon$  is drawn from  $\mathcal{N}(\mathbf{0}; \sigma^2/d\mathbf{1})$ . For all comparison methods, we set the step size  $\eta$  to  $1.0 \times 10^{-4}$ .

Figs. 8 and 9 illustrate  $\|\nabla f(\bar{\mathbf{x}})\|^2$  with respect to the number of rounds  $r$  when we vary data heterogeneity  $\zeta^2$  as  $\{0, 25, 50\}$  and set  $\sigma^2$  to one. The results show that Momentum Tracking converges in the same manner regardless of data heterogeneity  $\zeta^2$ . On the other hand, for DSGDm, QG-DSGDm, and DecentLaM,  $\|\nabla f(\bar{\mathbf{x}})\|^2$  increases as data heterogeneity  $\zeta^2$  increases. Hence, these results are consistent with our theoretical analysis that the convergence rate of Momentum Tracking is independent of data heterogeneity.

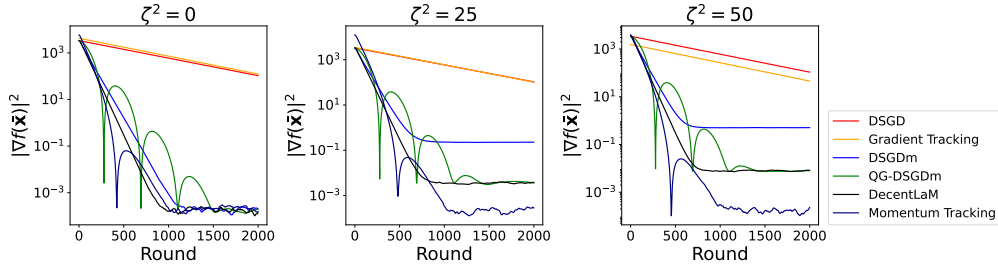


Figure 8: Comparison of the convergence in the initial training phase in the synthetic experiments.

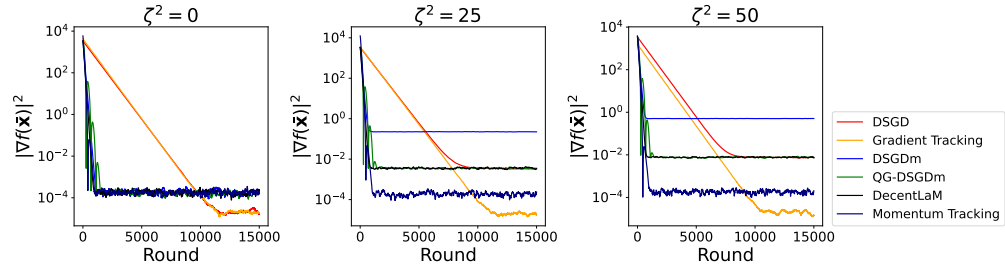


Figure 9: Comparison of the convergence in the synthetic experiments.

## D PROOF OF THEOREM 1

### D.1 TECHNICAL LEMMA

**Lemma 1.** For any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ ,  $\gamma > 0$ , it holds that

$$\|\mathbf{x} + \mathbf{y}\|^2 \leq (1 + \gamma)\|\mathbf{x}\|^2 + (1 + \gamma^{-1})\|\mathbf{y}\|^2. \quad (13)$$

**Lemma 2.** For any  $\mathbf{a}_1, \dots, \mathbf{a}_N$ , it holds that

$$\left\| \sum_{i=1}^N \mathbf{a}_i \right\|^2 \leq N \sum_{i=1}^N \|\mathbf{a}_i\|^2. \quad (14)$$

**Lemma 3.** For any  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ ,  $\gamma > 0$ , it holds that

$$2\langle \mathbf{x}, \mathbf{y} \rangle \leq \gamma\|\mathbf{x}\|^2 + \gamma^{-1}\|\mathbf{y}\|^2. \quad (15)$$

### D.2 MOMENTUM TRACKING IN MATRIX NOTATION

We define  $\mathbf{U}^{(r)}$ ,  $\mathbf{X}^{(r)}$ ,  $\mathbf{C}^{(r)}$ ,  $\nabla F(\mathbf{X}^{(r)}; \xi^{(r)})$ , and  $\nabla f(\mathbf{X}^{(r)})$  as follows:

$$\begin{aligned} \mathbf{U}^{(r)} &:= (\mathbf{u}_1^{(r)}, \dots, \mathbf{u}_N^{(r)}), \quad \mathbf{X}^{(r)} := (\mathbf{x}_1^{(r)}, \dots, \mathbf{x}_N^{(r)}), \quad \mathbf{C}^{(r)} := (\mathbf{c}_1^{(r)}, \dots, \mathbf{c}_N^{(r)}), \\ \nabla F(\mathbf{X}^{(r)}; \xi^{(r)}) &:= (\nabla F_1(\mathbf{x}_1^{(r)}; \xi_1^{(r)}), \dots, \nabla F_N(\mathbf{x}_N^{(r)}; \xi_N^{(r)})), \\ \nabla f(\mathbf{X}^{(r)}) &:= (\nabla f_1(\mathbf{x}_1^{(r)}), \dots, \nabla f_N(\mathbf{x}_N^{(r)})). \end{aligned}$$

Then, the update rule of Momentum Tracking can then be rewritten as follows:

$$\mathbf{U}^{(r+1)} = \beta \mathbf{U}^{(r)} + \nabla F(\mathbf{X}^{(r)}; \xi^{(r)}), \quad (16)$$

$$\mathbf{X}^{(r+1)} = \mathbf{X}^{(r)} \mathbf{W} - \eta (\mathbf{U}^{(r+1)} - \mathbf{C}^{(r)}), \quad (17)$$

$$\mathbf{C}^{(r+1)} = (\mathbf{C}^{(r)} - \mathbf{U}^{(r+1)}) \mathbf{W} + \mathbf{U}^{(r+1)}, \quad (18)$$

where  $\mathbf{U}^{(0)}$  and  $\mathbf{C}^{(0)}$  are initialized as follows:

$$\begin{aligned} \mathbf{U}^{(0)} &= \frac{1}{1-\beta} (\nabla F(\mathbf{X}^{(0)}; \xi^{(0)}) - \frac{1}{N} \nabla F(\mathbf{X}^{(0)}; \xi^{(0)}) \mathbf{1} \mathbf{1}^\top), \\ \mathbf{C}^{(0)} &= \frac{1}{1-\beta} (\nabla F(\mathbf{X}^{(0)}; \xi^{(0)}) - \frac{1}{N} \nabla F(\mathbf{X}^{(0)}; \xi^{(0)}) \mathbf{1} \mathbf{1}^\top). \end{aligned}$$

### D.3 ADDITIONAL NOTATION

We define the update rules of  $\mathbf{d}_i$  and  $\mathbf{e}_i$  as follows:

$$\mathbf{d}_i^{(r+1)} = \beta \mathbf{d}_i^{(r)} + \nabla f_i(\bar{\mathbf{x}}^{(r)}), \quad (19)$$

$$\mathbf{e}_i^{(r+1)} = \beta \mathbf{e}_i^{(r)} + \nabla f(\bar{\mathbf{x}}^{(r)}), \quad (20)$$

where  $\mathbf{d}_i^{(0)} = \frac{1}{1-\beta} (\nabla f_i(\bar{\mathbf{x}}^{(0)}) - \nabla f(\bar{\mathbf{x}}^{(0)}))$  and  $\mathbf{e}_i^{(0)} = \mathbf{0}$ . Note that it holds that  $\bar{\mathbf{d}}^{(r)} = \bar{\mathbf{e}}^{(r)}$  for any round  $r \geq 0$ . Then, we define  $\mathbf{D}$  and  $\mathbf{E}$  as follows:

$$\mathbf{D}^{(r)} := (\mathbf{d}_1^{(r)}, \dots, \mathbf{d}_N^{(r)}), \quad \mathbf{E}^{(r)} := (\mathbf{e}_1^{(r)}, \dots, \mathbf{e}_N^{(r)}).$$

The update rules of  $\mathbf{D}$  and  $\mathbf{E}$  can be written as follows:

$$\mathbf{D}^{(r+1)} = \beta \mathbf{D}^{(r)} + \nabla f(\bar{\mathbf{X}}^{(r)}), \quad (21)$$

$$\mathbf{E}^{(r+1)} = \beta \mathbf{E}^{(r)} + \frac{1}{N} \nabla f(\bar{\mathbf{X}}^{(r)}) \mathbf{1} \mathbf{1}^\top, \quad (22)$$

where  $\mathbf{D}^{(0)}$  and  $\mathbf{E}^{(0)}$  are initialized as follows:

$$\begin{aligned}\mathbf{D}^{(0)} &= \frac{1}{1-\beta}(\nabla f(\bar{\mathbf{X}}^{(0)}) - \frac{1}{N}\nabla f(\bar{\mathbf{X}}^{(0)})\mathbf{1}\mathbf{1}^\top), \\ \mathbf{E}^{(0)} &= \mathbf{0}.\end{aligned}$$

Note that  $\mathbf{d}_i$ ,  $\mathbf{e}_i$ ,  $\mathbf{D}$ , and  $\mathbf{E}$  are the only variables used in the proof that do not need to be computed in practice in Alg. 1. We define  $\Xi$ ,  $\mathcal{E}$ , and  $\mathcal{D}$  as follows:

$$\begin{aligned}\Xi^{(r)} &:= \frac{1}{N}\mathbb{E}\left\|\mathbf{X}^{(r)} - \bar{\mathbf{X}}^{(r)}\right\|_F^2, \\ \mathcal{E}^{(r)} &:= \frac{1}{N}\mathbb{E}\left\|\mathbf{D}^{(r+1)} - \mathbf{C}^{(r)} - \mathbf{E}^{(r+1)}\right\|_F^2, \\ \mathcal{D}^{(r)} &:= \frac{1}{N}\mathbb{E}\left\|\mathbf{D}^{(r+1)} - \mathbf{D}^{(r)} - \mathbf{E}^{(r+1)} + \mathbf{E}^{(r)}\right\|_F^2.\end{aligned}$$

Inspired by Yu et al. (2019), we define  $\bar{\mathbf{z}}$  as follows:

$$\bar{\mathbf{z}}^{(r)} := \begin{cases} \bar{\mathbf{x}}^{(r)}, & \text{if } r = 0 \\ \frac{1}{1-\beta}\bar{\mathbf{x}}^{(r)} - \frac{\beta}{1-\beta}\bar{\mathbf{x}}^{(r-1)}, & \text{otherwise} \end{cases}.$$

In the following, we define  $\pm a := a - a = 0$  for any  $a$  and  $\bar{a} := \frac{1}{N}\sum_{i=1}^N a_i$  for any  $a_1, \dots, a_N$ . Then,  $\mathbb{E}[\cdot]$  denotes the expectation over all randomness that occurs during training (i.e.,  $\{\xi_i^{(r)}\}_{i,r}$ ), and  $\mathbb{E}_r[\cdot]$  denotes the expectation over the randomness that occurs at round  $r$  (i.e.,  $\{\xi_i^{(r)}\}_i$ ).

#### D.4 USEFUL LEMMA

**Lemma 4.** For any round  $r \geq 0$ , it holds that  $\bar{\mathbf{c}}^{(r)} = \mathbf{0}$ .

*Proof.* For any round  $r \geq 0$ , we have

$$\begin{aligned} \sum_{i=1}^N \mathbf{c}_i^{(r+1)} &= \sum_{i=1}^N \sum_{j=1}^N W_{ij} (\mathbf{c}_j^{(r)} - \mathbf{u}_j^{(r+1)}) + \sum_{i=1}^N \mathbf{u}_i^{(r+1)} \\ &= \sum_{j=1}^N (\mathbf{c}_j^{(r)} - \mathbf{u}_j^{(r+1)}) \sum_{i=1}^N W_{ij} + \sum_{i=1}^N \mathbf{u}_i^{(r+1)}. \end{aligned}$$

Because  $\mathbf{W}$  is a mixing matrix, we obtain

$$\sum_{i=1}^N \mathbf{c}_i^{(r+1)} = \sum_{j=1}^N \mathbf{c}_j^{(r)}.$$

Since we have

$$\sum_{i=1}^N \mathbf{c}_i^{(0)} = \frac{1}{1-\beta} \sum_{i=1}^N \left( \nabla F_i(\mathbf{x}_i^{(0)}; \xi_i^{(0)}) - \frac{1}{N} \sum_{j=1}^N \nabla F_j(\mathbf{x}_j^{(0)}; \xi_j^{(0)}) \right) = \mathbf{0},$$

we obtain the statement.  $\square$

**Lemma 5.** For any round  $r \geq 0$ , it holds that

$$\bar{\mathbf{x}}^{(r+1)} = \bar{\mathbf{x}}^{(r)} - \eta \bar{\mathbf{u}}^{(r+1)}.$$

*Proof.* We have

$$\begin{aligned} \bar{\mathbf{x}}^{(r+1)} &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N W_{ij} \mathbf{x}_j^{(r)} - \eta (\bar{\mathbf{u}}^{(r+1)} - \bar{\mathbf{c}}^{(r)}) \\ &= \frac{1}{N} \sum_{j=1}^N \mathbf{x}_j^{(r)} \sum_{i=1}^N W_{ij} - \eta (\bar{\mathbf{u}}^{(r+1)} - \bar{\mathbf{c}}^{(r)}). \end{aligned}$$

The fact that  $\mathbf{W}$  is a mixing matrix gives us

$$\bar{\mathbf{x}}^{(r+1)} = \bar{\mathbf{x}}^{(r)} - \eta (\bar{\mathbf{u}}^{(r+1)} - \bar{\mathbf{c}}^{(r)}).$$

Then, using Lemma 4, we get the statement.  $\square$

**Lemma 6.** For any round  $r \geq 0$ , it holds that

$$\bar{\mathbf{z}}^{(r+1)} - \bar{\mathbf{z}}^{(r)} = -\frac{\eta}{1-\beta} \frac{1}{N} \sum_{i=1}^N \nabla F_i(\mathbf{x}_i^{(r)}; \xi_i^{(r)}).$$

*Proof.* For any  $r \geq 1$ , we have

$$\begin{aligned} \bar{\mathbf{z}}^{(r+1)} - \bar{\mathbf{z}}^{(r)} &= \frac{1}{1-\beta} (\bar{\mathbf{x}}^{(r+1)} - \bar{\mathbf{x}}^{(r)}) - \frac{\beta}{1-\beta} (\bar{\mathbf{x}}^{(r)} - \bar{\mathbf{x}}^{(r-1)}) \\ &= -\frac{\eta}{1-\beta} \bar{\mathbf{u}}^{(r+1)} + \frac{\eta\beta}{1-\beta} \bar{\mathbf{u}}^{(r)} \\ &= -\frac{\eta}{1-\beta} \frac{1}{N} \sum_{i=1}^N \nabla F_i(\mathbf{x}_i^{(r)}; \xi_i^{(r)}), \end{aligned}$$



where we use Lemma 5. When  $r = 0$ , we have

$$\begin{aligned}\bar{\mathbf{z}}^{(1)} - \bar{\mathbf{z}}^{(0)} &= \frac{1}{1-\beta}(\bar{\mathbf{x}}^{(1)} - \bar{\mathbf{x}}^{(0)}) \\ &= -\frac{\eta}{1-\beta}\bar{\mathbf{u}}^{(1)} \\ &= -\frac{\eta}{1-\beta}\frac{1}{N}\sum_{i=1}^N\nabla F_i(\mathbf{x}_i^{(0)}; \xi_i^{(0)}),\end{aligned}$$

where we use  $\bar{\mathbf{u}}^{(0)} = \mathbf{0}$ . This concludes the proof.  $\square$

**Lemma 7.** Suppose that Assumptions 1, 2, 3, and 4 hold. For any round  $R \geq 0$ , it holds that

$$\sum_{r=0}^R \mathbb{E} \left\| \bar{\mathbf{x}}^{(r)} - \bar{\mathbf{z}}^{(r)} \right\|^2 \leq \frac{\beta^2 \eta^2}{(1-\beta)^4} \sum_{r=0}^R \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 + \frac{\beta^2 \sigma^2 \eta^2}{N(1-\beta)^4} R.$$

*Proof.* From Lemma 5, we have

$$\begin{aligned}\mathbb{E} \left\| \bar{\mathbf{x}}^{(r)} - \bar{\mathbf{z}}^{(r)} \right\|^2 &= \mathbb{E} \left\| \frac{\eta\beta}{1-\beta} \bar{\mathbf{u}}^{(r)} \right\|^2 \\ &= \frac{\beta^2 \eta^2}{(1-\beta)^2} \mathbb{E} \left\| \sum_{k=0}^{r-1} \beta^{r-k-1} \frac{1}{N} \sum_{i=1}^N \nabla F_i(\mathbf{x}_i^{(k)}; \xi_i^{(k)}) \right\|^2,\end{aligned}$$

for any  $r \geq 1$ . Defining  $s^{(r)} := \sum_{k=0}^r \beta^{r-k}$ , we obtain

$$\begin{aligned}\mathbb{E} \left\| \bar{\mathbf{x}}^{(r)} - \bar{\mathbf{z}}^{(r)} \right\|^2 &= \frac{\beta^2 \eta^2}{(1-\beta)^2} s^{(r-1)^2} \mathbb{E} \left\| \sum_{k=0}^{r-1} \frac{\beta^{r-k-1}}{s^{(r-1)}} \frac{1}{N} \sum_{i=1}^N \nabla F_i(\mathbf{x}_i^{(k)}; \xi_i^{(k)}) \right\|^2 \\ &\stackrel{(a)}{\leq} \frac{\beta^2 \eta^2}{(1-\beta)^2} s^{(r-1)} \sum_{k=0}^{r-1} \beta^{r-k-1} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla F_i(\mathbf{x}_i^{(k)}; \xi_i^{(k)}) \right\|^2 \\ &\stackrel{(6)}{\leq} \frac{\beta^2 \eta^2}{(1-\beta)^2} s^{(r-1)} \sum_{k=0}^{r-1} \beta^{r-k-1} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(k)}) \right\|^2 + \frac{\beta^2 \sigma^2 \eta^2}{N(1-\beta)^2} s^{(r-1)} \sum_{k=0}^{r-1} \beta^{r-k-1},\end{aligned}$$

where we use Jensen's inequality in (a). Using  $s^{(r-1)} \leq \frac{1}{1-\beta}$ , we obtain

$$\mathbb{E} \left\| \bar{\mathbf{x}}^{(r)} - \bar{\mathbf{z}}^{(r)} \right\|^2 \leq \frac{\beta^2 \eta^2}{(1-\beta)^3} \sum_{k=0}^{r-1} \beta^{r-k-1} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(k)}) \right\|^2 + \frac{\beta^2 \sigma^2 \eta^2}{N(1-\beta)^4}.$$

Recursive addition yields

$$\begin{aligned}\sum_{r=1}^R \mathbb{E} \left\| \bar{\mathbf{x}}^{(r)} - \bar{\mathbf{z}}^{(r)} \right\|^2 &\leq \frac{\beta^2 \eta^2}{(1-\beta)^3} \sum_{r=1}^R \sum_{k=0}^{r-1} \beta^{r-k-1} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(k)}) \right\|^2 + \frac{\beta^2 \sigma^2 \eta^2}{N(1-\beta)^4} R \\ &= \frac{\beta^2 \eta^2}{(1-\beta)^3} \sum_{k=0}^{R-1} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(k)}) \right\|^2 \sum_{r=k+1}^R \beta^{r-k-1} + \frac{\beta^2 \sigma^2 \eta^2}{N(1-\beta)^4} R \\ &\leq \frac{\beta^2 \eta^2}{(1-\beta)^4} \sum_{k=0}^{R-1} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(k)}) \right\|^2 + \frac{\beta^2 \sigma^2 \eta^2}{N(1-\beta)^4} R,\end{aligned}$$

where we use  $\sum_{r=k+1}^R \beta^{r-k-1} \leq \frac{1}{1-\beta}$  in the last inequality. From the definition of  $\bar{\mathbf{z}}^{(0)}$ , we have  $\|\bar{\mathbf{x}}^{(0)} - \bar{\mathbf{z}}^{(0)}\|^2 = 0$ . This yields the statement.  $\square$

**Lemma 8.** Suppose that Assumptions 1, 2, 3, and 4 hold. For any round  $r \geq 0$ , it holds that

$$\mathbb{E} \left\| \bar{\mathbf{x}}^{(r+1)} - \bar{\mathbf{x}}^{(r)} \right\|^2 \leq 4L^2\eta^2\Xi^{(r)} + 2\beta^2\eta^2\mathbb{E} \left\| \bar{\mathbf{u}}^{(r)} \right\|^2 + 4\eta^2\mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 + \frac{\sigma^2\eta^2}{N}.$$

*Proof.* From Lemma 5, we have

$$\begin{aligned} \mathbb{E}_r \left\| \bar{\mathbf{x}}^{(r+1)} - \bar{\mathbf{x}}^{(r)} \right\|^2 &= \eta^2 \mathbb{E}_r \left\| \beta \bar{\mathbf{u}}^{(r)} + \frac{1}{N} \sum_{i=1}^N \nabla F_i(\mathbf{x}_i^{(r)}; \xi_i^{(r)}) \right\|^2 \\ &\stackrel{(6)}{\leq} \eta^2 \left\| \beta \bar{\mathbf{u}}^{(r)} + \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 + \frac{\sigma^2\eta^2}{N} \\ &\stackrel{(13)}{\leq} 2\beta^2\eta^2 \left\| \bar{\mathbf{u}}^{(r)} \right\|^2 + 2\eta^2 \underbrace{\left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2}_T + \frac{\sigma^2\eta^2}{N}. \end{aligned}$$

Then,  $T$  can be bounded from above as follows:

$$\begin{aligned} T &= \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \pm \nabla f_i(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\ &\stackrel{(13)}{\leq} 2 \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) - \nabla f_i(\bar{\mathbf{x}}^{(r)}) \right\|^2 + 2 \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\ &\stackrel{(14)}{\leq} \frac{2}{N} \sum_{i=1}^N \left\| \nabla f_i(\mathbf{x}_i^{(r)}) - \nabla f_i(\bar{\mathbf{x}}^{(r)}) \right\|^2 + 2 \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\ &\stackrel{(5)}{\leq} \frac{2L^2}{N} \sum_{i=1}^N \left\| \mathbf{x}_i^{(r)} - \bar{\mathbf{x}}^{(r)} \right\|^2 + 2 \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2. \end{aligned}$$

Then, we obtain the statement.  $\square$

**Lemma 9.** For any round  $r \geq 0$ , it holds that

$$\mathbb{E} \left\| \bar{\mathbf{e}}^{(r+1)} \right\|^2 \leq \frac{1}{1-\beta} \sum_{k=0}^r \beta^{r-k} \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2.$$

*Proof.* We have

$$\mathbb{E} \left\| \bar{\mathbf{e}}^{(r+1)} \right\|^2 = \mathbb{E} \left\| \sum_{k=0}^r \beta^{r-k} \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2,$$

where we use  $\bar{\mathbf{e}}^{(0)} = \mathbf{0}$ . Defining  $s^{(r)} := \sum_{k=0}^r \beta^{r-k}$ , we obtain

$$\begin{aligned} \mathbb{E} \left\| \bar{\mathbf{e}}^{(r+1)} \right\|^2 &= s^{(r)2} \mathbb{E} \left\| \sum_{k=0}^r \frac{\beta^{r-k}}{s^{(r)}} \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \\ &\leq s^{(r)} \sum_{k=0}^r \beta^{r-k} \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2, \end{aligned}$$

where we use Jensen's inequality. Using  $s^{(r)} \leq \frac{1}{1-\beta}$ , we obtain the statement.  $\square$

**Lemma 10.** Suppose that Assumptions 1, 2, 3, and 4 hold. For any round  $r \geq 0$ , it holds that

$$\frac{1}{N} \mathbb{E} \left\| \mathbf{D}^{(r+1)} - \mathbf{U}^{(r+1)} \right\|_F^2 \leq \frac{L^2}{1-\beta} \sum_{k=0}^r \beta^{r-k} \Xi^{(k)} + \frac{5\sigma^2}{(1-\beta)^3}.$$

*Proof.* We have

$$\begin{aligned} & \mathbb{E} \left\| \mathbf{D}^{(r+1)} - \mathbf{U}^{(r+1)} \right\|_F^2 \\ &= \mathbb{E} \left\| \sum_{k=0}^r \beta^{r-k} (\nabla f(\bar{\mathbf{X}}^{(k)}) - \nabla F(\mathbf{X}^{(k)}; \xi^{(k)})) + \beta^{r+1} (\mathbf{D}^{(0)} - \mathbf{U}^{(0)}) \right\|_F^2. \end{aligned}$$

Defining  $s^{(r)} := \sum_{k=0}^r \beta^{r-k}$ , we obtain

$$\begin{aligned} & \mathbb{E} \left\| \mathbf{D}^{(r+1)} - \mathbf{U}^{(r+1)} \right\|_F^2 \\ &= s^{(r+1)^2} \mathbb{E} \left\| \sum_{k=0}^r \frac{\beta^{r-k}}{s^{(r+1)}} (\nabla f(\bar{\mathbf{X}}^{(k)}) - \nabla F(\mathbf{X}^{(k)}; \xi^{(k)})) + \frac{\beta^{r+1}}{s^{(r+1)}} (\mathbf{D}^{(0)} - \mathbf{U}^{(0)}) \right\|_F^2 \\ &\stackrel{(a)}{\leq} s^{(r+1)} \sum_{k=0}^r \beta^{r-k} \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}^{(k)}) - \nabla F(\mathbf{X}^{(k)}; \xi^{(k)}) \right\|_F^2 + s^{(r+1)} \beta^{r+1} \mathbb{E} \left\| \mathbf{D}^{(0)} - \mathbf{U}^{(0)} \right\|_F^2 \\ &\stackrel{(14)}{\leq} s^{(r+1)} \sum_{k=0}^r \beta^{r-k} \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}^{(k)}) - \nabla F(\mathbf{X}^{(k)}; \xi^{(k)}) \right\|_F^2 \\ &\quad + \frac{2s^{(r+1)}}{(1-\beta)^2} \beta^{r+1} \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}^{(0)}) - \nabla F(\mathbf{X}^{(0)}; \xi^{(0)}) \right\|_F^2 \\ &\quad + \frac{2s^{(r+1)}}{(1-\beta)^2} \beta^{r+1} \mathbb{E} \left\| \frac{1}{N} \nabla f(\bar{\mathbf{X}}^{(0)}) \mathbf{1}\mathbf{1}^\top - \frac{1}{N} \nabla F(\mathbf{X}^{(0)}; \xi^{(0)}) \mathbf{1}\mathbf{1}^\top \right\|_F^2 \\ &\stackrel{(6)}{\leq} s^{(r+1)} \sum_{k=0}^r \beta^{r-k} \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}^{(k)}) - \nabla f(\mathbf{X}^{(k)}) \right\|_F^2 + s^{(r+1)} \sum_{k=0}^r \beta^{r-k} N \sigma^2 + \frac{4s^{(r+1)}}{(1-\beta)^2} \beta^{r+1} N \sigma^2, \end{aligned}$$

where we use Jensen's inequality for (a) and use  $\mathbf{X}^{(0)} = \bar{\mathbf{X}}^{(0)}$  for the last inequality. Then, using  $s^{(r)} \leq \frac{1}{1-\beta}$ , we obtain

$$\begin{aligned} & \mathbb{E} \left\| \mathbf{D}^{(r+1)} - \mathbf{U}^{(r+1)} \right\|_F^2 \\ &\leq \frac{1}{1-\beta} \sum_{k=0}^r \beta^{r-k} \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}^{(k)}) - \nabla f(\mathbf{X}^{(k)}) \right\|_F^2 + \frac{N \sigma^2}{(1-\beta)^2} + \frac{4N \sigma^2}{(1-\beta)^3} \beta^{r+1} \\ &\stackrel{\beta \in [0,1]}{\leq} \frac{1}{1-\beta} \sum_{k=0}^r \beta^{r-k} \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}^{(k)}) - \nabla f(\mathbf{X}^{(k)}) \right\|_F^2 + \frac{5N \sigma^2}{(1-\beta)^3} \\ &\stackrel{(5)}{\leq} \frac{L^2}{1-\beta} \sum_{k=0}^r \beta^{r-k} \mathbb{E} \left\| \bar{\mathbf{X}}^{(k)} - \mathbf{X}^{(k)} \right\|_F^2 + \frac{5N \sigma^2}{(1-\beta)^3}. \end{aligned}$$

This concludes the proof.  $\square$

**Lemma 11.** Suppose that Assumptions 1, 2, 3, and 4 hold. For any round  $r \geq 0$ , it holds that

$$\mathbb{E} \left\| \bar{\mathbf{u}}^{(r+1)} - \bar{\mathbf{d}}^{(r+1)} \right\|^2 \leq \frac{L^2}{1-\beta} \sum_{k=0}^r \beta^{r-k} \Xi^{(k)} + \frac{\sigma^2}{N(1-\beta)^2}$$

*Proof.* We have

$$\mathbb{E} \left\| \bar{\mathbf{u}}^{(r+1)} - \bar{\mathbf{d}}^{(r+1)} \right\|^2 = \mathbb{E} \left\| \sum_{k=0}^r \beta^{r-k} \left( \frac{1}{N} \sum_{i=1}^N \nabla F_i(\mathbf{x}_i^{(k)}; \xi_i^{(k)}) - \nabla f(\bar{\mathbf{x}}^{(k)}) \right) \right\|^2,$$

where we use  $\bar{\mathbf{u}}^{(0)} = \bar{\mathbf{d}}^{(0)} = \mathbf{0}$ . Defining  $s^{(r)} := \sum_{k=0}^r \beta^{r-k}$ , we obtain

$$\begin{aligned}
& \mathbb{E} \left\| \bar{\mathbf{u}}^{(r+1)} - \bar{\mathbf{d}}^{(r+1)} \right\|^2 \\
&= s^{(r)2} \mathbb{E} \left\| \sum_{k=0}^r \frac{\beta^{r-k}}{s^{(r)}} \left( \frac{1}{N} \sum_{i=1}^N \nabla F_i(\mathbf{x}_i^{(k)}; \xi_i^{(k)}) - \nabla f(\bar{\mathbf{x}}^{(k)}) \right) \right\|^2 \\
&\stackrel{(a)}{\leq} s^{(r)} \sum_{k=0}^r \beta^{r-k} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla F_i(\mathbf{x}_i^{(k)}; \xi_i^{(k)}) - \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \\
&\stackrel{(6)}{\leq} s^{(r)} \sum_{k=0}^r \beta^{r-k} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(k)}) - \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 + s^{(r)} \sum_{k=0}^r \beta^{r-k} \frac{\sigma^2}{N} \\
&\stackrel{(14)}{\leq} s^{(r)} \sum_{k=0}^r \beta^{r-k} \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left\| \nabla f_i(\mathbf{x}_i^{(k)}) - \nabla f_i(\bar{\mathbf{x}}^{(k)}) \right\|^2 + s^{(r)} \sum_{k=0}^r \beta^{r-k} \frac{\sigma^2}{N},
\end{aligned}$$

where we use Jensen's inequality in (a). Then, using  $s^{(r)} \leq \frac{1}{1-\beta}$ , we obtain

$$\begin{aligned}
& \mathbb{E} \left\| \bar{\mathbf{u}}^{(r+1)} - \bar{\mathbf{d}}^{(r+1)} \right\|^2 \\
&\leq \frac{1}{1-\beta} \sum_{k=0}^r \beta^{r-k} \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left\| \nabla f_i(\mathbf{x}_i^{(k)}) - \nabla f_i(\bar{\mathbf{x}}^{(k)}) \right\|^2 + \frac{\sigma^2}{N(1-\beta)^2} \\
&\stackrel{(5)}{\leq} \frac{L^2}{1-\beta} \sum_{k=0}^r \beta^{r-k} \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left\| \mathbf{x}_i^{(k)} - \bar{\mathbf{x}}^{(k)} \right\|^2 + \frac{\sigma^2}{N(1-\beta)^2}.
\end{aligned}$$

This concludes the proof.  $\square$

**Lemma 12.** Suppose that Assumptions 1, 2, 3, and 4 hold. For any round  $r \geq 0$ , it holds that

$$\mathbb{E} \left\| \bar{\mathbf{u}}^{(r+1)} \right\|^2 \leq \frac{2L^2}{1-\beta} \sum_{k=0}^r \beta^{r-k} \Xi^{(k)} + \frac{2}{1-\beta} \sum_{k=0}^r \beta^{r-k} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 + \frac{2\sigma^2}{N(1-\beta)^2}.$$

*Proof.* We have

$$\begin{aligned}
\mathbb{E} \left\| \bar{\mathbf{u}}^{(r+1)} \right\|^2 &= \mathbb{E} \left\| \bar{\mathbf{u}}^{(r+1)} \pm \bar{\mathbf{d}}^{(r+1)} \right\|^2 \\
&\stackrel{(13)}{\leq} 2\mathbb{E} \left\| \bar{\mathbf{u}}^{(r+1)} - \bar{\mathbf{d}}^{(r+1)} \right\|^2 + 2\mathbb{E} \left\| \bar{\mathbf{d}}^{(r+1)} \right\|^2.
\end{aligned}$$

From Lemmas 9 and 11, we obtain the statement.  $\square$

## D.5 MAIN PROOF

**Lemma 13** (Descent Lemma). *Suppose that Assumptions 1, 2, 3, and 4 hold. If the step size  $\eta$  satisfies*

$$\eta \leq \frac{1-\beta}{4L},$$

*then it holds that for any round  $r \geq 0$ ,*

$$\begin{aligned} \mathbb{E}f(\bar{\mathbf{z}}^{(r+1)}) &\leq \mathbb{E}f(\bar{\mathbf{z}}^{(r)}) + \frac{L^2\eta}{1-\beta} \mathbb{E} \left\| \bar{\mathbf{x}}^{(r)} - \bar{\mathbf{z}}^{(r)} \right\|^2 + \frac{L^2\eta}{1-\beta} \Xi^{(r)} \\ &\quad - \frac{\eta}{4(1-\beta)} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 - \frac{\eta}{4(1-\beta)} \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 + \frac{L\sigma^2\eta^2}{2N(1-\beta)^2}. \end{aligned}$$

*Proof.* From Assumption 3 and Lemma 6, we have

$$\begin{aligned} &\mathbb{E}_r f(\bar{\mathbf{z}}^{(r+1)}) \\ &\leq f(\bar{\mathbf{z}}^{(r)}) + \mathbb{E}_r \langle \nabla f(\bar{\mathbf{z}}^{(r)}), \bar{\mathbf{z}}^{(r+1)} - \bar{\mathbf{z}}^{(r)} \rangle + \frac{L}{2} \mathbb{E}_r \left\| \bar{\mathbf{z}}^{(r+1)} - \bar{\mathbf{z}}^{(r)} \right\|^2 \\ &= f(\bar{\mathbf{z}}^{(r)}) - \frac{\eta}{1-\beta} \left\langle \nabla f(\bar{\mathbf{z}}^{(r)}), \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\rangle + \frac{L\eta^2}{2(1-\beta)^2} \mathbb{E}_r \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}; \xi_i^{(r)}) \right\|^2 \\ &\stackrel{(6)}{\leq} f(\bar{\mathbf{z}}^{(r)}) - \frac{\eta}{1-\beta} \left\langle \nabla f(\bar{\mathbf{z}}^{(r)}), \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\rangle \\ &\quad + \frac{L\eta^2}{2(1-\beta)^2} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 + \frac{L\sigma^2\eta^2}{2N(1-\beta)^2} \\ &= f(\bar{\mathbf{z}}^{(r)}) + \underbrace{\frac{\eta}{1-\beta} \left\langle \nabla f(\bar{\mathbf{x}}^{(r)}) - \nabla f(\bar{\mathbf{z}}^{(r)}), \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\rangle}_{T_1} \\ &\quad - \underbrace{\frac{\eta}{1-\beta} \left\langle \nabla f(\bar{\mathbf{x}}^{(r)}), \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\rangle}_{T_2} + \underbrace{\frac{L\eta^2}{2(1-\beta)^2} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2}_{T_3} + \frac{L\sigma^2\eta^2}{2N(1-\beta)^2}. \end{aligned}$$

We can bound  $T_1$  from above as follows:

$$\begin{aligned} T_1 &\stackrel{(15), \gamma=2}{\leq} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) - \nabla f(\bar{\mathbf{z}}^{(r)}) \right\|^2 + \frac{1}{4} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 \\ &\stackrel{(5)}{\leq} L^2 \left\| \bar{\mathbf{x}}^{(r)} - \bar{\mathbf{z}}^{(r)} \right\|^2 + \frac{1}{4} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2. \end{aligned}$$

We can bound  $-T_2$  from above as follows:

$$\begin{aligned} -T_2 &= \frac{1}{2} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) - \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 - \frac{1}{2} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 - \frac{1}{2} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 \\ &\stackrel{(14)}{\leq} \frac{1}{2} \frac{1}{N} \sum_{i=1}^N \left\| \nabla f_i(\bar{\mathbf{x}}^{(r)}) - \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 - \frac{1}{2} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 - \frac{1}{2} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 \\ &\stackrel{(5)}{\leq} \frac{L^2}{2} \frac{1}{N} \sum_{i=1}^N \left\| \bar{\mathbf{x}}^{(r)} - \mathbf{x}_i^{(r)} \right\|^2 - \frac{1}{2} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 - \frac{1}{2} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2. \end{aligned}$$

Then, we can bound  $T_3$  from above as follows:

$$\begin{aligned}
T_3 &= \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \pm \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\
&\stackrel{(13)}{\leq} 2 \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) - \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 + 2 \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\
&\stackrel{(14)}{\leq} \frac{2}{N} \sum_{i=1}^N \left\| \nabla f_i(\mathbf{x}_i^{(r)}) - \nabla f_i(\bar{\mathbf{x}}^{(r)}) \right\|^2 + 2 \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\
&\stackrel{(5)}{\leq} \frac{2L^2}{N} \sum_{i=1}^N \left\| \mathbf{x}_i^{(r)} - \bar{\mathbf{x}}^{(r)} \right\|^2 + 2 \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2.
\end{aligned}$$

By combining them, we obtain

$$\begin{aligned}
&\mathbb{E}_r f(\bar{\mathbf{z}}^{(r+1)}) \\
&\leq f(\bar{\mathbf{z}}^{(r)}) + \frac{L^2\eta}{1-\beta} \left\| \bar{\mathbf{x}}^{(r)} - \bar{\mathbf{z}}^{(r)} \right\|^2 - \frac{\eta}{4(1-\beta)} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 \\
&\quad + \frac{L^2}{1-\beta} \left( \frac{1}{2} + \frac{L\eta}{1-\beta} \right) \eta \Xi^{(r)} - \frac{1}{1-\beta} \left( \frac{1}{2} - \frac{L\eta}{1-\beta} \right) \eta \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 + \frac{L\sigma^2\eta^2}{2N(1-\beta)^2}.
\end{aligned}$$

Using  $\eta \leq \frac{1-\beta}{4L}$ , we get the statement.  $\square$

**Lemma 14** (Recursion for  $\Xi$ ). *Suppose that Assumptions 1, 2, 3, and 4 hold. Then, it holds that for any round  $r \geq 0$ ,*

$$\Xi^{(r+1)} \leq (1 - \frac{p}{2})\Xi^{(r)} + \frac{9}{p}\eta^2\mathcal{E}^{(r)} + \frac{9}{Np}\eta^2\mathbb{E} \left\| \mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)} \right\|_F^2 + \frac{9}{Np}\eta^2\mathbb{E} \left\| \mathbf{E}^{(r+1)} \right\|_F^2.$$

*Proof.* Because  $\sum_{i=1}^N \|\mathbf{a}_i - \bar{\mathbf{a}}\|^2 \leq \sum_{i=1}^N \|\mathbf{a}_i\|^2$  for any  $\mathbf{a}_1, \dots, \mathbf{a}_N \in \mathbb{R}^d$ , we have

$$N\Xi^{(r)} = \mathbb{E} \left\| (\mathbf{X}^{(r)} - \bar{\mathbf{X}}^{(r-1)}) + (\bar{\mathbf{X}}^{(r-1)} - \bar{\mathbf{X}}^{(r)}) \right\|_F^2 \leq \mathbb{E} \left\| \mathbf{X}^{(r)} - \bar{\mathbf{X}}^{(r-1)} \right\|_F^2.$$

Then, we have

$$\begin{aligned}
\left\| \mathbf{X}^{(r+1)} - \bar{\mathbf{X}}^{(r+1)} \right\|_F^2 &\leq \left\| \mathbf{X}^{(r+1)} - \bar{\mathbf{X}}^{(r)} \right\|_F^2 \\
&= \left\| \mathbf{X}^{(r)}\mathbf{W} - \eta(\mathbf{U}^{(r+1)} - \mathbf{C}^{(r)}) - \bar{\mathbf{X}}^{(r)} \right\|_F^2 \\
&\stackrel{(13)}{\leq} (1+\gamma) \left\| \mathbf{X}^{(r)}\mathbf{W} - \bar{\mathbf{X}}^{(r)} \right\|_F^2 + (1+\gamma^{-1})\eta^2 \left\| \mathbf{U}^{(r+1)} - \mathbf{C}^{(r)} \right\|_F^2 \\
&\stackrel{(4)}{\leq} (1+\gamma)(1-p) \left\| \mathbf{X}^{(r)} - \bar{\mathbf{X}}^{(r)} \right\|_F^2 + (1+\gamma^{-1})\eta^2 \left\| \mathbf{U}^{(r+1)} - \mathbf{C}^{(r)} \right\|_F^2.
\end{aligned}$$

By substituting  $\gamma = \frac{p}{2}$  and using  $p \leq 1$ , we obtain

$$\begin{aligned}
&\left\| \mathbf{X}^{(r+1)} - \bar{\mathbf{X}}^{(r+1)} \right\|_F^2 \\
&\leq (1 - \frac{p}{2}) \left\| \mathbf{X}^{(r)} - \bar{\mathbf{X}}^{(r)} \right\|_F^2 + \frac{3}{p}\eta^2 \left\| \mathbf{U}^{(r+1)} - \mathbf{C}^{(r)} \right\|_F^2 \\
&= (1 - \frac{p}{2}) \left\| \mathbf{X}^{(r)} - \bar{\mathbf{X}}^{(r)} \right\|_F^2 + \frac{3}{p}\eta^2 \left\| \mathbf{U}^{(r+1)} \pm \mathbf{D}^{(r+1)} \pm \mathbf{E}^{(r+1)} - \mathbf{C}^{(r)} \right\|_F^2 \\
&\stackrel{(14)}{\leq} (1 - \frac{p}{2}) \left\| \mathbf{X}^{(r)} - \bar{\mathbf{X}}^{(r)} \right\|_F^2 \\
&\quad + \frac{9}{p}\eta^2 \left\| \mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)} \right\|_F^2 + \frac{9}{p}\eta^2 \left\| \mathbf{D}^{(r+1)} - \mathbf{C}^{(r)} - \mathbf{E}^{(r+1)} \right\|_F^2 + \frac{9}{p}\eta^2 \left\| \mathbf{E}^{(r+1)} \right\|_F^2.
\end{aligned}$$

This concludes the proof.  $\square$

**Lemma 15** (Recursion for  $\mathcal{E}$ ). *Suppose that Assumptions 1, 2, 3, and 4 hold. Then, it holds that for any round  $r \geq 0$ ,*

$$\begin{aligned} \mathcal{E}^{(r+1)} &\leq (1 - \frac{p}{2})\mathcal{E}^{(r)} + \frac{18\beta^2}{p}\mathcal{D}^{(r)} + \frac{24}{Np}\mathbb{E}\left\|\mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)}\right\|_F^2 \\ &\quad + \frac{144L^4}{p}\eta^2\Xi^{(r)} + \frac{72\beta^2L^2}{p}\eta^2\mathbb{E}\left\|\bar{\mathbf{u}}^{(r)}\right\|^2 + \frac{144L^2}{p}\eta^2\mathbb{E}\left\|\nabla f(\bar{\mathbf{x}}^{(r)})\right\|^2 + \frac{36L^2\sigma^2\eta^2}{Np}. \end{aligned}$$

*Proof.* We have

$$\begin{aligned} &\mathbb{E}\left\|\mathbf{D}^{(r+2)} - \mathbf{C}^{(r+1)} - \mathbf{E}^{(r+2)}\right\|_F^2 \\ &= \mathbb{E}\left\|\mathbf{D}^{(r+2)} - (\mathbf{C}^{(r)} - \mathbf{U}^{(r+1)})\mathbf{W} - \mathbf{U}^{(r+1)} - \mathbf{E}^{(r+2)} \pm \mathbf{D}^{(r+1)} \pm \mathbf{D}^{(r+1)}\mathbf{W} \pm \mathbf{E}^{(r+1)}\right\|_F^2 \\ &\stackrel{(13),(14)}{\leq} (1 + \gamma)\mathbb{E}\left\|(\mathbf{D}^{(r+1)} - \mathbf{C}^{(r)})\mathbf{W} - \mathbf{E}^{(r+1)}\right\|_F^2 \\ &\quad + 2(1 + \gamma^{-1})\mathbb{E}\left\|(\mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)})(\mathbf{W} - \mathbf{I})\right\|_F^2 \\ &\quad + 2(1 + \gamma^{-1})\mathbb{E}\left\|\mathbf{D}^{(r+2)} - \mathbf{D}^{(r+1)} + \mathbf{E}^{(r+1)} - \mathbf{E}^{(r+2)}\right\|_F^2 \\ &\stackrel{(4)}{\leq} (1 + \gamma)(1 - p)\mathbb{E}\left\|\mathbf{D}^{(r+1)} - \mathbf{C}^{(r)} - \mathbf{E}^{(r+1)}\right\|_F^2 \\ &\quad + 2(1 + \gamma^{-1})\mathbb{E}\left\|(\mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)})(\mathbf{W} - \mathbf{I})\right\|_F^2 \\ &\quad + 2(1 + \gamma^{-1})\mathbb{E}\left\|\mathbf{D}^{(r+2)} - \mathbf{D}^{(r+1)} + \mathbf{E}^{(r+1)} - \mathbf{E}^{(r+2)}\right\|_F^2, \end{aligned}$$

where we use Lemma 4 and  $\mathbf{E}^{(r+1)} = \frac{1}{N}\mathbf{D}^{(r+1)}\mathbf{1}\mathbf{1}^\top$  in the last inequality. Then, we have

$$\begin{aligned} &\mathbb{E}\left\|\mathbf{D}^{(r+2)} - \mathbf{C}^{(r+1)} - \mathbf{E}^{(r+2)}\right\|_F^2 \\ &\stackrel{(a)}{\leq} (1 + \gamma)(1 - p)\mathbb{E}\left\|\mathbf{D}^{(r+1)} - \mathbf{C}^{(r)} - \mathbf{E}^{(r+1)}\right\|_F^2 \\ &\quad + 2(1 + \gamma^{-1})\mathbb{E}\left\|\mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)}\right\|_F^2 \|\mathbf{W} - \mathbf{I}\|_{\text{op}}^2 \\ &\quad + 2(1 + \gamma^{-1})\mathbb{E}\left\|\mathbf{D}^{(r+2)} - \mathbf{D}^{(r+1)} + \mathbf{E}^{(r+1)} - \mathbf{E}^{(r+2)}\right\|_F^2 \\ &\stackrel{(b)}{\leq} (1 + \gamma)(1 - p)\mathbb{E}\left\|\mathbf{D}^{(r+1)} - \mathbf{C}^{(r)} - \mathbf{E}^{(r+1)}\right\|_F^2 + 8(1 + \gamma^{-1})\mathbb{E}\left\|\mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)}\right\|_F^2 \\ &\quad + 2(1 + \gamma^{-1})\mathbb{E}\left\|\mathbf{D}^{(r+2)} - \mathbf{D}^{(r+1)} + \mathbf{E}^{(r+1)} - \mathbf{E}^{(r+2)}\right\|_F^2, \end{aligned}$$

where  $\|\cdot\|_{\text{op}}$  denotes the operator norm. In (a), we use the following definition of the operator norm:  $\|\mathbf{W} - \mathbf{I}\|_{\text{op}} := \sup_{\hat{\mathbf{v}} \in \mathbb{R}^d \setminus \{\mathbf{0}\}} \frac{\|(\mathbf{W} - \mathbf{I})\hat{\mathbf{v}}\|}{\|\hat{\mathbf{v}}\|} \geq \frac{\|(\mathbf{W} - \mathbf{I})\mathbf{v}\|}{\|\mathbf{v}\|}$  for any  $\mathbf{v} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ . In (b), we use Gershgorin circle theorem and the fact that  $\mathbf{W}$  is a mixing matrix. Substituting  $\gamma = \frac{p}{2}$ , we obtain

$$\begin{aligned} &\mathbb{E}\left\|\mathbf{D}^{(r+2)} - \mathbf{C}^{(r+1)} - \mathbf{E}^{(r+2)}\right\|_F^2 \\ &\leq (1 - \frac{p}{2})\mathbb{E}\left\|\mathbf{D}^{(r+1)} - \mathbf{C}^{(r)} - \mathbf{E}^{(r+1)}\right\|_F^2 + \frac{24}{p}\mathbb{E}\left\|\mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)}\right\|_F^2 \\ &\quad + \frac{6}{p}\underbrace{\mathbb{E}\left\|\mathbf{D}^{(r+2)} - \mathbf{D}^{(r+1)} + \mathbf{E}^{(r+1)} - \mathbf{E}^{(r+2)}\right\|_F^2}_T. \end{aligned}$$

Then, we can bound  $T$  from above by expanding  $\mathbf{D}^{(r+2)}$ ,  $\mathbf{D}^{(r+1)}$ ,  $\mathbf{E}^{(r+2)}$ , and  $\mathbf{E}^{(r+1)}$  as follows:

$$\begin{aligned} T &\stackrel{(14)}{\leq} 3\beta^2 \mathbb{E} \left\| \mathbf{D}^{(r+1)} - \mathbf{D}^{(r)} + \mathbf{E}^{(r)} - \mathbf{E}^{(r+1)} \right\|_F^2 + 3\mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}^{(r+1)}) - \nabla f(\bar{\mathbf{X}}^{(r)}) \right\|_F^2 \\ &\quad + 3\mathbb{E} \left\| \frac{1}{N} \nabla f(\bar{\mathbf{X}}^{(r)}) \mathbf{1}\mathbf{1}^\top - \frac{1}{N} \nabla f(\bar{\mathbf{X}}^{(r+1)}) \mathbf{1}\mathbf{1}^\top \right\|_F^2 \\ &\stackrel{(5)}{\leq} 3\beta^2 \mathbb{E} \left\| \mathbf{D}^{(r+1)} - \mathbf{D}^{(r)} + \mathbf{E}^{(r)} - \mathbf{E}^{(r+1)} \right\|_F^2 + 6L^2 \mathbb{E} \left\| \bar{\mathbf{X}}^{(r+1)} - \bar{\mathbf{X}}^{(r)} \right\|_F^2. \end{aligned}$$

Using Lemma 8, we obtain

$$\begin{aligned} T &\leq 3\beta^2 \mathbb{E} \left\| \mathbf{D}^{(r+1)} - \mathbf{D}^{(r)} + \mathbf{E}^{(r)} - \mathbf{E}^{(r+1)} \right\|_F^2 \\ &\quad + N(24L^4\eta^2\Xi^{(r)} + 12\beta^2L^2\eta^2\mathbb{E} \left\| \bar{\mathbf{u}}^{(r)} \right\|^2 + 24L^2\eta^2\mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 + \frac{6L^2\sigma^2\eta^2}{N}). \end{aligned}$$

This concludes the proof.  $\square$

**Lemma 16** (Recursion for  $\mathcal{D}$ ). *Suppose that Assumptions 1, 2, 3, and 4 hold. Then, it holds that for any round  $r \geq 0$ ,*

$$\begin{aligned} \mathcal{D}^{(r+1)} &\leq \frac{2\beta^2}{1+\beta^2} \mathcal{D}^{(r)} + \frac{32L^4\eta^2}{1-\beta^2} \Xi^{(r)} + \frac{16L^2\beta^2\eta^2}{1-\beta^2} \mathbb{E} \left\| \bar{\mathbf{u}}^{(r)} \right\|^2 \\ &\quad + \frac{32L^2\eta^2}{1-\beta^2} \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 + \frac{8L^2\sigma^2\eta^2}{N(1-\beta^2)}. \end{aligned}$$

*Proof.* We have

$$\begin{aligned} &\mathbb{E} \left\| \mathbf{D}^{(r+2)} - \mathbf{D}^{(r+1)} - \mathbf{E}^{(r+2)} + \mathbf{E}^{(r+1)} \right\|_F^2 \\ &\stackrel{(13),(14)}{\leq} (1+\gamma)\beta^2 \mathbb{E} \left\| \mathbf{D}^{(r+1)} - \mathbf{D}^{(r)} + \mathbf{E}^{(r)} - \mathbf{E}^{(r+1)} \right\|_F^2 \\ &\quad + 2(1+\gamma^{-1})\mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}^{(r+1)}) - \nabla f(\bar{\mathbf{X}}^{(r)}) \right\|_F^2 \\ &\quad + 2(1+\gamma^{-1})\mathbb{E} \left\| \frac{1}{N} \nabla f(\bar{\mathbf{X}}^{(r)}) \mathbf{1}\mathbf{1}^\top - \frac{1}{N} \nabla f(\bar{\mathbf{X}}^{(r+1)}) \mathbf{1}\mathbf{1}^\top \right\|_F^2 \\ &\stackrel{(5)}{\leq} (1+\gamma)\beta^2 \mathbb{E} \left\| \mathbf{D}^{(r+1)} - \mathbf{D}^{(r)} + \mathbf{E}^{(r)} - \mathbf{E}^{(r+1)} \right\|_F^2 + 4L^2(1+\gamma^{-1})\mathbb{E} \left\| \bar{\mathbf{X}}^{(r+1)} - \bar{\mathbf{X}}^{(r)} \right\|_F^2. \end{aligned}$$

Substituting  $\gamma = \frac{1-\beta^2}{1+\beta^2}$ , we obtain

$$\begin{aligned} &\mathbb{E} \left\| \mathbf{D}^{(r+2)} - \mathbf{D}^{(r+1)} - \mathbf{E}^{(r+2)} + \mathbf{E}^{(r+1)} \right\|_F^2 \\ &\leq \frac{2\beta^2}{1+\beta^2} \mathbb{E} \left\| \mathbf{D}^{(r+1)} - \mathbf{D}^{(r)} + \mathbf{E}^{(r)} - \mathbf{E}^{(r+1)} \right\|_F^2 + \frac{8L^2}{1-\beta^2} \mathbb{E} \left\| \bar{\mathbf{X}}^{(r+1)} - \bar{\mathbf{X}}^{(r)} \right\|_F^2. \end{aligned}$$

Using Lemma 8, we obtain

$$\begin{aligned} &\mathbb{E} \left\| \mathbf{D}^{(r+2)} - \mathbf{D}^{(r+1)} - \mathbf{E}^{(r+2)} + \mathbf{E}^{(r+1)} \right\|_F^2 \\ &\leq \frac{2\beta^2}{1+\beta^2} \mathbb{E} \left\| \mathbf{D}^{(r+1)} - \mathbf{D}^{(r)} + \mathbf{E}^{(r)} - \mathbf{E}^{(r+1)} \right\|_F^2 \\ &\quad + N \left( \frac{32L^4\eta^2}{1-\beta^2} \Xi^{(r)} + \frac{16L^2\beta^2\eta^2}{1-\beta^2} \mathbb{E} \left\| \bar{\mathbf{u}}^{(r)} \right\|^2 + \frac{32L^2\eta^2}{1-\beta^2} \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 + \frac{8L^2\sigma^2\eta^2}{N(1-\beta^2)} \right). \end{aligned}$$

This concludes the proof.  $\square$



**Lemma 17** (Recursion for  $\Xi$ ,  $\mathcal{E}$ , and  $\mathcal{D}$ ). *We define  $t \in \mathbb{R}$  and  $A \in \mathbb{R}$  as follows:*

$$t := \frac{2\beta^2 p}{1 - \beta^2} + 4, \quad A := \frac{648}{1 - \frac{p}{t} - \frac{2\beta^2}{1 + \beta^2}}.$$

*Note that it holds that  $t \geq 4$  and  $A > 0$ . Suppose that Assumptions 1, 2, 3, and 4 hold, and step size  $\eta$  satisfies*

$$\eta \leq \frac{p}{8L\sqrt{324 + \frac{2A\beta^2}{1 - \beta^2}}}.$$

*Then, it holds that*

$$\begin{aligned} & \Xi^{(r+1)} + \frac{36}{p^2}\eta^2\mathcal{E}^{(r+1)} + \frac{A\beta^2}{p^3}\eta^2\mathcal{D}^{(r+1)} \\ & \leq (1 - \frac{p}{t})(\Xi^{(r)} + \frac{36}{p^2}\eta^2\mathcal{E}^{(r)} + \frac{A\beta^2}{p^3}\eta^2\mathcal{D}^{(r)}) \\ & \quad + \frac{1}{p}\eta^2\mathbb{E}\left\|\nabla f(\bar{\mathbf{x}}^{(r)})\right\|^2 + \frac{1}{Np}\left(9 + \frac{864}{p^2}\right)\eta^2\mathbb{E}\left\|\mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)}\right\|_F^2 \\ & \quad + \frac{L^2}{p^3}\left(2592\beta^2 + \frac{16A\beta^4}{1 - \beta^2}\right)\eta^4\mathbb{E}\left\|\bar{\mathbf{u}}^{(r)}\right\|^2 + \frac{9}{Np}\eta^2\mathbb{E}\left\|\mathbf{E}^{(r+1)}\right\|_F^2 + \frac{\sigma^2\eta^2}{p}. \end{aligned}$$

*Proof.* From Lemmas 14 and 15, we have

$$\begin{aligned} \Xi^{(r+1)} & \leq (1 - \frac{p}{2})\Xi^{(r)} + \frac{9}{p}\eta^2\mathcal{E}^{(r)} + \frac{9}{Np}\eta^2\mathbb{E}\left\|\mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)}\right\|_F^2 + \frac{9}{Np}\eta^2\mathbb{E}\left\|\mathbf{E}^{(r+1)}\right\|_F^2, \\ \frac{36}{p^2}\eta^2\mathcal{E}^{(r+1)} & \leq (1 - \frac{p}{2})\frac{36}{p^2}\eta^2\mathcal{E}^{(r)} + \frac{648\beta^2}{p^3}\eta^2\mathcal{D}^{(r)} + \frac{5184L^4}{p^3}\eta^4\Xi^{(r)} + \frac{2592\beta^2L^2}{p^3}\eta^4\mathbb{E}\left\|\bar{\mathbf{u}}^{(r)}\right\|^2 \\ & \quad + \frac{864}{Np^3}\eta^2\mathbb{E}\left\|\mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)}\right\|_F^2 + \frac{5184L^2}{p^3}\eta^4\mathbb{E}\left\|\nabla f(\bar{\mathbf{x}}^{(r)})\right\|^2 + \frac{1296L^2\sigma^2\eta^4}{Np^3}. \end{aligned}$$

Then, from Lemma 16, we have

$$\begin{aligned} \frac{A\beta^2}{p^3}\eta^2\mathcal{D}^{(r+1)} & \leq \frac{2A\beta^4}{(1 + \beta^2)p^3}\eta^2\mathcal{D}^{(r)} + \frac{32A\beta^2L^4}{(1 - \beta^2)p^3}\eta^4\Xi^{(r)} + \frac{16AL^2\beta^4}{(1 - \beta^2)p^3}\eta^4\mathbb{E}\left\|\bar{\mathbf{u}}^{(r)}\right\|^2 \\ & \quad + \frac{32A\beta^2L^2}{(1 - \beta^2)p^3}\eta^4\mathbb{E}\left\|\nabla f(\bar{\mathbf{x}}^{(r)})\right\|^2 + \frac{8A\beta^2L^2\sigma^2}{N(1 - \beta^2)p^3}\eta^4. \end{aligned}$$

Using  $\eta^2 \leq \frac{p^2}{L^2}$  and  $\eta^2 \leq \frac{p^2}{128L^2(162 + \frac{A\beta^2}{1 - \beta^2})}$ , we have

$$\left((1 - \frac{p}{2}) + \frac{5184L^4}{p^3}\eta^4 + \frac{32A\beta^2L^4}{(1 - \beta^2)p^3}\eta^4\right)\Xi^{(r)} \leq \left((1 - \frac{p}{2}) + \frac{L^2}{4p}\eta^2\right)\Xi^{(r)} \leq (1 - \frac{p}{4})\Xi^{(r)}.$$

In addition, we have

$$\begin{aligned} \left(\frac{9}{p}\eta^2 + (1 - \frac{p}{2})\frac{36}{p^2}\eta^2\right)\mathcal{E}^{(r)} & = (1 - \frac{p}{4})\frac{36}{p^2}\eta^2\mathcal{E}^{(r)}, \\ \left(\frac{648\beta^2}{p^3}\eta^2 + \frac{2A\beta^4}{(1 + \beta^2)p^3}\eta^2\right)\mathcal{D}^{(r)} & = \left(\frac{648}{A} + \frac{2\beta^2}{1 + \beta^2}\right)\frac{A\beta^2}{p^3}\eta^2\mathcal{D}^{(r)} = (1 - \frac{p}{t})\frac{A\beta^2}{p^3}\eta^2\mathcal{D}^{(r)}. \end{aligned}$$

Then, using  $t \geq 4$ , we obtain

$$\begin{aligned}
& \Xi^{(r+1)} + \frac{36}{p^2} \eta^2 \mathcal{E}^{(r+1)} + \frac{A\beta^2}{p^3} \eta^2 \mathcal{D}^{(r+1)} \\
& \leq (1 - \frac{p}{t}) (\Xi^{(r)} + \frac{36}{p^2} \eta^2 \mathcal{E}^{(r)} + \frac{A\beta^2}{p^3} \eta^2 \mathcal{D}^{(r)}) \\
& \quad + \frac{L^2}{p^3} \left( \frac{32A\beta^2}{1-\beta^2} + 5184 \right) \eta^4 \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\
& \quad + \frac{1}{Np} \left( 9 + \frac{864}{p^2} \right) \eta^2 \mathbb{E} \left\| \mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)} \right\|_F^2 \\
& \quad + \frac{9}{Np} \eta^2 \mathbb{E} \left\| \mathbf{E}^{(r+1)} \right\|_F^2 \\
& \quad + \frac{L^2}{p^3} \left( 2592\beta^2 + \frac{16A\beta^4}{1-\beta^2} \right) \eta^4 \mathbb{E} \left\| \bar{\mathbf{u}}^{(r)} \right\|^2 \\
& \quad + \frac{L^2\sigma^2}{Np^3} \left( 1296 + \frac{8A\beta^2}{1-\beta^2} \right) \eta^4.
\end{aligned}$$

Using  $\eta^2 \leq \frac{p^2}{32L^2(162 + \frac{A\beta^2}{1-\beta^2})}$ , we have

$$\frac{L^2}{p^3} \left( \frac{32A\beta^2}{1-\beta^2} + 5184 \right) \eta^4 \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \leq \frac{1}{p} \eta^2 \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2.$$

Using  $\eta^2 \leq \frac{p^2}{8L^2(162 + \frac{A\beta^2}{1-\beta^2})}$ , we obtain

$$\frac{L^2\sigma^2}{Np^3} \left( 1296 + \frac{8A\beta^2}{1-\beta^2} \right) \eta^4 \leq \frac{\sigma^2\eta^2}{Np} \leq \frac{\sigma^2\eta^2}{p}.$$

This concludes the proof.  $\square$

**Lemma 18.** We define  $t \in \mathbb{R}$  and  $A \in \mathbb{R}$  as follows:

$$t := \frac{2\beta^2 p}{1-\beta^2} + 4, \quad A := \frac{648}{1 - \frac{p}{t} - \frac{2\beta^2}{1+\beta^2}}.$$

Under the same assumptions as those in Lemma 17, it holds that

$$\sum_{r=0}^R \Xi^{(r)} \leq \frac{t}{p} \sum_{k=0}^R \Psi^{(k)} + \frac{145t\sigma^2\eta^2}{(1-\beta)^2 p^3} R.$$

where  $\Psi^{(r)}$  is defined as follows:

$$\begin{aligned}
\Psi^{(r)} &:= \frac{1}{p} \eta^2 \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 + \frac{1}{Np} \left( 9 + \frac{864}{p^2} \right) \eta^2 \mathbb{E} \left\| \mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)} \right\|_F^2 \\
&\quad + \frac{L^2}{p^3} \left( 2592\beta^2 + \frac{16A\beta^4}{1-\beta^2} \right) \eta^4 \mathbb{E} \left\| \bar{\mathbf{u}}^{(r)} \right\|^2 + \frac{9}{Np} \eta^2 \mathbb{E} \left\| \mathbf{E}^{(r+1)} \right\|_F^2.
\end{aligned}$$

*Proof.* We define  $\Theta^{(r)} := \Xi^{(r)} + \frac{36}{p^2} \eta^2 \mathcal{E}^{(r)} + \frac{A\beta^2}{p^3} \eta^2 \mathcal{D}^{(r)}$ . From Lemma 17, we obtain

$$\begin{aligned}
\Theta^{(r+1)} &\leq (1 - \frac{p}{t}) \Theta^{(r)} + \Psi^{(r)} + \frac{\sigma^2\eta^2}{p} \\
&\leq (1 - \frac{p}{t})^{r+1} \Theta^{(0)} + \sum_{k=0}^r (1 - \frac{p}{t})^{r-k} \Psi^{(k)} + \sum_{k=0}^r (1 - \frac{p}{t})^{r-k} \frac{\sigma^2\eta^2}{p}.
\end{aligned}$$

Using  $\sum_{k=0}^r (1 - \frac{p}{t})^{r-k} \leq \frac{t}{p}$ , we obtain

$$\Theta^{(r+1)} \leq (1 - \frac{p}{t})^{r+1} \Theta^{(0)} + \sum_{k=0}^r (1 - \frac{p}{t})^{r-k} \Psi^{(k)} + \frac{t\sigma^2\eta^2}{p^2}.$$

Then, for any  $R \geq 1$ , we obtain

$$\begin{aligned} \sum_{r=1}^R \Theta^{(r)} &\leq \sum_{r=1}^R (1 - \frac{p}{t})^r \Theta^{(0)} + \sum_{r=1}^R \sum_{k=0}^{r-1} (1 - \frac{p}{t})^{r-k-1} \Psi^{(k)} + \frac{t\sigma^2\eta^2}{p^2} R \\ &= \sum_{r=1}^R (1 - \frac{p}{t})^r \Theta^{(0)} + \sum_{k=0}^{R-1} \Psi^{(k)} \sum_{r=k+1}^R (1 - \frac{p}{t})^{r-k-1} + \frac{t\sigma^2\eta^2}{p^2} R \\ &\leq \frac{t}{p} \Theta^{(0)} + \frac{t}{p} \sum_{k=0}^{R-1} \Psi^{(k)} + \frac{t\sigma^2\eta^2}{p^2} R, \end{aligned}$$

where we use  $\sum_{r=1}^R (1 - \frac{p}{t})^r \leq \frac{t}{p}$  and  $\sum_{r=k+1}^R (1 - \frac{p}{t})^{r-k-1} \leq \frac{t}{p}$  in the last inequality. Then, from the definition of  $\mathcal{E}^{(r)}$ , we have

$$\begin{aligned} \mathcal{E}^{(0)} &= \frac{1}{N} \mathbb{E} \left\| \mathbf{D}^{(1)} - \mathbf{C}^{(0)} - \mathbf{E}^{(1)} \right\|_F^2 \\ &= \frac{1}{(1-\beta)^2} \frac{1}{N} \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}^{(0)}) - \frac{1}{N} \nabla f(\bar{\mathbf{X}}^{(0)}) \mathbf{1}\mathbf{1}^\top - \nabla F(\mathbf{X}^{(0)}; \xi^{(0)}) + \frac{1}{N} \nabla F(\mathbf{X}^{(0)}; \xi^{(0)}) \mathbf{1}\mathbf{1}^\top \right\|_F^2 \\ &\stackrel{(14)}{\leq} \frac{2}{(1-\beta)^2} \frac{1}{N} \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}^{(0)}) - \nabla F(\mathbf{X}^{(0)}; \xi^{(0)}) \right\|_F^2 \\ &\quad + \frac{2}{(1-\beta)^2} \frac{1}{N} \mathbb{E} \left\| \frac{1}{N} \nabla f(\bar{\mathbf{X}}^{(0)}) \mathbf{1}\mathbf{1}^\top - \frac{1}{N} \nabla F(\mathbf{X}^{(0)}; \xi^{(0)}) \mathbf{1}\mathbf{1}^\top \right\|_F^2 \\ &\stackrel{(6)}{\leq} \frac{4}{(1-\beta)^2} \sigma^2, \end{aligned}$$

where we use  $\mathbf{X}^{(0)} = \bar{\mathbf{X}}^{(0)}$  in the last inequality. From the definition of  $\mathcal{D}^{(r)}$ , we have

$$\mathcal{D}^{(0)} = \frac{1}{N} \mathbb{E} \left\| (\beta - 1) \mathbf{D}^{(0)} + \nabla f(\bar{\mathbf{X}}^{(0)}) - \frac{1}{N} \nabla f(\bar{\mathbf{X}}^{(0)}) \mathbf{1}\mathbf{1}^\top \right\|_F^2 = 0.$$

Then using  $\mathbf{X}^{(0)} = \bar{\mathbf{X}}^{(0)}$  (i.e.,  $\Xi^{(0)} = 0$ ), we have

$$\Theta^{(0)} \leq \frac{144\sigma^2\eta^2}{(1-\beta)^2 p^2}.$$

Here, the above upper bounds of  $\mathcal{E}^{(0)}$  and  $\mathcal{D}^{(0)}$  are attributed to how we choose the initial values  $\mathbf{u}_i^{(0)}$ ,  $\mathbf{c}_i^{(0)}$ ,  $\mathbf{d}_i^{(0)}$ , and  $\mathbf{e}_i^{(0)}$  for  $i \in V$ . Then, combining them, we obtain

$$\begin{aligned} \sum_{r=1}^R \Theta^{(r)} &\leq \frac{t}{p} \sum_{k=0}^{R-1} \Psi^{(k)} + \frac{144t\sigma^2\eta^2}{(1-\beta)^2 p^3} + \frac{t\sigma^2\eta^2}{p^2} R \\ &\leq \frac{t}{p} \sum_{k=0}^{R-1} \Psi^{(k)} + \frac{145t\sigma^2\eta^2}{(1-\beta)^2 p^3} R, \end{aligned}$$

where we use  $p \in (0, 1]$ ,  $\beta \in [0, 1)$ , and  $R \geq 1$  in the last inequality. Then, using  $\Theta^{(r)} \geq \Xi^{(r)}$  and  $\Xi^{(0)} = 0$ , we obtain the statement.  $\square$

**Lemma 19.** We define  $t \in \mathbb{R}$  and  $A \in \mathbb{R}$  as follows:

$$t := \frac{2\beta^2 p}{1 - \beta^2} + 4, \quad A := \frac{648}{1 - \frac{p}{t} - \frac{2\beta^2}{1+\beta^2}}.$$

Note that it holds that  $t \geq 4$  and  $A > 0$ . Suppose that the same assumptions as those in Lemma 17 hold. Then, if step size  $\eta$  satisfies

$$\eta \leq \min\left\{\frac{p}{4L\sqrt{324 + \frac{2A\beta^2}{1-\beta^2}}}, \frac{(1-\beta)p}{2L\sqrt{t(5 + \frac{432}{p^2})}}, \frac{(1-\beta)p}{8L\sqrt{5t}}\right\},$$

it holds that

$$4L^2 \sum_{r=0}^R \Xi^{(r)} \leq \frac{1}{2} \sum_{k=0}^R \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 + \frac{40L^2 t}{(1-\beta)^3 p^2} \left( 10 + \frac{29}{p} + \frac{864}{p^2} \right) \sigma^2 \eta^2 (R+1).$$

*Proof.* We define  $\Psi^{(r)}$  as follows:

$$\begin{aligned} \Psi^{(r)} := & \frac{1}{p} \eta^2 \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 + \frac{1}{Np} \left( 9 + \frac{864}{p^2} \right) \eta^2 \mathbb{E} \left\| \mathbf{U}^{(r+1)} - \mathbf{D}^{(r+1)} \right\|_F^2 \\ & + \frac{L^2}{p^3} \left( 2592\beta^2 + \frac{16A\beta^4}{1-\beta^2} \right) \eta^4 \mathbb{E} \left\| \bar{\mathbf{u}}^{(r)} \right\|^2 + \frac{9}{Np} \eta^2 \mathbb{E} \left\| \mathbf{E}^{(r+1)} \right\|_F^2. \end{aligned}$$

Using  $\bar{\mathbf{e}}^{(r)} = \mathbf{e}_i^{(r)}$  and Lemmas 9, 10, and 12, we obtain

$$\begin{aligned} \Psi^{(r)} \leq & \frac{1}{p} \eta^2 \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\ & + \eta^2 \frac{9}{(1-\beta)p} \sum_{k=0}^r \beta^{r-k} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \\ & + \frac{2L^2}{(1-\beta)p^3} \left( 2592\beta^2 + \frac{16A\beta^4}{1-\beta^2} \right) \eta^4 \left( \sum_{k=0}^{r-1} \beta^{r-k-1} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \right) \\ & + \frac{L^2}{(1-\beta)p} \left( 9 + \frac{864}{p^2} \right) \eta^2 \left( \sum_{k=0}^r \beta^{r-k} \Xi^{(k)} \right) \\ & + \frac{2L^4}{(1-\beta)p^3} \left( 2592\beta^2 + \frac{16A\beta^4}{1-\beta^2} \right) \eta^4 \left( \sum_{k=0}^{r-1} \beta^{r-k-1} \Xi^{(k)} \right) \\ & + \frac{5}{(1-\beta)^3 p} \left( 9 + \frac{864}{p^2} \right) \sigma^2 \eta^2 \\ & + \frac{2L^2}{N(1-\beta)^2 p^3} \left( 2592\beta^2 + \frac{16A\beta^4}{1-\beta^2} \right) \sigma^2 \eta^4, \end{aligned}$$

for any round  $r \geq 1$ . Then, we obtain

$$\begin{aligned}
\sum_{r=1}^R \Psi^{(r)} &\leq \frac{1}{p} \eta^2 \sum_{r=1}^R \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\
&\quad + \eta^2 \frac{9}{(1-\beta)p} \sum_{r=1}^R \sum_{k=0}^r \beta^{r-k} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \\
&\quad + \frac{2L^2}{(1-\beta)p^3} \left( 2592\beta^2 + \frac{16A\beta^4}{1-\beta^2} \right) \eta^4 \left( \sum_{r=1}^R \sum_{k=0}^{r-1} \beta^{r-k-1} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \right) \\
&\quad + \frac{L^2}{(1-\beta)p} \left( 9 + \frac{864}{p^2} \right) \eta^2 \left( \sum_{r=1}^R \sum_{k=0}^r \beta^{r-k} \Xi^{(k)} \right) \\
&\quad + \frac{2L^4}{(1-\beta)p^3} \left( 2592\beta^2 + \frac{16A\beta^4}{1-\beta^2} \right) \eta^4 \left( \sum_{r=1}^R \sum_{k=0}^{r-1} \beta^{r-k-1} \Xi^{(k)} \right) \\
&\quad + \frac{5}{(1-\beta)^3 p} \left( 9 + \frac{864}{p^2} \right) \sigma^2 \eta^2 R \\
&\quad + \frac{2L^2}{N(1-\beta)^2 p^3} \left( 2592\beta^2 + \frac{16A\beta^4}{1-\beta^2} \right) \sigma^2 \eta^4 R.
\end{aligned}$$

Then, we obtain

$$\begin{aligned}
\sum_{r=1}^R \Psi^{(r)} &\leq \frac{1}{p} \eta^2 \sum_{r=1}^R \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\
&\quad + \eta^2 \frac{9}{(1-\beta)p} \sum_{k=0}^R \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \sum_{r=\max\{1,k\}}^R \beta^{r-k} \\
&\quad + \frac{2L^2}{(1-\beta)p^3} \left( 2592\beta^2 + \frac{16A\beta^4}{1-\beta^2} \right) \eta^4 \left( \sum_{k=0}^{R-1} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \sum_{r=k+1}^R \beta^{r-k-1} \right) \\
&\quad + \frac{L^2}{(1-\beta)p} \left( 9 + \frac{864}{p^2} \right) \eta^2 \left( \sum_{k=0}^R \Xi^{(k)} \sum_{r=\max\{1,k\}}^R \beta^{r-k} \right) \\
&\quad + \frac{2L^4}{(1-\beta)p^3} \left( 2592\beta^2 + \frac{16A\beta^4}{1-\beta^2} \right) \eta^4 \left( \sum_{k=0}^{R-1} \Xi^{(k)} \sum_{r=k+1}^R \beta^{r-k-1} \right) \\
&\quad + \frac{5}{(1-\beta)^3 p} \left( 9 + \frac{864}{p^2} \right) \sigma^2 \eta^2 R \\
&\quad + \frac{2L^2}{N(1-\beta)^2 p^3} \left( 2592\beta^2 + \frac{16A\beta^4}{1-\beta^2} \right) \sigma^2 \eta^4 R.
\end{aligned}$$

Using  $\sum_{r=k+1}^R \beta^{r-k-1} \leq \frac{1}{1-\beta}$  and  $\sum_{r=\max\{1,k\}}^R \beta^{r-k} \leq \frac{1}{1-\beta}$ , we obtain

$$\begin{aligned}
\sum_{r=1}^R \Psi^{(r)} &\leq \frac{1}{p} \eta^2 \sum_{r=1}^R \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\
&\quad + \eta^2 \frac{9}{(1-\beta)^2 p} \sum_{k=0}^R \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \\
&\quad + \frac{2L^2 \beta^2}{(1-\beta)^2 p^3} \left( 2592 + \frac{16A\beta^2}{1-\beta^2} \right) \eta^4 \left( \sum_{k=0}^{R-1} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \right) \\
&\quad + \frac{L^2}{(1-\beta)^2 p} \left( 9 + \frac{864}{p^2} \right) \eta^2 \left( \sum_{k=0}^R \Xi^{(k)} \right) \\
&\quad + \frac{2L^4 \beta^2}{(1-\beta)^2 p^3} \left( 2592 + \frac{16A\beta^2}{1-\beta^2} \right) \eta^4 \left( \sum_{k=0}^{R-1} \Xi^{(k)} \right) \\
&\quad + \frac{5}{(1-\beta)^3 p} \left( 9 + \frac{864}{p^2} \right) \sigma^2 \eta^2 R \\
&\quad + \frac{2L^2 \beta^2}{N(1-\beta)^2 p^3} \left( 2592 + \frac{16A\beta^2}{1-\beta^2} \right) \sigma^2 \eta^4 R.
\end{aligned}$$

Then, using  $\bar{\mathbf{u}}^{(0)} = \mathbf{0}$  and Lemmas 9 and 10, we have

$$\Psi^{(0)} \leq \frac{1}{p} \eta^2 \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(0)}) \right\|^2 + \frac{9}{(1-\beta)p} \eta^2 \left\| \nabla f(\bar{\mathbf{x}}^{(0)}) \right\|^2 + \frac{5}{(1-\beta)^3 p} \left( 9 + \frac{864}{p^2} \right) \sigma^2 \eta^2.$$

Then, we obtain

$$\begin{aligned}
\sum_{r=0}^R \Psi^{(r)} &\leq \frac{1}{p} \eta^2 \sum_{r=0}^R \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\
&\quad + \eta^2 \frac{18}{(1-\beta)^2 p} \sum_{k=0}^R \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \\
&\quad + \frac{2L^2 \beta^2}{(1-\beta)^2 p^3} \left( 2592 + \frac{16A\beta^2}{1-\beta^2} \right) \eta^4 \left( \sum_{k=0}^{R-1} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \right) \\
&\quad + \frac{L^2}{(1-\beta)^2 p} \left( 9 + \frac{864}{p^2} \right) \eta^2 \left( \sum_{k=0}^R \Xi^{(k)} \right) \\
&\quad + \frac{2L^4 \beta^2}{(1-\beta)^2 p^3} \left( 2592 + \frac{16A\beta^2}{1-\beta^2} \right) \eta^4 \left( \sum_{k=0}^{R-1} \Xi^{(k)} \right) \\
&\quad + \frac{5}{(1-\beta)^3 p} \left( 9 + \frac{864}{p^2} \right) \sigma^2 \eta^2 (R+1) \\
&\quad + \frac{2L^2 \beta^2}{N(1-\beta)^2 p^3} \left( 2592 + \frac{16A\beta^2}{1-\beta^2} \right) \sigma^2 \eta^4 R.
\end{aligned}$$

Using  $\eta^2 \leq \frac{p^2}{32L^2(162 + \frac{A\beta^2}{1-\beta^2})}$ , we have

$$\frac{2L^2 \beta^2}{(1-\beta)^2 p^3} \left( 2592 + \frac{16A\beta^2}{1-\beta^2} \right) \eta^4 \left( \sum_{k=0}^{R-1} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \right) \leq \frac{\beta^2}{(1-\beta)^2 p} \eta^2 \left( \sum_{k=0}^{R-1} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \right).$$

Using  $\eta^2 \leq \frac{p^2}{32L^2(162 + \frac{A\beta^2}{1-\beta^2})}$ , we have

$$\frac{2L^4 \beta^2}{(1-\beta)^2 p^3} \left( 2592 + \frac{16A\beta^2}{1-\beta^2} \right) \eta^4 \left( \sum_{k=0}^{R-1} \Xi^{(k)} \right) \leq \frac{L^2 \beta^2}{(1-\beta)^2 p} \eta^2 \left( \sum_{k=0}^{R-1} \Xi^{(k)} \right).$$

Using  $\eta^2 \leq \frac{p^2}{32L^2(162 + \frac{A\beta^2}{1-\beta^2})}$ , we have

$$\frac{2L^2\beta^2}{N(1-\beta)^2p^3} \left( 2592 + \frac{16A\beta^2}{1-\beta^2} \right) \sigma^2\eta^4R \leq \frac{\beta^2}{N(1-\beta)^2p} \sigma^2\eta^2R.$$

Then, using  $\beta \in [0, 1)$  and  $N \geq 1$ , we obtain

$$\begin{aligned} \sum_{r=0}^R \Psi^{(r)} &\leq \frac{1}{p} \eta^2 \sum_{r=0}^R \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\ &\quad + \eta^2 \frac{19}{(1-\beta)^2p} \sum_{k=0}^R \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \\ &\quad + \frac{L^2}{(1-\beta)^2p} \left( 10 + \frac{864}{p^2} \right) \eta^2 \left( \sum_{k=0}^R \Xi^{(k)} \right) \\ &\quad + \frac{5}{(1-\beta)^3p} \left( 10 + \frac{864}{p^2} \right) \sigma^2\eta^2(R+1). \end{aligned}$$

Using  $\beta \in [0, 1)$  and Lemma 18, we obtain

$$\begin{aligned} \sum_{r=0}^R \Xi^{(r)} &\leq \frac{20t}{(1-\beta)^2p^2} \eta^2 \sum_{k=0}^R \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \\ &\quad + \frac{tL^2}{(1-\beta)^2p^2} \left( 10 + \frac{864}{p^2} \right) \eta^2 \left( \sum_{k=0}^R \Xi^{(k)} \right) \\ &\quad + \frac{5t}{(1-\beta)^3p^2} \left( 10 + \frac{864}{p^2} \right) \sigma^2\eta^2(R+1) + \frac{145t\sigma^2\eta^2}{(1-\beta)^2p^3} R \\ &\leq \frac{20t}{(1-\beta)^2p^2} \eta^2 \sum_{k=0}^R \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 \\ &\quad + \frac{tL^2}{(1-\beta)^2p^2} \left( 10 + \frac{864}{p^2} \right) \eta^2 \left( \sum_{k=0}^R \Xi^{(k)} \right) \\ &\quad + \frac{5t}{(1-\beta)^3p^2} \left( 10 + \frac{29}{p} + \frac{864}{p^2} \right) \sigma^2\eta^2(R+1). \end{aligned}$$

Then, using  $\eta^2 \leq \frac{(1-\beta)^2p^2}{4tL^2(5 + \frac{432}{p^2})}$ , we obtain

$$\frac{1}{2} \sum_{r=0}^R \Xi^{(r)} \leq \frac{20t}{(1-\beta)^2p^2} \eta^2 \sum_{k=0}^R \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 + \frac{5t}{(1-\beta)^3p^2} \left( 10 + \frac{29}{p} + \frac{864}{p^2} \right) \sigma^2\eta^2(R+1).$$

Multiplying  $8L^2$ , we obtain

$$4L^2 \sum_{r=0}^R \Xi^{(r)} \leq \frac{160L^2t}{(1-\beta)^2p^2} \eta^2 \sum_{k=0}^R \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) \right\|^2 + \frac{40L^2t}{(1-\beta)^3p^2} \left( 10 + \frac{29}{p} + \frac{864}{p^2} \right) \sigma^2\eta^2(R+1).$$

Using  $\eta^2 \leq \frac{(1-\beta)^2p^2}{320L^2t}$ , we obtain the statement.  $\square$

**Lemma 20.** We define  $t \in \mathbb{R}$  as follows:

$$t := \frac{2\beta^2p}{1-\beta^2} + 4.$$

Suppose that the assumptions of Lemma 19 hold. Then, if step size  $\eta$  satisfies

$$\eta \leq \frac{(1-\beta)^2}{2\sqrt{2}L},$$

it holds that

$$\begin{aligned} \frac{1}{2(R+1)} \sum_{r=0}^R \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 &\leq \frac{4(1-\beta)}{\eta(R+1)} \left( f(\bar{\mathbf{z}}^{(0)}) - f^* \right) + \frac{2L\sigma^2\eta}{N(1-\beta)} \\ &\quad + \frac{L^2}{(1-\beta)^3} \left( \frac{40t}{p^2} \left( 10 + \frac{29}{p} + \frac{864}{p^2} \right) + \frac{4\beta^2}{N(1-\beta)} \right) \sigma^2\eta^2. \end{aligned}$$

*Proof.* Using Lemma 13 and Assumption 1, we have

$$\begin{aligned} \sum_{r=0}^R \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 &\leq \frac{4(1-\beta)}{\eta} \left( f(\bar{\mathbf{z}}^{(0)}) - f^* \right) + 4L^2 \sum_{r=0}^R \mathbb{E} \left\| \bar{\mathbf{x}}^{(r)} - \bar{\mathbf{z}}^{(r)} \right\|^2 + 4L^2 \sum_{r=0}^R \Xi^{(r)} \\ &\quad - \sum_{r=0}^R \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 + \frac{2L\sigma^2\eta}{N(1-\beta)} (R+1). \end{aligned}$$

From Lemma 7, we have

$$4L^2 \sum_{r=0}^R \mathbb{E} \left\| \bar{\mathbf{x}}^{(r)} - \bar{\mathbf{z}}^{(r)} \right\|^2 \leq \frac{4L^2\beta^2\eta^2}{(1-\beta)^4} \sum_{r=0}^R \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 + \frac{4L^2\beta^2\sigma^2\eta^2}{N(1-\beta)^4} R.$$

Combining them yields

$$\begin{aligned} &\sum_{r=0}^R \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\ &\leq \frac{4(1-\beta)}{\eta} \left( f(\bar{\mathbf{z}}^{(0)}) - f^* \right) + 4L^2 \sum_{r=0}^R \Xi^{(r)} - \left( 1 - \frac{4L^2\beta^2\eta^2}{(1-\beta)^4} \right) \sum_{r=0}^R \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}_i^{(r)}) \right\|^2 \\ &\quad + \frac{2L\sigma^2\eta}{N(1-\beta)} (R+1) + \frac{4L^2\beta^2\sigma^2\eta^2}{N(1-\beta)^4} R. \end{aligned}$$

Using  $\eta^2 \leq \frac{(1-\beta)^4}{8L^2}$  and  $\beta < 1$ , we obtain

$$\begin{aligned} &\sum_{r=0}^R \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\ &\leq \frac{4(1-\beta)}{\eta} \left( f(\bar{\mathbf{z}}^{(0)}) - f^* \right) + 4L^2 \sum_{r=0}^R \Xi^{(r)} + \frac{2L\sigma^2\eta}{N(1-\beta)} (R+1) + \frac{4L^2\beta^2\sigma^2\eta^2}{N(1-\beta)^4} R. \end{aligned}$$

Using Lemma 19, we obtain

$$\begin{aligned} &\frac{1}{2} \sum_{r=0}^R \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\ &\leq \frac{4(1-\beta)}{\eta} \left( f(\bar{\mathbf{z}}^{(0)}) - f^* \right) + \frac{2L\sigma^2\eta}{N(1-\beta)} (R+1) \\ &\quad + \frac{L^2}{(1-\beta)^3} \left( \frac{40t}{p^2} \left( 10 + \frac{29}{p} + \frac{864}{p^2} \right) + \frac{4\beta^2}{N(1-\beta)} \right) \sigma^2\eta^2 (R+1). \end{aligned}$$

This concludes the proof.  $\square$

**Lemma 21.** We define  $t \in \mathbb{R}$  and  $A \in \mathbb{R}$  as follows:

$$t := \frac{2\beta^2 p}{1-\beta^2} + 4, \quad A := \frac{648}{1 - \frac{p}{t} - \frac{2\beta^2}{1+\beta^2}}.$$

Then, it holds that

$$\frac{(1-\beta)^2 p^2}{16L\sqrt{\frac{7836\beta^2}{(1-\beta^2)^3 p} + 282}} \leq \min \left\{ \frac{1-\beta}{4L}, \frac{p}{8L\sqrt{324 + \frac{2A\beta^2}{1-\beta^2}}}, \frac{(1-\beta)p^2}{2L\sqrt{t(5p^2 + 432)}} \frac{(1-\beta)p}{8L\sqrt{5t}}, \frac{(1-\beta)^2}{2\sqrt{2}L} \right\}.$$



*Proof.* Because  $\sqrt{\frac{7836\beta^2}{(1-\beta^2)^3p}} + 282 > 1$ ,  $p \in (0, 1]$ , and  $\beta \in [0, 1)$ , we have

$$\frac{(1-\beta)^2 p^2}{16L\sqrt{\frac{7836\beta^2}{(1-\beta^2)^3p}} + 282} \leq \min \left\{ \frac{1-\beta}{4L}, \frac{(1-\beta)^2}{2\sqrt{2}L} \right\}.$$

From  $p \leq 1$ , we have

$$t - 3 = \frac{2\beta^2 p}{1-\beta^2} + 1 \geq \frac{1+\beta^2}{1-\beta^2} p = \frac{p}{1-\frac{2\beta^2}{1+\beta^2}}.$$

Then, we obtain

$$1 - \frac{p}{t} - \frac{2\beta^2}{1+\beta^2} \geq \frac{p}{t-3} - \frac{p}{t} = \frac{3p}{t(t-3)} \geq \frac{3p}{t^2}.$$

Then, we obtain

$$A \leq \frac{216t^2}{p}.$$

Using the above inequality, we obtain

$$\frac{A\beta^2}{1-\beta^2} + 162 \leq \frac{216\beta^2 t^2}{p(1-\beta^2)} + 162.$$

From the definition of  $t$ , we obtain

$$\begin{aligned} \frac{A\beta^2}{1-\beta^2} + 30t + 162 &\leq \frac{216\beta^2 t^2}{p(1-\beta^2)} + 30t + 162 \\ &= \frac{216\beta^2}{p(1-\beta^2)} \left( \frac{2\beta^2 p}{1-\beta^2} + 4 \right)^2 + 30 \left( \frac{2\beta^2 p}{1-\beta^2} + 4 \right) + 162 \\ &= \frac{216\beta^2}{p(1-\beta^2)} \left( \frac{4\beta^4 p^2}{(1-\beta^2)^2} + \frac{16\beta^2 p}{1-\beta^2} + 16 \right) + 30 \left( \frac{2\beta^2 p}{1-\beta^2} + 4 \right) + 162 \\ &\leq \frac{7836\beta^2}{p(1-\beta^2)^3} + 282, \end{aligned}$$

where we use  $\beta \in [0, 1)$  and  $p \in (0, 1]$  in the last inequality. Then, we obtain

$$\begin{aligned} \frac{(1-\beta)^2 p^2}{16L\sqrt{\frac{7836\beta^2}{(1-\beta^2)^3p}} + 282} &\leq \frac{(1-\beta)^2 p^2}{16L\sqrt{\frac{A\beta^2}{1-\beta^2} + 30t + 162}} \\ &\leq \min \left\{ \frac{p}{8L\sqrt{324 + \frac{2A\beta^2}{1-\beta^2}}}, \frac{(1-\beta)p^2}{2L\sqrt{t(5p^2 + 432)}} \frac{(1-\beta)p}{8L\sqrt{5t}} \right\}. \end{aligned}$$

This concludes the proof.  $\square$

**Lemma 22** (Convergence Rate for Non-convex Case). *Suppose that Assumptions 1, 2, 3, and 4 hold. Then, for any  $R \geq 1$ , there exists a step size  $\eta$  such that it holds that*

$$\begin{aligned} &\frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\ &\leq \mathcal{O} \left( \sqrt{\frac{r_0 \sigma^2 L}{NR}} + \left( \frac{r_0^2 \sigma^2 L^2}{p^4 R^2 (1-\beta)} \left( 1 + \frac{p\beta^2}{1-\beta} \right) \right)^{\frac{1}{3}} + \frac{Lr_0}{(1-\beta)p^2 R} \sqrt{1 + \frac{\beta^2}{(1-\beta^2)^3 p}} \right), \end{aligned}$$

where  $r_0 := f(\bar{\mathbf{x}}^{(0)}) - f^*$ .

*Proof.* From Lemmas 20 and 21, if the step size  $\eta$  satisfies the following:

$$\eta \leq \frac{(1-\beta)^2 p^2}{16L\sqrt{\frac{7836\beta^2}{(1-\beta^2)^3 p} + 282}},$$

then we have

$$\begin{aligned} & \frac{1}{2(R+1)} \sum_{r=0}^R \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\ & \leq \frac{4}{\tilde{\eta}(R+1)} \left( f(\bar{\mathbf{z}}^{(0)}) - f^* \right) + \frac{2L\sigma^2\tilde{\eta}}{N} + \underbrace{\frac{L^2}{1-\beta} \left( \frac{40t}{p^2} \left( 10 + \frac{29}{p} + \frac{864}{p^2} \right) + \frac{4\beta^2}{N(1-\beta)} \right)}_T \sigma^2 \tilde{\eta}^2, \end{aligned}$$

where we define  $\tilde{\eta} := \frac{\eta}{1-\beta}$ . Then, we can bound  $T$  from above as follows:

$$\begin{aligned} T &= \frac{L^2}{1-\beta} \left( \frac{40}{p^2} \left( \frac{2\beta^2 p}{1-\beta^2} + 4 \right) \left( 10 + \frac{29}{p} + \frac{864}{p^2} \right) + \frac{4\beta^2}{N(1-\beta)} \right) \\ &\stackrel{p \in (0,1]}{\leq} \frac{L^2}{1-\beta} \left( \frac{36120}{p^4} \left( \frac{2\beta^2 p}{1-\beta^2} + 4 \right) + \frac{4\beta^2}{N(1-\beta)} \right) \\ &\stackrel{p \in (0,1], \beta \in [0,1]}{\leq} \frac{36120L^2}{(1-\beta)p^4} \left( \frac{3\beta^2 p}{1-\beta} + 4 \right). \end{aligned}$$

Then, we obtain

$$\begin{aligned} & \frac{1}{2(R+1)} \sum_{r=0}^R \mathbb{E} \left\| \nabla f(\bar{\mathbf{x}}^{(r)}) \right\|^2 \\ & \leq \frac{4}{\tilde{\eta}(R+1)} \left( f(\bar{\mathbf{z}}^{(0)}) - f^* \right) + \frac{2L\sigma^2\tilde{\eta}}{N} + \frac{36120L^2}{(1-\beta)p^4} \left( \frac{3\beta^2 p}{1-\beta} + 4 \right) \sigma^2 \tilde{\eta}^2. \end{aligned}$$

Using Lemma 17 in the previous work (Koloskova et al., 2020), we obtain the statement.  $\square$

## E HYPERPARAMETER SETTINGS

Tables 6, 7, 8, 9, and 10 list the hyperparameter settings for each dataset. We evaluated the performance of each comparison method for different step sizes and selected the step size that achieved the highest accuracy on the validation dataset.

Table 6: Experimental settings for FashionMNIST.

Neural network architecture	LeNet (LeCun et al., 1998)
Normalization	Group normalization (Wu & He, 2018)
Step size	{0.005, 0.001, 0.0005}
L2 penalty	0.001
Batch size	100
Data augmentation	RandomCrop
Total number of epochs	500

Table 7: Experimental settings for SVHN.

Neural network architecture	LeNet (LeCun et al., 1998)
Normalization	Group normalization (Wu & He, 2018)
Step size	{0.005, 0.001, 0.0005}
L2 penalty	0.001
Batch size	100
Data augmentation	RandomCrop
Total number of epochs	500

Table 8: Experimental settings for CIFAR-10.

Neural network architecture	LeNet (LeCun et al., 1998)
Normalization	Group normalization (Wu & He, 2018)
Step size	{0.005, 0.001, 0.0005}
L2 penalty	0.001
Batch size	100
Data augmentation	RandomCrop, RandomHorizontalFlip
Total number of epochs	500

Table 9: Experimental settings for CIFAR-10 with VGG-11.

Neural network architecture	VGG-11 (Simonyan & Zisserman, 2015)
Normalization	Group normalization (Wu & He, 2018)
Step size	{0.05, 0.01, 0.005}
Step size decay	/10 at epoch 500 and 750.
L2 penalty	0.001
Batch size	100
Data augmentation	RandomCrop, RandomHorizontalFlip, RandomErasing
Total number of epochs	1000

Table 10: Experimental settings for CIFAR-10 with ResNet-34.

Neural network architecture	ResNet-34 (He et al., 2016)
Normalization	Group normalization (Wu & He, 2018)
Step size	{0.05, 0.01, 0.005}
Step size decay	/10 at epoch 375 and 563.
L2 penalty	0.001
Batch size	100
Data augmentation	RandomCrop, RandomHorizontalFlip, RandomErasing
Total number of epochs	750