# FASTER ADAPTIVE MOMENTUM-BASED FEDERATED METHODS FOR DISTRIBUTED COMPOSITION OPTI MIZATION

Anonymous authors

Paper under double-blind review

# ABSTRACT

Federated learning is a popular distributed learning paradigm in machine learning. Meanwhile, composition optimization is an effective hierarchical learning model, which appears in many machine learning applications such as meta learning and robust learning. More recently, although a few federated composition optimization algorithms have been proposed, they still suffer from high sample and communication complexities. In the paper, thus, we propose a class of faster adaptive federated compositional optimization algorithms (i.e., MFCGD and AdaMFCGD) to solve the nonconvex distributed composition problems, which builds on the momentum-based variance reduced and local-SGD techniques. In particular, our adaptive algorithm (i.e., AdaMFCGD) uses a unified adaptive matrix to flexibly incorporate various adaptive learning rates. Moreover, we provide a solid theoretical analysis for our algorithms under non-i.i.d. setting, and prove our algorithms obtain a lower sample and communication complexities simultaneously than the existing federated composition optimization algorithms. Specifically, our algorithms obtain lower sample complexity of  $\tilde{O}(\epsilon^{-3})$  with lower communication complexity of  $O(\epsilon^{-2})$  in finding an  $\epsilon$ -stationary solution. We conduct numerical experiments on robust federated learning and distributed meta learning tasks to demonstrate the efficiency of our algorithms.

# 029 030 031

039 040

006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027

028

# 1 INTRODUCTION

Composition optimization is an effective hierarchical model in machine learning, which is widely used to many applications such as reinforcement learning Wang et al. (2017b); Huo et al. (2018), meta learning Wang et al. (2021), robust federated learning Huang et al. (2021a) and deep AUC maximization Yuan et al. (2022). In the paper, we study the following distributed composition optimization problem:

$$\min_{x \in \mathbb{R}^d} \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{\xi^m} \bigg[ f^m \Big( \mathbb{E}_{\zeta^m} \big[ g^m(x; \zeta^m) \big]; \xi^m \Big) \bigg], \tag{1}$$

041  
042 where 
$$F(x) := \frac{1}{M} \sum_{m=1}^{M} f^m(x)$$
 and  $f^m(x) = \mathbb{E}_{\xi^m} \left[ f^m \left( \mathbb{E}_{\zeta^m} [g^m(x; \zeta^m)]; \xi^m \right) \right]$ . Here  $y^m = \sum_{m \in \mathbb{N}} f^m(x; \zeta^m) = \sum_{m \in \mathbb{N}} f^m(x$ 

 $\begin{array}{ll} \mathbf{043} & g^m(x) = \mathbb{E}_{\zeta^m \sim \mathcal{S}^m} \left[ g^m(x; \zeta^m) \right] \text{ and } f^m(y^m) = \mathbb{E}_{\xi^m \sim \mathcal{D}^m} \left[ f^m(y^m; \xi^m) \right] \text{ for any } m \in [M] \text{ denote} \\ \text{ the inner and outer objective functions respectively in m-th client. Here } \xi^m \text{ and } \zeta^m \text{ for any } m \in [M] \text{ are independent random variables follow unknown distributions } \mathcal{D}^m \text{ and } \mathcal{S}^m \text{ respectively. For} \\ \text{ any } m, j \in [M] \text{ possibly } \mathcal{D}^m \neq \mathcal{D}^j, \ \mathcal{S}^m \neq \mathcal{S}^j \text{ and } \mathcal{D}^m \neq \mathcal{S}^j. \text{ Applications of Problem (1)} \\ \text{ involve many machine learning problems with a compositional structure, which include model- \\ \text{ agnostic meta learning Tutunov et al. (2020); Chen et al. (2020b); Wang et al. (2021), reinforcement \\ \text{ learning Wang et al. (2017b); Huo et al. (2018) and sparse additive models Wang et al. (2017a). In \\ \text{ the following, we give two specific applications that can be formulated as the distributed composition \\ \text{ Problem (1).} \end{array}$ 

**1). Task-Distributed Meta Learning.** Meta learning is to learn some properties in the optimal model to improve model performances with more experiences, i.e., learning to learn Andrychowicz et al. (2016). Model-Agnostic Meta Learning (MAML) Finn et al. (2017) is a class of popular meta

054 Table 1: Sample and Communication complexities comparison of the representative federated **compositional optimization** algorithms in finding an  $\epsilon$ -stationary point of the distributed composi-056 tion optimization problem (1), i.e.,  $\mathbb{E} \|\nabla F(x)\| \leq \epsilon$  or its equivalent variants. ALR denotes adaptive learning rate. 057

-	Algorithm	Reference	Sample Complexity	Communication Complexity	ALR
	ComFedL	Huang et al. (2021a)	$O(\epsilon^{-8})$	$O(\epsilon^{-4})$	
-	Local-MOML	Wang et al. (2021)	$O(\epsilon^{-5})$	$O(\epsilon^{-3})$	
-	FEDNEST	Tarzanagh et al. (2022)	$\tilde{O}(\epsilon^{-4})$	$\tilde{O}(\epsilon^{-4})$	
-	Local-SCGDM	Gao et al. (2022)	$O(\epsilon^{-4})$	$O(\epsilon^{-3})$	
-	MFCGD	Ours	$ ilde{O}(\epsilon^{-3})$	$ ilde{O}(\epsilon^{-2})$	
-	AdaMFCGD	Ours	$ ilde{O}(\epsilon^{-3})$	$ ilde{O}(\epsilon^{-2})$	$\checkmark$

learning methods, which is to find a common initialization that can adapt to a desired model for a set of new tasks after taking several gradient descent steps. In the paper, we consider a class of task-distributed MAMLs, where a set of tasks  $\{\mathcal{T}_m\}_{m=1}^M$  are drawn from a certain task distribution and each task is assigned in each client. Specifically, we solve the following task-distributed MAML problem:

$$\min_{x \in \mathbb{R}^d} \frac{1}{M} \sum_{m=1}^M f^m \left( x - \eta \nabla f^m(x) \right), \tag{2}$$

(3)

074 where  $f^m(x) = \mathbb{E}_{\xi^m \sim \mathcal{D}^m}[f(x;\xi^m)]$ , and random variable  $\xi^m$  follows the unknown distribution 075  $\mathcal{D}^m$ , and  $\eta > 0$  is a learning rate. Let  $f^m(y^m) = f^m(g^m(x))$  and  $y^m = g^m(x) = x - \eta \nabla f^m(x)$ , 076 the above problem (2) is a special case of the above composition problem (1).

077 2). Distributionally Robust Federated Learning. Federated learning (FL) McMahan et al. (2017); 078 Kairouz et al. (2019); Li et al. (2021a) is a distributed and privacy preserving machine learning 079 method to learn a global model collaboratively from decentralized data distributed over a network of 080 devices. To tackle the data heterogeneity from different devices, some robust FL algorithms Mohri 081 et al. (2019); Reisizadeh et al. (2020); Deng et al. (2020b) have been studied. In the paper, as in 082 Huang et al. (2021a), we consider solving the following distributed composition problem to reach 083 distributionally robust FL, defined as

084 085

058

067

068

069

070

071 072 073

087

089 090

091

where  $g^m(x) = \mathbb{E}[g^m(x;\xi^m)]$  denotes the loss function in the *m*-th client, and  $f(\cdot)$  is a monotonically increasing function. Clearly, the problem (3) is a special case of the above problem (1).

 $\min_{x \in \mathbb{R}^d} \frac{1}{M} \sum_{m=1}^M f\Big(\mathbb{E}\big[g^m(x;\xi^m)\big]\Big),$ 

Although recently many compositional gradient algorithms have been proposed to solve the composition problems, few distributed algorithms focus on solving the distributed composition opti-092 mization problems. More recently, Huang et al. (2021a); Wang et al. (2021); Gao et al. (2022); Tarzanagh et al. (2022) proposed some federated compositional gradient algorithms for the distributed stochastic composition problems. However, few adaptive algorithm focuses on the composition 094 optimization problems under the distributed setting. Meanwhile, these existing federated compo-095 sition optimization methods suffer from large sample and communication complexities (Please see 096 Table 1). Then there exists a natural question:

098 099

100

Could we develop faster and adaptive federated learning methods to solve the distributed composition optimization problem (1)?

In the paper, we provide an affirmative answer to the above question and propose a class of 102 faster momentum-based federated compositional gradient descent algorithms (i.e., MFCGD and 103 AdaMFCGD) to solve Problem (1), which build on the local Stochastic Gradient Descent (SGD) 104 and momentum-based variance reduced techniques to obtain lower sample and communication com-105 plexities simultaneously. Our main contributions are as follows:

106 107

(1) We propose a class of faster adaptive momentum-based federated compositional gradient descent algorithms (i.e., MFCGD and AdaMFCGD) to solve the nonconvex distributed composition problems, which build on the momentum-based variance reduced and local-SGD techniques. In particular, our adaptive algorithm (i.e., AdaMFCGD) uses a unified adaptive matrix to flexibly incorporate various adaptive learning rates.

- (2) We provide a solid convergence analysis framework for our algorithms under non-i.i.d. setting, and prove that our algorithms obtain simultaneously lower sample complexity of  $\tilde{O}(\epsilon^{-3})$  and lower communication complexity of  $\tilde{O}(\epsilon^{-2})$  than the existing federated composition methods for finding an  $\epsilon$ -stationary solution (Please see Table 1).
  - (3) Experimental results demonstrate efficiency of our algorithms on the task-distributed meta learning and robust federated learning tasks.

# 119 2 RELATED WORKS

108

110

111

116

117 118

In this section, we overview some representative composition optimization, federated optimization and adaptive optimization methods, respectively.

122 123 2.1 COMPOSITION OPTIMIZATION

Composition optimization has been widely applied to many applications such as reinforcement 124 learning Wang et al. (2017b), model-agnostic meta Learning Tutunov et al. (2020) and risk man-125 agement Huo et al. (2018). Recently, many compositional gradient-based methods have recently 126 been proposed to solve these composition optimization problems. For example, stochastic com-127 positional gradient methods Wang et al. (2017a;b); Ghadimi et al. (2020) have been proposed to 128 solve these problems. Subsequently, some variance-reduced compositional algorithms Huo et al. 129 (2018); Lin et al. (2018); Zhang & Xiao (2019) have been proposed for composition optimization. 130 Tutunov et al. (2020); Chen et al. (2020b) presented a class of momentum-based compositional 131 gradient methods for stochastic composition optimization. More recently, Jiang et al. (2022) pro-132 posed a class of efficient momentum-based variance reduced methods for non-convex stochastic composition optimization. Huang & Gao (2022) studied the stochastic composition optimization on 133 Riemannian manifolds. 134

For the distributed setting, Huang et al. (2021a) firstly studied federated learning algorithm for the general distributed composition optimization. Meanwhile, Wang et al. (2021) studied personalized federated learning algorithm based on the composition optimization. Subsequently, Gao et al. (2022); Tarzanagh et al. (2022) proposed some accelerated federated learning algorithms for the distributed composition optimization.

# 140 2.2 FEDERATED LEARNING

Federated Learning (FL) McMahan et al. (2017); Li et al. (2021a); Zhang et al. (2022) is a promising 142 distributed machine learning framework for collaboratively training the global model without shar-143 ing the local data to obtain the privacy-preserving learning solutions, and is widely used in many 144 applications such as healthcare informatics Xu et al. (2021) and automatic diagnosis of COVID-19 145 Yang et al. (2021). McMahan et al. (2017) first studied FL and proposed the FedAvg algorith-146 m for FL based on local-SGD algorithms Stich (2019), where each client conducts multiple steps 147 of SGD with its local data and then sends the learned model to the server for averaging. Subsequently, Li et al. (2019); Karimireddy et al. (2019); Deng & Mahdavi (2021) have studied the 148 convergence properties of the local-SGD and FedAvg algorithms or their variations. To acceler-149 ate the vanilla local-SGD and FedAvg algorithms, various accelerated FL algorithms Yuan & Ma 150 (2020); Karimireddy et al. (2020); Khanduri et al. (2021); Chen et al. (2020a) have been develope-151 d and studied. For example, Karimireddy et al. (2020) proposed a stochastic controlled averaging 152 algorithm for FL by adopting the variance-reduced technique of SARAH Nguyen et al. (2017)/SPI-153 DER Fang et al. (2018). Subsequently, Khanduri et al. (2021) proposed a faster federated algorithm 154 based on momentum-based variance reduced technique of STORM Cutkosky & Orabona (2019) and 155 ProxHSGD Tran-Dinh et al. (2022), which obtains lower sample and communication complexities 156 simultaneously.

To solve the data heterogeneity in FL, Mohri et al. (2019); Deng et al. (2020b) proposed some effective robust FL algorithms by learning the worst-case loss based on the minimax optimization problems. To further incorporate personalization in FL, some personalized federated learning models Fallah et al. (2020); Deng et al. (2020a); Li et al. (2021b) have been developed and studied. For example, Li et al. (2021b) proposed an effective and efficient personalized FL algorithm (i.e., Ditto) by learning a regularized local model for each client.

# 162 2.3 ADAPTIVE OPTIMIZATION METHODS

163 Adaptive optimization methods Duchi et al. (2011); Kingma & Ba (2014) are a class of efficient 164 optimization methods due to using adaptive learning rates in machine learning, and they have been 165 widely studied in machine learning community. For example, AdaGrad Duchi et al. (2011) is the first 166 adaptive gradient method. Adam Kingma & Ba (2014) is a popular variation of AdaGrad algorithm 167 based on the momentum technique, which is the default optimization algorithm for training largescale machine learning models. Meanwhile, some variants of Adam algorithm Reddi et al. (2019); 168 Chen et al. (2019) have been proposed to obtain a convergence guarantee under the nonconvex 169 setting. To further improve the performance of Adam algorithm, recently some new its variants 170 such as AdamW Loshchilov & Hutter (2018) have been developed. More recently, some accelerated 171 adaptive gradient methods Cutkosky & Orabona (2019); Huang et al. (2021b) have been proposed 172 based on the momentum-based variance reduced techniques. In parallel, some adaptive gradient 173 methods Reddi et al. (2020); Chen et al. (2020c) are proposed for distributed optimization. For 174 example, Reddi et al. (2020) proposed a class of adaptive federated algorithms for FL by using 175 adaptive learning rates at the server side. 176

# 177 NOTATIONS

179 Let [M] denote the set  $\{1, 2, \dots, M\}$ .  $\|\cdot\|$  denotes the  $\ell_2$  norm for vectors and Frobenius norm for 180 matrices.  $\langle x, y \rangle$  denotes the inner product of two vectors x and y. For vectors x and y,  $x^r$  (r > 0)181 denotes the element-wise power operation, x/y denotes the element-wise division and  $\max(x, y)$ 182 denotes the element-wise maximum.  $I_d$  denotes a d-dimensional identity matrix.  $A \succ 0$  denotes 183 that A is a positive definite matrix.  $a_t = O(b_t)$  denotes that  $a_t \le cb_t$  for some constant c > 0. The 184 notation  $\tilde{O}(\cdot)$  hides logarithmic terms.  $\prod_C [x] = \arg \min_{||w|| \le C} ||x - w||^2$  denote a projection onto the ball with radius C > 0.

186

178

187 188

203

# **3** FEDERATED COMPOSITIONAL GRADIENT DESCENT ALGORITHMS

In this section, we propose a class of faster momentum-based federated compositional gradient de-189 scent algorithms (i.e., MFCGD and AdaMFCGD) to solve the problem (1), which builds on the 190 local-SGD and momentum-based variance reduced techniques. Specifically, the local-SGD tech-191 nique reduce the communication complexity and the momentum-based variance reduced technique 192 reduce the sample complexity without relying on large batches. Meanwhile, our AdaMFCGD al-193 gorithm uses the unified adaptive matrix to flexibly incorporate various adaptive learning rates in 194 updating variables. Specifically, Algorithm 1 provides a procedure framework of our MFCGD and 195 AdaMFCGD algorithms. 196

In Algorithm 1, when mod(t,q) = 0 (i.e., **synchronization** step), the server receives the local variables  $\{x_t^m\}_{m=1}^M$  and local gradients  $\{w_t^m\}_{m=1}^M$  from the clients, and then averages them to obtain the averaged variables  $\{\bar{x}_t\}$  and averaged gradients  $\{\bar{w}_t\}$ . Based on these averaged gradients  $\{\bar{w}_t\}$ , we can generate some adaptive matrices  $\{A_t\}_{t\geq 1}$  (i.e., adaptive learning rates). Note that for our non-adaptive MFCGD algorithm, we only set  $A_t = I_d$  for all  $t \geq 1$  in Algorithm 1. Besides one example given at the line 6 of Algorithm 1, we can also generate many other adaptive matrices. For example, we can generate adaptive matrix  $A_t$  as the norm-type of Adam, defined as

$$_{t} = \vartheta_{t}a_{t-1} + (1 - \vartheta_{t})\|\bar{w}_{t}\|, \quad A_{t} = \operatorname{diag}(a_{t} + \rho), \tag{4}$$

where  $0 < \vartheta_t \le 1$ . Note that we can directly choose  $\alpha_t$ ,  $\beta_t$  or  $\varrho_t$  instead of  $\vartheta_t$  to reduce the number of tuning parameters in our algorithm. Next, based on these adaptive matrices, we can update the variable x in the server, then sent it to each client.

When  $mod(t, q) \neq 0$  (i.e., **asynchronization** step), the clients receive the updated variables  $\{\bar{x}_{t+1}\}$ and the generated adaptive matrices  $\{A_t\}$  from the server. Then the clients use the momentum-based variance reduced technique of STORM Cutkosky & Orabona (2019) and ProxHSGD Tran-Dinh et al. (2022) to update the stochastic gradients based on local data: for  $m \in [M]$ 

212 
$$h_{t+1}^m = g^m(x_{t+1}^m; \zeta_{t+1}^m) + (1 - \alpha_{t+1}) \left( h_t^m - g^m(x_t^m; \zeta_{t+1}^m) \right)$$

213 
$$u_{t+1}^m = \prod_{C_a} \left[ \nabla g^m(x_{t+1}^m; \zeta_{t+1}^m) + (1 - \beta_{t+1}) \left( u_t^m - \nabla g^m(x_t^m; \zeta_{t+1}^m) \right) \right]$$

214 
$$[j_{l+1}, j_{l+1}, j_{l+$$

215 
$$v_{t+1}^m = \prod_{C_f} \left[ \nabla f^m(h_{t+1}^m; \xi_{t+1}^m) + (1 - \varrho_{t+1}) \left( v_t^m - \nabla f(h_t^m; \xi_{t+1}^m) \right) \right],$$

216 Algorithm 1 MFCGD and AdaMFCGD Algorithms 217 1: Input: T, q, tuning parameters  $\{\gamma, \eta_t, \alpha_t, \beta_t, \varrho_t\}$  and initial input  $x_1 \in \mathbb{R}^d$ ; 218 2: initialize: Set  $x_1^m = x_1$  for  $m \in [M]$ , and draw 2q independent samples  $\{\xi_{1,i}^m\}_{i=1}^q$  and 219  $\{\zeta_{1,j}^m\}_{j=1}^q$ , and then compute  $h_1^m = \frac{1}{q} \sum_{j=1}^q g^m(x_1^m; \zeta_{1,j}^m), u_1^m = \frac{1}{q} \sum_{j=1}^q \nabla g^m(x_1^m; \zeta_{1,j}^m)$  and  $v_1^m = \frac{1}{q} \sum_{j=1}^q \nabla f(h_1^m; \xi_{1,j}^m)$  for all  $m \in [M]$ ; Generate adaptive matrix  $A_1 \in \mathbb{R}^{d \times d}$ . 220 221 3: for t = 1 to T do 222  $\begin{array}{ll} \text{if} \mod(t,q)=0 \text{ then} \\ \bar{w}_t = \frac{1}{M} \sum_{m=1}^M w_t^m \text{ and } \bar{x}_t = \frac{1}{M} \sum_{m=1}^M x_t^m; \\ \text{Generate the adaptive matrix } A_t \in \mathbb{R}^{d \times d}; \end{array}$ 223 4: 224 5: 225 6: 226 One example of  $A_t$  by using update rule ( $a_0 = 0, 0 < \vartheta_t < 1, \rho > 0$ .) Compute  $a_t = \vartheta_t a_{t-1} + (1 - \vartheta_t) \overline{w}_t^2$ ,  $A_t = \text{diag}(\sqrt{a_t} + \rho)$ ; 227  $x_{t+1}^m = \bar{x}_{t+1} = \arg\min_{x \in \mathbb{R}^d} \left\{ \langle x, \bar{w}_t \rangle + \frac{1}{2\eta_t \gamma} (x - \bar{x}_t)^T A_t (x - \bar{x}_t) \right\}; \text{ (Sent them to}$ 228 7: 229 Clients) 230 8: else 231 for each client  $m \in [M]$  (in parallel) do 9: 232  $w_t^m = (u_t^m)^T v_t^m;$ 10:  $x_{t+1}^{m} = \arg\min_{x \in \mathbb{R}^{d}} \left\{ \langle x, w_{t}^{m} \rangle + \frac{1}{2\eta_{t}\gamma} \left( x - x_{t}^{m} \right)^{T} A_{t} \left( x - x_{t}^{m} \right) \right\};$ 233 11: 12:  $A_{t+1} = A_t;$ 235 end for 13: 236 14: end if 237 15: for each client  $m \in [M]$  (in parallel) do 238 Draw two independent samples  $\xi_{t+1}^m$  and  $\zeta_{t+1}^m$ ; 16: 239  $h_{t+1}^m = g^m(x_{t+1}^m; \zeta_{t+1}^m) + (1 - \alpha_{t+1}) \left( h_t^m - g^m(x_t^m; \zeta_{t+1}^m) \right);$ 17: 240  $u_{t+1}^m = \prod_{C_g} \left[ \nabla g^m(x_{t+1}^m; \zeta_{t+1}^m) + (1 - \beta_{t+1}) \left( u_t^m - \nabla g^m(x_t^m; \zeta_{t+1}^m) \right) \right];$ 18: 241  $v_{t+1}^m = \Pi_{C_f} \Big[ \nabla f^m(h_{t+1}^m; \xi_{t+1}^m) + (1 - \varrho_{t+1}) \big( v_t^m - \nabla f(h_t^m; \xi_{t+1}^m) \big) \Big];$ 242 19: 243  $w_{t+1}^m = (u_{t+1}^m)^T v_{t+1}^m;$ 20: 244 end for 21: 245 22: end for 246 23: **Output:** Chosen uniformly random from  $\{\bar{x}_t\}_{t=1}^T$ , where  $\bar{x}_t = \frac{1}{M} \sum_{m=1}^M x_t^m$ . 247 248

where  $\alpha_{t+1} \in (0,1)$ ,  $\beta_{t+1} \in (0,1)$  and  $\varrho_{t+1} \in (0,1)$ . Here the projection functions  $\prod_{C_g} \lfloor \cdot \rfloor$  and  $\prod_{C_f} \lfloor \cdot \rfloor$  ensure that the estimated stochastic gradients  $u_{t+1}^m$  and  $v_{t+1}^m$  are bounded, i.e.,  $\|u_{t+1}^m\| \leq C_g$  and  $\|v_{t+1}^m\| \leq C_f$  for any  $t \geq 1$ . Based on the estimated stochastic gradients and adaptive matrices, the clients update the variables  $\{x_t^m\}_{m=1}^M$ , defined as

$$x_{t+1}^{m} = x_{t}^{m} - \gamma \eta_{t} A_{t}^{-1} w_{t}^{m} = \arg \min_{x \in \mathbb{R}^{d}} \left\{ \langle x, w_{t}^{m} \rangle + \frac{1}{2\eta_{t} \gamma} \left( x - x_{t}^{m} \right)^{T} A_{t} \left( x - x_{t}^{m} \right) \right\},$$
(5)

where  $\gamma > 0$  and  $\eta_t > 0$ . In our algorithms, all clients use the same adaptive matrix generated from the server as in Chen et al. (2020c). Note that the existing adaptive FL algorithms such as local-AMSGrad Chen et al. (2020c) only builds on some specific adaptive learning rates such as AMSGrad Reddi et al. (2019). However, our algorithms can use the unified adaptive matrix to flexibly incorporate various adaptive learning rates.

# 4 CONVERGENCE ANALYSIS

249

250

251

252

257

258

259

260 261

262

In this section, we study the convergence properties of our MFCGD and AdaMFCGD algorithms
 under some mild assumptions. All related proofs are provided in the Appendix A. We first review
 some useful lemmas and assumptions.

**Assumption 1.** (*Lipschitz Gradients*) For any  $m \in [M]$ , there exist constants  $L_f$  and  $L_g$  for  $\nabla f^m(y;\xi^m), \nabla g^m(x;\zeta^m)$  respectively satisfying

268  
269  

$$\begin{aligned} \|\nabla g^m(x_1, \zeta^m) - \nabla g^m(x_2, \zeta^m)\| \le L_g \|x_1 - x_2\|, \ \forall x_1, x_2 \in \mathbb{R}^d, \\ \|\nabla f^m(y_1; \zeta^m) - \nabla f^m(y_2; \zeta^m)\| \le L_f \|y_1 - y_2\|, \ \forall y_1, y_2 \in \mathbb{R}^p. \end{aligned}$$

**Assumption 2.** (Bounded Gradients) For any  $m \in [M]$ , gradient  $\nabla q^m(x; \zeta^m)$  and Jacobian matrix  $\nabla f^m(y;\xi^m)$  have the upper bounds  $C_q$  and  $C_f$  respectively, i.e., 

$$\|\nabla g^m(x;\zeta^m)\| \le C_q, \|\nabla f^m(y;\zeta^m)\| \le C_f, \forall x \in \mathbb{R}^d, y \in \mathbb{R}^p.$$

**Assumption 3.** (Bounded Variances) For any  $m \in [M]$ , functions  $f^m(y;\xi^m)$  and  $g^m(x;\zeta^m)$ and its gradients are unbiased and the bounded variances, i.e., we have  $\mathbb{E}[g^m(x;\zeta^m)] = g^m(x)$ ,  $\mathbb{E}[\nabla g^m(x;\zeta^m)] = \nabla g^m(x), \ \mathbb{E}[\nabla f^m(y;\xi^m)] = \nabla f^m(y) \text{ and }$ 

$$\mathbb{E}\|g^m(x;\zeta^m) - g^m(x)\|^2 \le \sigma^2, \ \mathbb{E}\|\nabla g^m(x;\zeta^m) - \nabla g^m(x)\|^2 \le \sigma^2.$$

$$\mathbb{E}\|\nabla f^m(y;\xi^m) - \nabla f^m(y)\|^2 \le \sigma^2, \quad \forall x \in \mathbb{R}^d, \ y \in \mathbb{R}^p$$

where  $\sigma > 0$ .

**Assumption 4.** F(x) has a lower bound, i.e.,  $F^* = \inf_{x \in \mathbb{R}^d} F(x)$ . 

**Assumption 5.** In our algorithms, the adaptive matrix  $A_t$  for all  $t \ge 1$  satisfies  $A_t \succeq \rho I_d$ , where  $\rho > 0$  is an appropriate positive number. 

**Assumption 6.** For any  $m, j \in [M]$ ,  $x \in \mathbb{R}^d$  and  $y \in \mathbb{R}^p$ , we have  $\|\nabla f^m(y) - \nabla f^j(y)\| \leq \delta_f$ ,  $\|\nabla g^m(x) - \nabla g^j(x)\| \le \delta_g$  and  $\|g^m(x) - g^j(x)\| \le \delta_g$ , where  $\delta_f > 0$  and  $\delta_q > 0$  are constants. 

Assumptions 1 ensures the smoothness of functions  $f^m(y;\xi^m)$ ,  $g^m(x;\zeta^m)$  for any  $m \in [M]$ , Assumption 2 ensures the bounded gradients (or Jacobian matrix) of functions  $f^m(y;\xi^m)$  and  $g^m(x;\zeta^m)$  for any  $m \in [M]$ . Assumption 3 ensures the bounded variances of stochastic gradi-ent or value of functions  $f^m(y;\xi^m)$  and  $g^m(x;\zeta^m)$  for any  $m \in [M]$ . Assumption 4 guarantees the feasibility of the problem (1). Assumptions 1-4 have been commonly used in the convergence analysis of the stochastic composition algorithms Wang et al. (2017a;b). Assumption 5 has been commonly used in the existing adaptive methods Huang et al. (2021b). Assumption 6 is the stan-dard condition constrained the data heterogeneity in non-i.i.d FL setting Li et al. (2019). In fact, we can obtain the part results of Assumption 6 based on Assumptions 1-2. For example, we have  $\|\nabla f^m(y) - \nabla f^j(y)\| \le 2\sigma + 2C_f$ , where the last inequality holds by Assumptions 1-2. Similarly, we have  $\|\nabla g^m(y) - \nabla g^j(y)\| \le 2\sigma + 2C_q$  based on Assumptions 1-2. 

### 4.1 CONVERGENCE PROPERTIES OF ADAMFCGD ALGORITHM

In this subsection, we provide the convergence properties of our AdaMFCGD algorithm.

**Theorem 1.** Assume the sequence  $\{\bar{x}_t\}_{t=1}^T$  be generated from AdaMFCGD algorithm. Under the above Assumptions, and let  $\eta_t = \frac{k}{(n+t)^{1/3}}$  for all  $t \ge 0$ ,  $\alpha_{t+1} = c_1\eta_t^2$ ,  $\beta_{t+1} = c_2\eta_t^2$ ,  $\varrho_{t+1} = c_3\eta_t^2$ , 
$$\begin{split} n &\geq \max\left(2, k^3, (c_1k)^3, (c_2k)^3, (c_3k)^3, \frac{(24k\gamma qL_{fg}C_{fg})^3}{\rho^3}\right), k > 0, c_1 \geq \frac{2}{3k^3} + B, c_2 \geq \frac{2}{3k^3} + 5C_f^2, \\ c_1^2 + c_2^2 &\leq \frac{(24)^4 q^2 \gamma^4 L_{fg}^4 C_{fg}^4}{9\rho^4}, c_3 \geq \frac{2}{3k^3} + 5C_g^2, \frac{\rho(c_1^2 + c_3^2)^{1/4}}{12\sqrt{5q}L_{fg}C_{fg}} \leq \gamma \leq \min\left(\frac{3\rho qL_{fg}C_{fg}}{4(C_g^2 + L_g^2 + 2L_f^2C_g^2)}, \frac{n^{1/3}\rho}{2Lk}\right), \\ B \geq 20C_g^2 L_f^2 + \frac{c_2^2 C_g^2 L_f^2}{216q^3 \gamma^3 L_{fg}^3 C_{fg}^3} + \frac{\Theta \rho^2(c_1^2 + c_3^2)}{30q^2 \gamma^4 C_{fg}^2 L_{fg}^2 C_g^2}, \Theta = \left(5C_f^2 L_g^2 + \frac{c_2^2 C_g^2 L_f^2}{864q^3 \gamma^3 L_{fg}^3 C_{fg}^3}\right) \frac{\rho^2}{(24)^2 L_{fg}^2 C_{fg}^2} + \frac{\gamma\rho}{6qL_{fg}C_{fg}} \left(C_g^2 + L_g^2 + 2L_f^2 C_g^2\right) and \Theta + \frac{BC_g^2 \rho^2}{(24)^2 L_{fg}^2 C_{fg}^2} \leq \frac{5\rho^2}{48}, we have \end{split}$$

$$\frac{1}{T}\sum_{t=1}^{T} \mathbb{E}\|\nabla F(\bar{x}_t)\| \le \left(\frac{\sqrt{2G}n^{1/6}}{T^{1/2}} + \frac{\sqrt{2G}}{T^{1/3}}\right) \sqrt{\frac{1}{T}\sum_{t=1}^{T} \mathbb{E}\|A_t\|^2},\tag{6}$$

$$\begin{array}{l} \textbf{313} \qquad \text{where } C_{fg}^2 = \max(C_f^2, C_g^2), \ L_{fg}^2 = L_f^2 C_g^2 + L_g^2, \ G = \frac{4(F(\bar{x}_1) - F^*)}{k\rho\gamma} + \frac{12n^{1/3}\sigma^2}{qk^2\rho^2} + 4k^2 \Big(\frac{\hat{\delta}^2}{4\gamma^2 L_{fg}^2} + \frac{\delta^2}{qk^2\rho^2} + \frac{\delta^2}{qk^2\rho^2} \Big) \\ \textbf{314} \qquad \qquad \frac{\left(c_1^2 + c_2^2 + c_3^2\right)\sigma^2}{3\rho\gamma qL_{fg}C_{fg}}\Big) \ln(n+T) \ \text{and} \ \hat{\delta}^2 = 2c_1^2 L_f^2 \sigma^2 + c_3^2 \sigma^2 + 4c_3^2 \delta_f^2 + 4c_3^2 L_f^2 \delta_g^2 + c_2^2 \sigma^2 + 3c_2^2 \delta_g^2. \end{array}$$

**Remark 1.** Under the above Assumption 2, we have  $\left\|\frac{1}{M}\sum_{m=1}^{M} \left(\nabla g^{m}(\bar{x}_{t})\right)^{T} \nabla f^{m}(g^{m}(\bar{x}_{t}))\right\| \leq C_{f}C_{g}$ . When the adaptive matrix  $A_{t}$  be generated from the line 6 of Algorithm 1, we have  $\sqrt{\frac{1}{T}\sum_{t=1}^{T} \mathbb{E}||A_t||^2} \le 2(C_f^2 C_g^2 + \rho).$  Without loss of generality, let  $k = O(1), \rho = O(1), c_1 = O(1), c_2 = O(1), c_3 = O(1), c_4 = O(1), c_5 = O(1), c_6 = O(1), c_7 = O(1), c_8 = O(1),$  $c_2 = O(1), c_3 = O(1)$  and  $n = O(q^3)$ , we have and  $G = \tilde{O}(1)$ . Let  $q = T^{1/3}$  and 

322  
323 
$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E} \|\nabla F(\bar{x}_t)\| \le \tilde{O}\left(\frac{\sqrt{q}}{\sqrt{T}} + \frac{1}{T^{1/3}}\right) = \tilde{O}\left(\frac{1}{T^{1/3}}\right) \le \epsilon,$$
(7)

then we have  $T = \tilde{O}(\epsilon^{-3})$ . Since our AdaMFCGD algorithm requires 2 samples at each iteration expect for the first iteration requires 2q samples, it has a sample complexity of  $2q + 2T = \tilde{O}(\epsilon^{-3})$ . Thus, our AdaMFCGD algorithm requires  $\tilde{O}(\epsilon^{-3})$  sample (or gradient) complexity and  $\frac{T}{q} = T^{2/3} = \tilde{O}(\epsilon^{-2})$  communication complexity to find an  $\epsilon$ -stationary point of the Problem (1).

**Remark 2.** From Theorem 1, our AdaMFCGD algorithm simultaneously have lower sample and communication complexities than the existing federated compositional optimization algorithms (Please see Table 1). Moreover, our AdaMFCGD algorithm simultaneously have lower sample and communication complexities than the existing adaptive single-level FL algorithms such as the local-AMSGrad Chen et al. (2020c) algorithm that needs sample complexity of  $O(\epsilon^{-4})$  and communication complexity of  $O(\epsilon^{-3})$  for finding an  $\epsilon$ -stationary point of the distributed single-level optimization problem, i.e., the above problem (1) with  $g^m(x) = x$  for all  $m \in [M]$ .

# 4.2 CONVERGENCE PROPERTIES OF MFCGD ALGORITHM

In this subsection, we provide the convergence properties of our non-adaptive **MFCGD** algorithm, i.e., set  $A_t = I_d$  for all  $t \ge 1$ .

**Theorem 2.** Assume the sequence  $\{\bar{x}_t\}_{t=1}^T$  be generated from **MFCGD** algorithm, i.e.,  $A_t = I_d$  for all  $t \ge 1$  in Algorithm 1. Under the above Assumptions, and let  $\eta_t = \frac{k}{(n+t)^{1/3}}$  for all  $t \ge 0$ ,  $\alpha_{t+1} = c_1\eta_t^2$ ,  $\beta_{t+1} = c_2\eta_t^2$ ,  $\varrho_{t+1} = c_3\eta_t^2$ ,  $n \ge \max(2, k^3, (c_1k)^3, (c_2k)^3, (c_3k)^3, (24k\gamma qL_{fg}C_{fg})^3)$ , k > 0,  $c_1 \ge \frac{2}{3k^3} + B$ ,  $c_2 \ge \frac{2}{3k^3} + 5C_f^2$ ,  $c_1^2 + c_2^2 \le \frac{(24)^4q^2\gamma^4L_{fg}C_{fg}^4}{9}$ ,  $c_3 \ge \frac{2}{3k^3} + 5C_g^2$ ,  $\frac{(c_1^2+c_3^2)^{1/4}}{12\sqrt{5q}L_{fg}C_{fg}} \le \gamma \le \min\left(\frac{3qL_{fg}C_{fg}}{4(C_g^2+L_g^2+2L_f^2C_g^2)}, \frac{n^{1/3}}{2Lk}\right)$ ,  $B \ge 20C_g^2L_f^2 + \frac{c_2^2C_g^2L_f^2}{216q^3\gamma^3L_{fg}^3C_{fg}^3} + \frac{\Theta(c_1^2+c_3^2)}{30q^2\gamma^4C_{fg}^2L_{fg}^2C_{fg}^2}$ ,  $\Theta = \left(5C_f^2L_g^2 + \frac{c_2^2C_g^2L_f^2}{864q^3\gamma^3L_{fg}^3C_{fg}^3}\right)\frac{1}{(24)^2L_{fg}^2C_{fg}^2} + \frac{\gamma}{6qL_{fg}C_{fg}}\left(C_g^2 + L_g^2 + 2L_f^2C_g^2\right)$  and  $\Theta + \frac{BC_g^2}{(24)^2L_{fg}^2C_{fg}^2} \le \frac{5}{48}$ , we have

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E} \|\nabla F(\bar{x}_t)\| \le \frac{\sqrt{2G}n^{1/6}}{T^{1/2}} + \frac{\sqrt{2G}}{T^{1/3}},\tag{8}$$

where  $C_{fg}^2 = \max(C_f^2, C_g^2)$ ,  $L_{fg}^2 = L_f^2 C_g^2 + L_g^2$ ,  $G = \frac{4(F(\bar{x}_1) - F^*)}{k\gamma} + \frac{12n^{1/3}\sigma^2}{qk^2} + 4k^2 \Big(\frac{\hat{\delta}^2}{4\gamma^2 L_{fg}^2} + \frac{(c_1^2 + c_2^2 + c_3^2)\sigma^2}{3\gamma q L_{fg} C_{fg}}\Big) \ln(n+T)$  and  $\hat{\delta}^2 = 2c_1^2 L_f^2 \sigma^2 + c_3^2 \sigma^2 + 4c_3^2 \delta_f^2 + 4c_3^2 L_f^2 \delta_g^2 + c_2^2 \sigma^2 + 3c_2^2 \delta_g^2$ .

**Remark 3.** The proof of Theorem 2 can totally follow the proofs of the above Theorem 1 with the parameter  $\rho = 1$ . Without loss of generality, let k = O(1),  $c_1 = O(1)$ ,  $c_2 = O(1)$ ,  $c_3 = O(1)$  and  $n = O(q^3)$ , we have and  $G = \tilde{O}(1)$ . Let  $q = T^{1/3}$  and

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\|\nabla F(\bar{x}_t)\| \le \tilde{O}\left(\frac{\sqrt{q}}{\sqrt{T}} + \frac{1}{T^{1/3}}\right) = \tilde{O}\left(\frac{1}{T^{1/3}}\right) \le \epsilon,\tag{9}$$

then we have  $T = \tilde{O}(\epsilon^{-3})$ . Since our MFCGD algorithm requires 2 samples at each iteration expect for the first iteration requires 2q samples, it has a sample complexity of  $2q + 2T = \tilde{O}(\epsilon^{-3})$ . As the above AdaMFCGD algorithm, our MFCGD algorithm also obtain lower sample complexity of  $\tilde{O}(\epsilon^{-3})$  and communication complexity of  $\tilde{O}(\epsilon^{-2})$  in finding an  $\epsilon$ -stationary solution of the Problem (1).

# 5 NUMERICAL EXPERIMENTS

In this section, we apply some numerical experiments to demonstrate efficiency of our MFCGD and
 AdaMFCGD algorithms on the robust federated learning and distributed meta learning tasks. Note
 that the experiment on distributed meta learning task is given in the Appendix B. In the experiments,
 we compare our algorithms with the existing federated composition optimization algorithms in Table
 1 for solving distributed composition optimization problems.



Figure 1: The performances of various FCO methods are evaluated on a synthetic imbalanced dataset based on MNIST, with a focus on addressing the distributed problem (10). The results are visualized using four plots: the first two show the accuracy(%) and loss on the training set, while the last two show the accuracy(%) and loss on the test set. These plots provide insights into how the FL methods perform on imbalanced data and their effectiveness in improving model performance while ensuring fairness across different clients. 393

### 5.1 **ROBUST FEDERATED LEARNING**

388

389

390

391

392

394 395

397

399

In this subsection, we evaluate the efficacy of our algorithms by performing a distributionally robust federated learning task defined in (3). Specifically, this robust federated learning problem can be rewritten as into a distributed composition optimization problem as in Huang et al. (2021a),

$$\min_{x \in \mathbb{R}^d} \frac{1}{M} \sum_{m=1}^M f\Big(g^m(x)/\lambda\Big),\tag{10}$$

where  $f(\cdot) = \exp(\cdot/\lambda)$  and  $\lambda > 0$  is a regularization parameter. In fact, we can also use some other 404 monotonically increasing functions instead of  $f(\cdot)$ . 405

406 In the experiments, we tackle a multi-class classification problem on the MNISTLeCun et al. (2010) 407 dataset with a 3-layer Convolutional Neural Network (CNN) which is widely used in this task. The experiments are conducted on a network comprising 10 clients and 1 server. To introduce data im-408 balance across clients, we randomly select one client to have a larger dataset of 5000 images, while 409 the remaining clients have a significantly smaller dataset of only 20 images. This unequal distri-410 bution of data aims to create a challenging scenario where the algorithm must focus on the hardest 411 and most important task, namely the client with the dominant number of images, to achieve good 412 performance. In this way, we can test the algorithm's adaptability and ability to handle this highly 413 imbalanced dataset. In these experiments, we performed a grid search to identify the optimal hy-414 perparameters for each method, which we discussed in detail in subsection B.2. We set the learning 415 rate to 0.01 for all methods and used in our Algorithm 1 for adaptive matrix generation. We fixed 416 the total number of training steps to 500 and set the asynchronization step q to 5 if not specified.

417 From Figure 1, we can find that our MFCGD and AdaMFCGD algorithms have a faster convergence 418 rate and more stable optimization processes compared to the other composition federated optimiza-419 tion algorithms, partly contributes to momentum-based variance reduced strategy. Specifically, the 420 experimental results show that our algorithms outperformed the existing composition federated op-421 timization approaches, such as ComFedL Huang et al. (2021a), FEDNEST Tarzanagh et al. (2022), 422 and Local-SCGDM Gao et al. (2022), in terms of both accuracy rates and cross-entropy losses. Moreover, the comparison between MFCGD and AdaMFCGD methods highlighted the advantage of 423 the unified adaptive matrix, which flexibly incorporates various adaptive learning rates. Meanwhile, 424 Figure 2 demonstrates the robustness of our AdaMFCGD algorithm by varying the asynchronization 425 step q. It is worth noting that our AdaMFCGD achieves its optimal performance when q = 1, as 426 this implies that an adaptive matrix is calculated in each iteration, allowing the momentum-based 427 variance reduction technique to fully showcase its potential. 428

Figure 2 shows the robustness of our AdaMFCGD algorithm by varying the regularization parameter 429  $\lambda$ . The results indicate that our AdaMFCGD algorithm achieves good test accuracy and low test loss 430 with different  $\lambda$ . Moreover, when decreasing  $\lambda$ , our AdaMFCGD algorithm converges much faster, 431 especially when running multiple local epochs. As in Huang et al. (2021a),  $\lambda$  is a penalty parameter



Figure 2: Comparing the accuracy(%) (left) and cross-entropy loss (right) on different synchronization step q and different regularization parameter  $\lambda$ , respectively, on our AdaMFCGD algorithm

in distributionally robust FL that penalizes divergence between  $r = (r_1, \dots, r_M)$  and  $1/M = (1/M, \dots, 1/M)$ , where  $r_m \in (0, 1)$  for any  $m \in [M]$  denotes the proportion of *m*-th client in the entire mode. A larger  $\lambda$  means places more emphasis on bringing the original proportion parameter  $r_m$  closer to the average weight 1/M. Since the optimal weights in our dataset are far from 1/M, we should select  $\lambda$  relatively small. The results in Figure 2 also confirm that our AdaMFCGD algorithm converges much faster with smaller  $\lambda$ . Specifically, our AdaMFCGD algorithm converges much faster, when choosing  $\lambda = 0.2$  or  $\lambda = 0.5$  compared to that  $\lambda = 2$  in terms of both accuracy rate, cross-entropy loss on true labels and convergence speed.

452 453 454

455

463

464

474

442

443 444 445

446

447

448

449

450

451

# 6 CONCLUSION

In the paper, we proposed a class of faster adaptive momentum-based federated compositional gradient descent methods to solve the nonconvex distributed composition problems based on the localSGD and momentum-based variance reduced techniques. In particular, our adaptive algorithm (i.e.,
AdaMFCGD) uses a unified adaptive matrix to flexibly incorporate various adaptive learning rates
further accelerate algorithm. Moreover, we established a solid convergence analysis framework for
our methods, and proved that they obtain lower sample and communication complexities simultaneously than the existing federated composition optimization methods.

# References

- Marcin Andrychowicz, Misha Denil, Sergio Gómez Colmenarejo, Matthew W Hoffman, David
  Pfau, Tom Schaul, Brendan Shillingford, and Nando de Freitas. Learning to learn by gradient
  descent by gradient descent. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pp. 3988–3996, 2016.
- Yair Censor and Arnold Lent. An iterative row-action method for interval convex programming.
   *Journal of Optimization theory and Applications*, 34(3):321–353, 1981.
- Yair Censor and Stavros Andrea Zenios. Proximal minimization algorithm withd-functions. *Journal of Optimization Theory and Applications*, 73(3):451–464, 1992.
- Cheng Chen, Ziyi Chen, Yi Zhou, and Bhavya Kailkhura. Fedcluster: Boosting the convergence of federated learning via cluster-cycling. In *2020 IEEE International Conference on Big Data (Big Data)*, pp. 5017–5026. IEEE, 2020a.
- Tianyi Chen, Yuejiao Sun, and Wotao Yin. Solving stochastic compositional optimization is nearly
   as easy as solving stochastic optimization. *arXiv preprint arXiv:2008.10847*, 2020b.
- Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong. On the convergence of a class of adam-type algorithms for non-convex optimization. In *7th International Conference on Learning Representations (ICLR)*, 2019.
- Xiangyi Chen, Xiaoyun Li, and Ping Li. Toward communication efficient adaptive gradient method. In *Proceedings of the 2020 ACM-IMS on Foundations of Data Science Conference*, pp. 119–128, 2020c.

486 487	Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex sgd. <i>Advances in neural information processing systems</i> , 32, 2019.
400 489	Yuyang Deng and Mehrdad Mahdavi. Local stochastic gradient descent ascent: Convergence anal-
490 491	Statistics, pp. 1387–1395. PMLR, 2021.
492	
493	Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Adaptive personalized federated
494	icarining. <i>arxiv preprint arxiv.2003.13401</i> , 2020a.
495	Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Distributionally robust federated
496 497	averaging. Advances in Neural Information Processing Systems, 33:15111–15122, 2020b.
498 499	John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. <i>Journal of machine learning research</i> , 12(7), 2011.
500 501	Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning: A meta- learning approach. <i>arXiv preprint arXiv:2002.07948</i> , 2020.
502	Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near optimal non-convey op
503 504	timization via stochastic path-integrated differential estimator. In Advances in Neural Information Processing Systems, pp. 689–699, 2018
505	<i>Trocessing Bystems</i> , pp. 009-009, 2010.
506	Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation
507	of deep networks. In International Conference on Machine Learning, pp. 1126–1135. PMLR, 2017
508	2017.
509	Hongchang Gao, Junyi Li, and Heng Huang. On the convergence of local stochastic compositional
510	gradient descent with momentum. In International Conference on Machine Learning, pp. 7017-
510	7035. PMLR, 2022.
513	Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Mini-batch stochastic approximation meth-
514	ods for nonconvex stochastic composite optimization. <i>Mathematical Programming</i> , 155(1-2): 267–305–2016
515	207 505, 2010.
517 518	Saeed Ghadimi, Andrzej Ruszczynski, and Mengdi Wang. A single timescale stochastic approxima- tion method for nested stochastic optimization. <i>SIAM Journal on Optimization</i> , 30(1):960–979, 2020.
519	
520 521	Feihu Huang and Shangqian Gao. Riemannian gradient methods for stochastic composition prob- lems. <i>Neural Networks</i> , 153:224–234, 2022.
522 523 524	Feihu Huang, Junyi Li, and Heng Huang. Compositional federated learning: Applications in distributionally robust averaging and meta learning. <i>arXiv preprint arXiv:2106.11264</i> , 2021a.
525 526	Feihu Huang, Junyi Li, and Heng Huang. Super-adam: faster and universal framework of adaptive gradients. <i>Advances in Neural Information Processing Systems</i> , 34:9074–9085, 2021b.
527	They was the Din Cy. Ii Liu and Hang Huang Accelerated with a few stackastic and with
528	chouyuan nuo, Bin Gu, Ji Liu, and Heng Huang. Accelerated method for stochastic composition optimization with nonsmooth regularization. In <i>Proceedings of the AAAI Conference on Artificial</i>
529	Intelligence, volume 32, 2018.
530	<b>3</b> , <b>1</b>
531	Wei Jiang, Bokun Wang, Yibo Wang, Lijun Zhang, and Tianbao Yang. Optimal algorithms for s-
532 533	ing, pp. 10195–10216. PMLR, 2022.
534	Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Ariun Nitin
535 536	Bhagoji, Keith Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances
537	and open problems in federated learning. arXiv preprint arXiv:1912.04977, 2019.
538 539	Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian Stich, and Martin Jaggi. Error feedback fixes signsgd and other gradient compression schemes. In <i>International Conference on Machine Learning</i> , pp. 3252–3261. PMLR, 2019.

540	Sai Dranaath Karimiraddy, Satyan Kala, Mahryar Mahri, Sashank Daddi, Sahastian Stich, and Anan
541	sal Praneeth Karimireddy, Salyen Kale, Menryar Monri, Sasnank Reddi, Sebastian Stich, and Anan- da Theartha Surash, Sasffold, Stochastia controllad averaging for federated learning. In Interna
542	tional Conference on Machine Learning, pp. 5132–5143 PMLR 2020
543	nonai Conjerence on Machine Learning, pp. 5152–5145. 1 MLR, 2020.
544	Prashant Khanduri, Pranay Sharma, Haibo Yang, Mingyi Hong, Jia Liu, Ketan Rajawat, and Pramod
545	Varshney. Stem: A stochastic two-sided momentum algorithm achieving near-optimal sample and
546	communication complexities for federated learning. Advances in Neural Information Processing
547	Systems, 34:6050–6061, 2021.
548	D' 1. 'I D Z' and a l I' and D a Altan A mathed for the back and a fail of a V'
549	Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint
550	<i>urxiv</i> .1412.0900, 2014.
551	Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
552	2009.
553	
554	Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. ATT Labs [Online].
555	Available: http://yann.lecun.com/exdb/mnist, 2, 2010.
556	Oinhin Li Zavi Wan Zhaomin Wu Siyu Hu Naiha Wang Yuan Li Yu Liu and Dingshang Ha
557	A survey on federated learning systems: vision, hype and reality for data privacy and protection
558	IFFF Transactions on Knowledge and Data Engineering 2021a
550	TEEE Transactions on Knowledge and Data Engineering, 2021a.
560	Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated
561	learning through personalization. In International Conference on Machine Learning, pp. 6357-
501	6368. PMLR, 2021b.
562	
505	Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of federed on non-iid data, arXiv proprint arXiv 1007 02180, 2010
504	ledavg on non-nd data. arxiv preprint arxiv:1907.02189, 2019.
505	Tianvi Lin, Chenyou Fan, Mengdi Wang, and Michael L Jordan. Improved sample complexity for
500	stochastic compositional variance reduced gradient. arXiv preprint arXiv:1806.00458, 2018.
569	
560	Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In International Confer-
570	ence on Learning Representations, 2018.
571	Brendan McMahan Fider Moore Daniel Ramage Seth Hampson and Blaise Aguera y Arcas
572	Communication-efficient learning of deep networks from decentralized data. In Artificial Intelli-
572	gence and Statistics, pp. 1273–1282. PMLR, 2017.
574	
575	Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In Interna-
575	tional Conference on Machine Learning, pp. 4615–4625. PMLR, 2019.
570	Low M Neuron Jie Lin Vetra Scheinhaus and Martin Telefe Seach A neural method for machine
578	Lanii vi inguyen, jie Liu, Kaiya Schenheig, and Martin Takac. Saran: A novel method for machine learning problems using stochastic recursive gradient. In International Conference on Machine
570	Learning problems using stochastic recursive gradient. In International Conjerence on Machine
580	Leanning, pp. 2015 2021. 1 milit, 2017.
500	Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,
582	Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. arXiv preprint arX-
582	iv:2003.00295, 2020.
584	Cashank I Daddi Satuan Kala and Saniin Kumar On the commence of stars and the set
585	sashalik J Keudi, Salyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. <i>arXiv</i>
505	preprina ar Aiv. 1904.09237, 2017.
587	Amirhossein Reisizadeh, Farzan Farnia, Ramtin Pedarsani, and Ali Jadbabaie. Robust federated
588	learning: The case of affine distribution shifts. In NeurIPS, 2020.
580	
505	Sebastian Urban Stich. Local sgd converges fast and communicates little. In International Confer-
501	ence on Learning Representations (ICLR), 2019.
502	Dayoud Ataee Tarzanagh Mingchen Li Christos Thrampoulidis and Samet Oymak Fednest Fed-
592	erated bilevel, minimax, and compositional optimization. In International Conference on Machine
555	<i>Learning</i> , pp. 21146–21179. PMLR, 2022.

594 595 596	Quoc Tran-Dinh, Nhan H Pham, Dzung T Phan, and Lam M Nguyen. A hybrid stochastic opti- mization framework for composite nonconvex optimization. <i>Mathematical Programming</i> , 191 (2):1005–1071, 2022.
597 598 599	Rasul Tutunov, Minne Li, Jun Wang, and Haitham Bou-Ammar. Compositional adam: An adaptive compositional solver. <i>arXiv preprint arXiv:2002.03755</i> , 2020.
600 601	Bokun Wang, Zhuoning Yuan, Yiming Ying, and Tianbao Yang. Memory-based optimization meth- ods for model-agnostic meta-learning. <i>arXiv preprint arXiv:2106.04911</i> , 2021.
602 603 604 605	Mengdi Wang, Ethan X Fang, and Han Liu. Stochastic compositional gradient descent: algorithms for minimizing compositions of expected-value functions. <i>Mathematical Programming</i> , 161(1-2): 419–449, 2017a.
606 607	Mengdi Wang, Ji Liu, and Ethan X Fang. Accelerating stochastic composition optimization. <i>The Journal of Machine Learning Research</i> , 18(1):3721–3743, 2017b.
608 609 610	Jie Xu, Benjamin S Glicksberg, Chang Su, Peter Walker, Jiang Bian, and Fei Wang. Federated learning for healthcare informatics. <i>Journal of Healthcare Informatics Research</i> , 5(1):1–19, 2021.
611 612 613	Qian Yang, Jianyi Zhang, Weituo Hao, Gregory P Spell, and Lawrence Carin. Flop: Federated learning on medical datasets using partial networks. In <i>Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery &amp; Data Mining</i> , pp. 3845–3853, 2021.
614 615 616	Honglin Yuan and Tengyu Ma. Federated accelerated stochastic gradient descent. Advances in Neural Information Processing Systems, 33:5332–5344, 2020.
617 618	Zhuoning Yuan, Zhishuai Guo, Nitesh Chawla, and Tianbao Yang. Compositional training for end- to-end deep auc maximization. In <i>International Conference on Learning Representations</i> , 2022.
619 620 621	Junyu Zhang and Lin Xiao. Multi-level composite stochastic optimization via nested variance re- duction. <i>arXiv preprint arXiv:1908.11468</i> , 2019.
622 623 624	Lefeng Zhang, Tianqing Zhu, Ping Xiong, Wanlei Zhou, and S Yu Philip. A robust game-theoretical federated learning framework with joint differential privacy. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 35(4):3333–3346, 2022.
625 626	
627	
628	
629	
630	
631	
632	
633	
634	
635	
636	
637	
638	
639	
640	
641	
642	
643	
644	
645	
646	
647	

# A CONVERGENCE ANALYSIS

In this section, we provide the detailed convergence analysis of our algorithms.

We first introduce some useful notations:  $\bar{w}_t = \frac{1}{M} \sum_{m=1}^M w_t^m$ ,  $\bar{x}_t = \frac{1}{M} \sum_{m=1}^M x_t^m$ ,

<sup>657</sup> Next, we review and provide some useful lemmas.

**Lemma 1.** Given M vectors  $\{u^m\}_{m=1}^M$ , the following inequalities satisfy:  $||u^m + u^j||^2 \le (1 + c)||u^m||^2 + (1 + \frac{1}{c})||u^j||^2$  for any c > 0, and  $||\sum_{m=1}^M u^m||^2 \le M \sum_{m=1}^M ||u^m||^2$ .

 $F(x) = \frac{1}{M} \sum_{m=1}^{M} f^{m}(g^{m}(x)), \quad \nabla F(x) = \frac{1}{M} \sum_{m=1}^{M} \left( \nabla g^{m}(x) \right)^{T} \nabla f^{m}(g^{m}(x)).$ 

**661 Lemma 2.** Given a finite sequence  $\{u^m\}_{m=1}^M$ , and  $\bar{u} = \frac{1}{M} \sum_{m=1}^M u^m$ , the following inequality **662** satisfies  $\sum_{m=1}^M \|u^m - \bar{u}\|^2 \le \sum_{m=1}^M \|u^m\|^2$ .

Given a  $\rho$ -strongly convex function  $\varphi(x)$ , we define a prox-function (Bregman distance) Censor & Lent (1981); Censor & Zenios (1992) associated with  $\varphi(x)$  as follows:

$$D(z,x) = \varphi(z) - \left[\varphi(x) + \langle \nabla \varphi(x), z - x \rangle\right].$$
(11)

Then we define a generalized projection problem as in Ghadimi et al. (2016): 

$$x^{+} = \arg\min_{z \in \mathcal{X}} \left\{ \langle z, w \rangle + \frac{1}{\gamma} D(z, x) + h(z) \right\},$$
(12)

where  $\mathcal{X} \subseteq \mathbb{R}^d$ ,  $w \in \mathbb{R}^d$  and  $\gamma > 0$ . In the paper, we consider h(x) = 0. Meanwhile, we also define a generalized projected gradient (a.k.a., gradient mapping):

$$\mathcal{G}_{\mathcal{X}}(x, w, \gamma) = \frac{x - x^+}{\gamma}.$$
(13)

**Lemma 3.** (Lemma 1 in Ghadimi et al. (2016)) Let  $x^+$  be given in (12). Then, for any  $x \in \mathcal{X}$ ,  $w \in \mathbb{R}^d$  and  $\gamma > 0$ , we have

$$\langle w, \mathcal{G}_{\mathcal{X}}(x, w, \gamma) \rangle \ge \rho \| \mathcal{G}_{\mathcal{X}}(x, w, \gamma) \|^2 + \frac{1}{\gamma} \big[ h(x^+) - h(x) \big], \tag{14}$$

where  $\rho > 0$  depends on  $\rho$ -strongly convex function  $\varphi(x)$ .

When h(x) = 0, in the above lemma 3, we have

$$\langle w, \mathcal{G}_{\mathcal{X}}(x, w, \gamma) \rangle \ge \rho \| \mathcal{G}_{\mathcal{X}}(x, w, \gamma) \|^2.$$
 (15)

**Lemma 4.** (*Restatement of Lemma 1*) Given the above Assumptions 1-2, the function F(x) is L-smooth, i.e., for any  $x_1, x_2 \in \mathbb{R}^d$ , we have

$$\|\nabla F(x_1) - \nabla F(x_2)\|^2 \le L^2 \|x_1 - x_2\|^2,$$
(16)

where  $L = \sqrt{2C_{f}^{2}L_{g}^{2} + 2C_{g}^{4}L_{f}^{2}}$ .

*Proof.* Based on Assumptions 1-2, the deterministic functions  $f^m(y) = \mathbb{E}[f^m(y;\xi^m)] g^m(x) = \mathbb{E}[f^m(x;\zeta^m)]$  and its gradients also satisfy the Lipschitz gradients and bounded gradients. For example, for any  $y_1, y_2 \in \mathbb{R}^n$ 

$$\|\nabla f^{m}(y_{1}) - \nabla f^{m}(y_{1})\| = \|\mathbb{E} [\nabla f^{m}(y_{1};\xi^{m}) - \nabla f^{m}(y_{1};\xi^{m})]\|$$
  
 
$$\leq \mathbb{E} \|\nabla f^{m}(y_{1};\xi^{m}) - \nabla f^{m}(y_{1};\xi^{m})\| \leq L_{f} \|y_{1} - y_{2}\|,$$
(17)

where the first inequality holds by Jensen's inequality, and the last inequality holds by Assumption 1.

$$\begin{aligned} & \text{Since } F(x) = \frac{1}{M} \sum_{m=1}^{M} f^m(g^m(x)), \text{ we have} \\ & \|\nabla F(x_1) - \nabla F(x_2)\|^2 \\ & = \left\|\frac{1}{M} \sum_{m=1}^{M} (\nabla g^m(x_1))^T \nabla f^m(g^m(x_1)) - \frac{1}{M} \sum_{m=1}^{M} (\nabla g^m(x_2))^T \nabla f^m(g^m(x_2))\right\|^2 \\ & \leq \frac{1}{M} \sum_{m=1}^{M} \|(\nabla g^m(x_1))^T \nabla f^m(g^m(x_1)) - (\nabla g^m(x_2))^T \nabla f^m(g^m(x_2))\|^2 \\ & \leq \frac{1}{M} \sum_{m=1}^{M} \|(\nabla g^m(x_1))^T \nabla f^m(g^m(x_1)) - (\nabla g^m(x_2))^T \nabla f^m(g^m(x_1)) + (\nabla g^m(x_2))^T \nabla f^m(g^m(x_1))) \\ & - (\nabla g^m(x_2))^T \nabla f^m(g^m(x_2))\|^2 \\ & \leq \frac{1}{M} \sum_{m=1}^{M} 2C_f^2 \|\nabla g^m(x_1) - \nabla g^m(x_2)\|^2 + \frac{1}{M} \sum_{m=1}^{M} 2C_g^2 \|\nabla f^m(g^m(x_1)) - \nabla f^m(g^m(x_2))\|^2 \\ & \leq 2C_f^2 L_g^2 \|x_1 - x_2\|^2 + 2C_g^4 L_f^2 \|x_1 - x_2\|^2 = (2C_f^2 L_g^2 + 2C_g^4 L_f^2) \|x_1 - x_2\|^2, \end{aligned}$$
where the second last and the last inequalities hold by Assumptions 1-2.

**Lemma 5.** (*Restatement of Lemma 2*) Assume the gradient estimator  $\{\bar{w}_t\}_{t=1}^T$  generated from Algorithm 1, where  $w_t = \frac{1}{M} \sum_{m=1}^M w_t^m$ , we have

$$\|\bar{w}_t - \nabla F(\bar{x}_t)\|^2 \le \frac{1}{M} \sum_{m=1}^M \left( 2C_f^2 \|u_t^m - \nabla g^m(\bar{x}_t)\|^2 + 4C_g^2 \|v_t^m - \nabla f^m(h_t^m)\|^2 + 4C_g^2 L_f^2 \|h_t^m - g^m(\bar{x}_t)\|^2 \right)$$
(19)

*Proof.* Since  $\bar{w}_t = \frac{1}{M} \sum_{m=1}^{M} (u_t^m)^T v_t^m$ , we have

$$\begin{split} \|\bar{w}_{t} - \nabla F(\bar{x}_{t})\|^{2} \\ &= \|\frac{1}{M} \sum_{m=1}^{M} (u_{t}^{m})^{T} v_{t}^{m} - \frac{1}{M} \sum_{m=1}^{M} (\nabla g^{m}(\bar{x}_{t}))^{T} \nabla f^{m}(g^{m}(\bar{x}_{t}))\|^{2} \\ &= \|\frac{1}{M} \sum_{m=1}^{M} (u_{t}^{m})^{T} v_{t}^{m} - \frac{1}{M} \sum_{m=1}^{M} (\nabla g^{m}(\bar{x}_{t}))^{T} v_{t}^{m} + \frac{1}{M} \sum_{m=1}^{M} (\nabla g^{m}(\bar{x}_{t}))^{T} v_{t}^{m} - \frac{1}{M} \sum_{m=1}^{M} (\nabla g^{m}(\bar{x}_{t}))^{T} \nabla f^{m}(g^{m}(\bar{x}_{t}))\|^{2} \\ &\leq \frac{1}{M} \sum_{m=1}^{M} 2C_{f}^{2} \|u_{t}^{m} - \nabla g^{m}(\bar{x}_{t})\|^{2} + \frac{1}{M} \sum_{m=1}^{M} 2C_{g}^{2} \|v_{t}^{m} - \nabla f^{m}(g^{m}(\bar{x}_{t}))\|^{2} \\ &= \frac{2C_{f}^{2}}{M} \sum_{m=1}^{M} \|u_{t}^{m} - \nabla g^{m}(\bar{x}_{t})\|^{2} + \frac{2C_{g}^{2}}{M} \sum_{m=1}^{M} \|v_{t}^{m} - \nabla f^{m}(h_{t}^{m}) + \nabla f^{m}(h_{t}^{m}) - \nabla f^{m}(g^{m}(\bar{x}_{t}))\|^{2} \\ &\leq \frac{2C_{f}^{2}}{M} \sum_{m=1}^{M} \|u_{t}^{m} - \nabla g^{m}(\bar{x}_{t})\|^{2} + \frac{4C_{g}^{2}}{M} \sum_{m=1}^{M} \|v_{t}^{m} - \nabla f^{m}(h_{t}^{m})\|^{2} + \frac{4C_{g}^{2}L_{f}^{2}}{M} \sum_{m=1}^{M} \|h_{t}^{m} - g^{m}(\bar{x}_{t})\|^{2}, \end{split}$$

$$\tag{20}$$

where the first inequality is due to Assumptions 1-2 and the above Lemma 1.

**Lemma 6.** Suppose that the sequence  $\{\bar{x}_t\}_{t=1}^T$  be generated from Algorithm 1, where  $\bar{x}_t = \frac{1}{M} \sum_{m=1}^M x_t^m$ . Let  $0 < \gamma \le \frac{\rho}{2L\eta_t}$ , then we have

$$F(\bar{x}_{t+1}) \leq F(\bar{x}_t) + \frac{1}{M} \sum_{m=1}^{M} \left( \frac{2C_f^2 \eta_t \gamma}{\rho} \| u_t^m - \nabla g^m(\bar{x}_t) \|^2 + \frac{4C_g^2 \eta_t \gamma}{\rho} \| v_t^m - \nabla f^m(h_t^m) \|^2 + \frac{4C_g^2 L_f^2 \eta_t \gamma}{\rho} \| h_t^m - g^m(\bar{x}_t) \|^2 \right) - \frac{\rho}{2\eta_t \gamma} \| \bar{x}_{t+1} - \bar{x}_t \|^2.$$
(21)

 *Proof.* According to the above Lemma 4, the function F(x) is L-smooth. Thus we have

$$F(\bar{x}_{t+1}) \leq F(\bar{x}_t) + \langle \nabla F(\bar{x}_t), \bar{x}_{t+1} - \bar{x}_t \rangle + \frac{L}{2} \|\bar{x}_{t+1} - \bar{x}_t\|^2$$

$$= F(\bar{x}_t) + \underbrace{\langle \bar{w}_t, \bar{x}_{t+1} - \bar{x}_t \rangle}_{=T_1} + \underbrace{\langle \nabla F(\bar{x}_t) - \bar{w}_t, \bar{x}_{t+1} - \bar{x}_t \rangle}_{=T_2} + \frac{L}{2} \|\bar{x}_{t+1} - \bar{x}_t\|^2.$$
(22)

According to Assumption 5, i.e.,  $A_t \succ \rho I_d$  for any  $t \ge 1$ , the mirror function  $\varphi_t(x) = \frac{1}{2}x^T A_t x$  is  $\rho$ -strongly convex, then we can define a Bregman distance as in Ghadimi et al. (2016),

$$D_t(x,\bar{x}_t) = \varphi_t(x) - \left[\varphi_t(\bar{x}_t) + \langle \nabla \varphi_t(\bar{x}_t), x - \bar{x}_t \rangle\right] = \frac{1}{2} (x - \bar{x}_t)^T A_t(x - \bar{x}_t).$$
(23)

When  $t = s_t = q \lfloor t/q \rfloor + 1$ , according to the line 7 of Algorithm 1, we have  $\bar{x}_{t+1} = \arg \min_{x \in \mathbb{R}^d} \left\{ \langle \bar{w}_t, x \rangle + \frac{1}{2\eta_t \gamma} (x - \bar{x}_t)^T A_t (x - \bar{x}_t) \right\}$ . By using Lemma 1 in Ghadimi et al. (2016) to the problem  $\bar{x}_{t+1} = \arg \min_{x \in \mathbb{R}^d} \left\{ \langle \bar{w}_t, x \rangle + \frac{1}{2\eta_t \gamma} (x - \bar{x}_t)^T A_t (x - \bar{x}_t)^T A_t (x - \bar{x}_t) \right\}$ , we can obtain

$$\langle \bar{w}_t, \frac{1}{\eta_t \gamma} (\bar{x}_t - \bar{x}_{t+1}) \rangle \ge \rho \| \frac{1}{\eta_t \gamma} (\bar{x}_t - \bar{x}_{t+1}) \|^2.$$
 (24)

When  $t \in (s_t, s_t + q)$ , according to the line 11 of Algorithm 1, we have  $x_{t+1}^m = \arg \min_{x \in \mathbb{R}^d} \left\{ \langle w_t^m, x \rangle + \frac{1}{2\eta_t \gamma} (x - x_t^m)^T A_t (x - x_t^m) \right\}$ . Similarly, we have

$$\langle w_t^m, \frac{1}{\eta_t \gamma} (x_t^m - x_{t+1}^m) \rangle \ge \rho \| \frac{1}{\eta_t \gamma} (x_t^m - x_{t+1}^m) \|^2.$$
 (25)

Then we have

$$\frac{1}{M} \sum_{m=1}^{M} \langle w_t^m, \frac{1}{\eta_t \gamma} (x_t^m - x_{t+1}^m) \rangle \ge \rho \frac{1}{M} \sum_{m=1}^{M} \| \frac{1}{\eta_t \gamma} (x_t^m - x_{t+1}^m) \|^2 \\
\ge \rho \| \frac{1}{\eta_t \gamma} \frac{1}{M} \sum_{m=1}^{M} (x_t^m - x_{t+1}^m) \|^2 = \rho \| \frac{1}{\eta_t \gamma} (\bar{x}_t - \bar{x}_{t+1}) \|^2.$$
(26)

Thus we have

$$\langle w_t^m, \frac{1}{\eta_t \gamma} (\bar{x}_t - \bar{x}_{t+1}) \rangle \ge \rho \| \frac{1}{\eta_t \gamma} (\bar{x}_t - \bar{x}_{t+1}) \|^2.$$
 (27)

Since  $\bar{w}_t = \frac{1}{M} \sum_{m=1}^{M} w_t^m$ , averaging the above inequality (27) from m = 1 to M, we can obtain

$$\langle \bar{w}_t, \frac{1}{\eta_t \gamma} (\bar{x}_t - \bar{x}_{t+1}) \rangle = \frac{1}{M} \sum_{m=1}^M \langle w_t^m, \frac{1}{\eta_t \gamma} (\bar{x}_t - \bar{x}_{t+1}) \rangle$$

$$\geq \rho \frac{1}{M} \sum_{m=1}^M \| \frac{1}{\eta_t \gamma} (\bar{x}_t - \bar{x}_{t+1}) \|^2 = \rho \| \frac{1}{\eta_t \gamma} (\bar{x}_t - \bar{x}_{t+1}) \|^2.$$
(28)

 Then we have for any  $t \in [s_t, s_t + q)$ ,

$$T_1 = \langle \bar{w}_t, \bar{x}_{t+1} - \bar{x}_t \rangle \le -\frac{\rho}{\eta_t \gamma} \| \bar{x}_{t+1} - \bar{x}_t \|^2.$$
(29)

Since  $s_t = q \lfloor t/q \rfloor + 1$  and all  $t \in [s_t, s_t + q)$ , clearly, we have, for all  $t \ge 1$ 

$$T_1 = \langle \bar{w}_t, \bar{x}_{t+1} - \bar{x}_t \rangle \le -\frac{\rho}{\eta_t \gamma} \| \bar{x}_{t+1} - \bar{x}_t \|^2.$$
(30)

Next, consider the bound of the term  $T_2$ , we have

$$T_2 = \langle \nabla F(\bar{x}_t) - \bar{w}_t, \bar{x}_{t+1} - \bar{x}_t \rangle$$

$$\leq \|\nabla F(\bar{x}_{t}) - \bar{w}_{t}\| \cdot \|\bar{x}_{t+1} - \bar{x}_{t}\| \\\leq \frac{\eta_{t}\gamma}{\rho} \|\nabla F(\bar{x}_{t}) - \bar{w}_{t}\|^{2} + \frac{\rho}{4\eta_{t}\gamma} \|\bar{x}_{t+1} - \bar{x}_{t}\|^{2},$$
(31)

where the first inequality is due to the Cauchy-Schwarz inequality and the last is due to Young's inequality. By combining the above inequalities (22), (30) with (31), we obtain

$$F(\bar{x}_{t+1}) \leq F(\bar{x}_{t}) + \langle \nabla F(\bar{x}_{t}) - \bar{w}_{t}, \bar{x}_{t+1} - \bar{x}_{t} \rangle + \langle \bar{w}_{t}, \bar{x}_{t+1} - \bar{x}_{t} \rangle + \frac{L}{2} \| \bar{x}_{t+1} - \bar{x}_{t} \|^{2}$$

$$\leq F(\bar{x}_{t}) + \frac{\eta_{t}\gamma}{\rho} \| \nabla F(\bar{x}_{t}) - \bar{w}_{t} \|^{2} + \frac{\rho}{4\eta_{t}\gamma} \| \bar{x}_{t+1} - \bar{x}_{t} \|^{2} - \frac{\rho}{\eta_{t}\gamma} \| \bar{x}_{t+1} - \bar{x}_{t} \|^{2} + \frac{L}{2} \| \bar{x}_{t+1} - \bar{x}_{t} \|^{2}$$

$$= F(\bar{x}_{t}) + \frac{\eta_{t}\gamma}{\rho} \| \nabla F(\bar{x}_{t}) - \bar{w}_{t} \|^{2} - \frac{\rho}{2\eta_{t}\gamma} \| \bar{x}_{t+1} - \bar{x}_{t} \|^{2} - (\frac{\rho}{4\eta_{t}\gamma} - \frac{L}{2}) \| \bar{x}_{t+1} - \bar{x}_{t} \|^{2}$$

$$\leq F(\bar{x}_{t}) + \frac{\eta_{t}\gamma}{\rho} \| \nabla F(\bar{x}_{t}) - \bar{w}_{t} \|^{2} - \frac{\rho}{2\eta_{t}\gamma} \| \bar{x}_{t+1} - \bar{x}_{t} \|^{2}$$

$$\leq F(\bar{x}_{t}) + \frac{1}{M} \sum_{m=1}^{M} \left( \frac{2C_{f}^{2}\eta_{t}\gamma}{\rho} \| u_{t}^{m} - \nabla g^{m}(\bar{x}_{t}) \|^{2} + \frac{4C_{g}^{2}\eta_{t}\gamma}{\rho} \| v_{t}^{m} - \nabla f^{m}(h_{t}^{m}) \|^{2} + \frac{4C_{g}^{2}L_{f}^{2}\eta_{t}\gamma}{\rho} \| h_{t}^{m} - g^{m}(\bar{x}_{t}) \|^{2} \right) - \frac{\rho}{2\eta_{t}\gamma} \| \bar{x}_{t+1} - \bar{x}_{t} \|^{2}, \qquad (32)$$

where the second last inequality is due to  $0 < \gamma \leq \frac{\rho}{2L\eta_t}$ , and the last inequality holds by Lemma 5.

**Lemma 7.** Under the above assumptions, and assume the stochastic gradient estimators  $\{h_t^m, u_t^m, v_t^m\}_{t=1}^T$  be generated from Algorithm 1, we have, for any  $m \in [M]$ 

$$\mathbb{E}\|h_{t+1}^m - g^m(x_{t+1}^m)\|^2 \le (1 - \alpha_{t+1})\mathbb{E}\|h_t^m - g^m(x_t^m)\|^2 + 2\alpha_{t+1}^2\sigma^2 + 2C_g^2\mathbb{E}\|x_{t+1}^m - x_t^m\|^2,$$
(33)

$$\mathbb{E}\|u_{t+1}^m - \nabla g^m(x_{t+1}^m)\|^2 \le (1 - \beta_{t+1})\mathbb{E}\|u_t^m - \nabla g^m(x_t^m)\|^2 + 2\beta_{t+1}^2 \sigma^2 + 2L_g^2 \mathbb{E}\|x_{t+1}^m - x_t^m\|^2.$$
(34)

$$\mathbb{E} \|v_{t+1}^m - \nabla f^m(h_{t+1}^m)\|^2 \le (1 - \varrho_{t+1}) \mathbb{E} \|v_t^m - \nabla f^m(h_t^m)\|^2 + 4L_f^2 C_g^2 \mathbb{E} \|x_{t+1}^m - x_t^m\|^2 + 2\varrho_{t+1}^2 \sigma^2 + 8\alpha_{t+1}^2 L_f^2 \mathbb{E} \|h_t^m - g^m(x_t^m)\|^2 + 8L_f^2 \alpha_{t+1}^2 \sigma^2, \quad (35)$$

*Proof.* Without loss of generality, we only prove the above inequality (35), and it is similar to the other inequalities. Since  $v_{t+1}^m = \prod_{C_f} \left[ \nabla f^m(h_{t+1}^m; \xi_{t+1}^m) + (1 - \rho_{t+1}) \left( v_t^m - \nabla f^m(h_t^m; \xi_{t+1}^m) \right) \right]$ , we

864	have
865 866	$\mathbb{E} \  v_{t+1}^m - \nabla f^m (h_{t+1}^m) \ ^2$
867	$= \mathbb{E} \left\  \Pi_{C_{t}} \left[ \nabla f^{m}(h_{t+1}^{m}; \mathcal{E}_{t+1}^{m}) + (1 - \rho_{t+1}) (v_{t}^{m} - \nabla f^{m}(h_{t}^{m}; \mathcal{E}_{t+1}^{m})) \right] - \Pi_{C_{t}} \left[ \nabla f^{m}(h_{t+1}^{m}) \right] \right\ ^{2}$
868	$ = \sum_{i=1}^{m} \left[ \sum_{j=1}^{m} (h_{i+1}^{m}; \xi_{i+1}^{m}) + (1 - \rho_{i+1}) (p_{i}^{m} - \nabla f^{m}(h_{i}^{m}; \xi_{i+1}^{m})) - \nabla f^{m}(h_{i+1}^{m}) \right]^{2} $
869	$= \mathbb{E} \left[ (1 - o_{t+1})(v_t^m - \nabla f^m(h_t^m)) - o_{t+1}(\nabla f^m(h_{t+1}^m) - \nabla f^m(h_{t+1}^m)) \right]$
871	$ = \left[ \left( 1 - c_{t} \right) \left( \nabla f^{m}(h^{m}, f^{m}) - \nabla f^{m}(h^{m}, f^{m}) - \nabla f^{m}(h^{m}, f^{m}) - \nabla f^{m}(h^{m}, h^{m}) \right) \right]^{2} $
872	$+ (1 - \varrho_{t+1}) \left( \sqrt{f^{(n)}(n_{t+1}; \xi_{t+1})} - \sqrt{f^{(n)}(n_t; \xi_{t+1})} - \sqrt{f^{(n)}(n_{t+1})} + \sqrt{f^{(n)}(n_t)} \right) \ $
873	$= (1 - \varrho_{t+1})^2 \mathbb{E} \ v_t^m - \nabla f^m(h_t^m)\ ^2 + \mathbb{E} \ \varrho_{t+1}(\nabla f^m(h_{t+1}^m) - \nabla f^m(h_{t+1}^m;\xi_{t+1}^m))$
874 875	$-(1-\varrho_{t+1})\left(\nabla f^m(h_{t+1}^m;\xi_{t+1}^m)-\nabla f^m(h_t^m;\xi_{t+1}^m)-\nabla f^m(h_{t+1}^m)+\nabla f^m(h_t^m)\right)\right)^2$
876	$\leq (1 - \varrho_{t+1})^2 \mathbb{E} \  v_t^m - \nabla f^m(h_t^m) \ ^2 + 2\varrho_{t+1}^2 \mathbb{E} \  \nabla f^m(h_{t+1}^m) - \nabla f^m(h_{t+1}^m; \xi_{t+1}^m) \ ^2$
877	$+2(1-\varrho_{t+1})^2 \left\ \nabla f^m(h_{t+1}^m;\xi_{t+1}^m) - \nabla f^m(h_t^m;\xi_{t+1}^m) - \nabla f^m(h_{t+1}^m) + \nabla f^m(h_t^m)\right\ ^2$
878 879	$\leq (1 - \varrho_{t+1})^2 \mathbb{E} \ v_t^m - \nabla f^m(h_t^m)\ ^2 + 2\varrho_{t+1}^2 \sigma^2 + 2(1 - \varrho_{t+1})^2 \ \nabla f^m(h_{t+1}^m; \xi_{t+1}^m) - \nabla f^m(h_t^m; \xi_{t+1}^m)\ ^2$
880	$\leq (1 - \rho_{t+1})^2 \mathbb{E} \  v_t^m - \nabla f^m(h_t^m) \ ^2 + 2\rho_{t+1}^2 \sigma^2 + 2(1 - \rho_{t+1})^2 L_f^2 \mathbb{E} \  h_{t+1}^m - h_t^m \ ^2, \tag{36}$
881	where the third equality holds by the following fact:
882	$\mathbb{E}_{\ell^{m}}\left[\rho_{t+1}(\nabla f^{m}(h_{t+1}^{m}) - \nabla f^{m}(h_{t+1}^{m};\xi_{t+1}^{m})) - (1 - \rho_{t+1})(\nabla f^{m}(h_{t+1}^{m};\xi_{t+1}^{m}) - \nabla f^{m}(h_{t}^{m};\xi_{t+1}^{m})\right]$
884	$-\nabla f^{m}(h^{m}_{t+1}) + \nabla f^{m}(h^{m}_{t+1})] = 0.$
885	and the second last inequality holds by the inequality $\mathbb{E}\ \mathcal{L} - \mathbb{E}[\mathcal{L}]\ ^2 < \mathbb{E}\ \mathcal{L}\ ^2$ and Assumption 3:
886	the last inequality is due to Assumption 1.
888	Since $h_{t+1}^m = g^m(x_{t+1}^m; \zeta_{t+1}^m) + (1 - \alpha_{t+1}) (h_t^m - g^m(x_t^m; \zeta_{t+1}^m))$ , we have
889	$\mathbb{E}\ h_{t+1}^m - h_t^m\ ^2 = \mathbb{E}\ q^m(x_{t+1}^m;\zeta_{t+1}^m) - q^m(x_t^m;\zeta_{t+1}^m) - \alpha_{t+1}(h_t^m - q^m(x_t^m;\zeta_{t+1}^m))\ ^2$
890	$< 2\mathbb{E} \ q^{m}(x_{t+1}^{m};\zeta_{t+1}^{m}) - q^{m}(x_{t}^{m};\zeta_{t+1}^{m})\ ^{2} + 2\alpha_{t+1}^{2}\mathbb{E} \ h_{t}^{m} - q^{m}(x_{t}^{m};\zeta_{t+1}^{m})\ ^{2}$
891 892	$\leq 2C_{+}^{2} \ x_{++}^{m} - x_{+}^{m}\ ^{2} + 2\alpha_{++}^{2} \mathbb{E}\ h_{+}^{m} - a^{m}(x_{++}^{m};\zeta_{+++}^{m})\ ^{2}$
893	$= 2C^2 \ r^m_{t+1} - r^m_{t+1}\ ^2 + 2\alpha_{t+1}^2 \ \mathbb{E}\ h^m - a^m(r^m_{t+1})\ ^2 + a^m(r^m_{t+1}) - a^m(r^m_{t+1})\ ^2$
894	$= 2C_{g} \ w_{t+1} - w_{t}\  + 2\alpha_{t+1} \ w_{t} - g(w_{t}, \varsigma_{t+1}) + g(w_{t}) - g(w_{t}) - g(w_{t}) \  $ $\leq 2C^{2} \ w_{t} - w_{t}\ ^{2} + 4\alpha^{2} \ w_{t}\  + 4\alpha^{2} \ w_{t} - w_{t}\ ^{2} + 4\alpha^{2} \ w_{t}\ ^{2} + 4\alpha^{2} \ w_{t} - w_{t}\ ^{2} + 4\alpha^{2} \ w_{t}\ ^$
895 896	$\leq 2C_g \ x_{t+1} - x_t\  + 4\alpha_{t+1} \ u_t - g(x_t)\  + 4\alpha_{t+1} \ y(x_t, \zeta_{t+1}) + g(x_t)\  $ $\leq 2C_g \ x_{t+1} - x_t\  + 4\alpha_{t+1} \ u_t - g(x_t)\  + 4\alpha_{t+1} \ y(x_t, \zeta_{t+1}) + g(x_t)\  $ $\leq 2C_g \ x_{t+1} - x_t\  + 4\alpha_{t+1} \ u_t - g(x_t)\  + 4\alpha_{t+1} \ u_t - g(x$
897	$\leq 2C_g \ x_{t+1} - x_t\  + 4\alpha_{t+1}\mathbb{E}\ n_t - g\ (x_t)\  + 4\alpha_{t+1}\delta , \qquad (57)$
898	where the second inequality holds by Assumption , .
899	Combining the above inequalities (36) with (37), we have $\overline{m} = \frac{m}{m} \frac{m}$
900 901	$\mathbb{E}\ v_{t+1}^{m} - \nabla f^{m}(h_{t+1}^{m})\ ^{2}$
902	$\leq (1 - \varrho_{t+1})^2 \mathbb{E} \ v_t^m - \nabla f^m(h_t^m)\ ^2 + 2\varrho_{t+1}^2 \sigma^2 + 2(1 - \varrho_{t+1})^2 L_f^2 \mathbb{E} \ h_{t+1}^m - h_t^m\ ^2$
903	$\leq (1 - \varrho_{t+1}) \mathbb{E} \  v_t^m - \nabla f^m(h_t^m) \ ^2 + 2\varrho_{t+1}^2 \sigma^2 + 4L_f^2 C_g^2 \  x_{t+1}^m - x_t^m \ ^2$
904	$+ 8\alpha_{t+1}^2 L_f^2 \mathbb{E} \ h_t^m - g^m(x_t^m)\ ^2 + 8L_f^2 \alpha_{t+1}^2 \sigma^2,$
905 906	where the last inequality holds by $0 < \varrho_{t+1} \le 1$ .
907	
908	<b>Lemma 8.</b> Based on the above Assumptions 1-2 and 6, we have
909	M 1 M M
310	$\sum \mathbb{E} \left\  \nabla m(1m) - \frac{1}{2} \sum \nabla r(i(1i)) \right\ ^2 < 0.12 \sum \mathbb{E} \left\  1m - m(-1) \right\ ^2 + 4MS^2 + 4MS^2 + 4MS^2$

$$\begin{split} &\sum_{m=1}^{M} \mathbb{E} \left\| \nabla f^{m}(h_{t}^{m}) - \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(h_{t}^{j}) \right\|^{2} \leq 8L_{f}^{2} \sum_{m=1}^{M} \mathbb{E} \|h_{t}^{m} - g^{m}(\bar{x}_{t})\|^{2} + 4M\delta_{f}^{2} + 4ML_{f}^{2}\delta_{g}^{2}, \\ &\sum_{m=1}^{M} \mathbb{E} \left\| \nabla g^{m}(x_{t}^{m}) - \frac{1}{M} \sum_{j=1}^{M} \nabla g^{j}(x_{t}^{j}) \right\|^{2} \leq 6L_{g}^{2} \sum_{m=1}^{M} \mathbb{E} \|x_{t}^{m} - \bar{x}_{t}\|^{2} + 3M\delta_{g}^{2} \\ &\sum_{m=1}^{M} \mathbb{E} \left\| g^{m}(x_{t}^{m}) - \frac{1}{M} \sum_{j=1}^{M} g^{j}(x_{t}^{j}) \right\|^{2} \leq 6C_{g}^{2} \sum_{m=1}^{M} \mathbb{E} \|x_{t}^{m} - \bar{x}_{t}\|^{2} + 3M\delta_{g}^{2}. \end{split}$$

$$\begin{aligned} Proof. \ \text{Consider the term } & \sum_{m=1}^{M} \mathbb{E} \| \nabla f^{m}(h_{t}^{m}) - \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(h_{t}^{j}) \|^{2}, \text{ we have} \\ & \sum_{m=1}^{M} \mathbb{E} \| \nabla f^{m}(h_{t}^{m}) - \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(h_{t}^{j}) \|^{2} \\ & = \sum_{m=1}^{M} \mathbb{E} \| \nabla f^{m}(h_{t}^{m}) - \nabla f^{m}(g^{m}(\bar{x}_{t})) + \nabla f^{m}(g^{m}(\bar{x}_{t})) - \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(g^{m}(\bar{x}_{t})) + \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(g^{m}(\bar{x}_{t})) \\ & - \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(g^{j}(\bar{x}_{t})) + \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(g^{j}(\bar{x}_{t})) - \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(h_{t}^{j}) \|^{2} \\ & \leq \sum_{m=1}^{M} 4\mathbb{E} \| \nabla f^{m}(h_{t}^{m}) - \nabla f^{m}(g^{m}(\bar{x}_{t})) \|^{2} + \sum_{m=1}^{M} 4\mathbb{E} \| \nabla f^{m}(g^{m}(\bar{x}_{t})) - \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(g^{m}(\bar{x}_{t})) \|^{2} \\ & + \sum_{m=1}^{M} 4\mathbb{E} \| \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(g^{m}(\bar{x}_{t})) - \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(g^{j}(\bar{x}_{t})) \|^{2} + \sum_{m=1}^{M} 4\mathbb{E} \| \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(g^{j}(\bar{x}_{t})) - \frac{1}{M} \sum_{j=1}^{M} \nabla f^{j}(h_{t}^{j}) \|^{2} \\ & \leq 4L_{f}^{2} \sum_{m=1}^{M} \mathbb{E} \| h_{t}^{m} - g^{m}(\bar{x}_{t}) \|^{2} + 4 \sum_{m=1}^{M} \frac{1}{M} \sum_{j=1}^{M} \mathbb{E} \| \nabla f^{m}(g^{m}(\bar{x}_{t})) - \nabla f^{j}(g^{m}(\bar{x}_{t})) \|^{2} \\ & + 4L_{f}^{2} \sum_{m=1}^{M} \frac{1}{M} \sum_{j=1}^{M} \| g^{m}(\bar{x}_{t}) - g^{j}(\bar{x}_{t}) \|^{2} + 4L_{f}^{2} \sum_{j=1}^{M} \frac{1}{M} \sum_{m=1}^{M} \mathbb{E} \| g^{m}(\bar{x}_{t}) - h_{t}^{m} \|^{2} \\ & \leq 8L_{f}^{2} \sum_{m=1}^{M} \frac{1}{M} \sum_{j=1}^{M} \| g^{m}(\bar{x}_{t}) - g^{j}(\bar{x}_{t}) \|^{2} + 4L_{f}^{2} \sum_{j=1}^{M} \frac{1}{M} \sum_{m=1}^{M} \frac{1}{M} \sum_{m=1}^{M} \mathbb{E} \| g^{m}(\bar{x}_{t}) - h_{t}^{m} \|^{2} \\ & \leq 8L_{f}^{2} \sum_{m=1}^{M} \mathbb{E} \| h_{t}^{m} - g^{m}(\bar{x}_{t}) \|^{2} + 4M\delta_{f}^{2} + 4ML_{f}^{2}\delta_{g}^{2}, \end{aligned} \tag{38} \\ & \text{where the last inequality holds by Assumption 6} \end{aligned}$$

where the last inequality holds by Assumption 6.

Next, we have

$$\begin{split} &\sum_{m=1}^{M} \mathbb{E} \left\| \nabla g^{m}(x_{t}^{m}) - \frac{1}{M} \sum_{j=1}^{M} \nabla g^{j}(x_{t}^{j}) \right\|^{2} \\ &= \sum_{m=1}^{M} \mathbb{E} \left\| \nabla g^{m}(x_{t}^{m}) - \nabla g^{m}(\bar{x}_{t}) + \nabla g^{m}(\bar{x}_{t}) - \frac{1}{M} \sum_{j=1}^{M} \nabla g^{j}(\bar{x}_{t}) + \frac{1}{M} \sum_{j=1}^{M} \nabla g^{j}(\bar{x}_{t}) - \frac{1}{M} \sum_{j=1}^{M} \nabla g^{j}(x_{t}^{j}) \right\|^{2} \\ &\leq \sum_{m=1}^{M} 3\mathbb{E} \left\| \nabla g^{m}(x_{t}^{m}) - \nabla g^{m}(\bar{x}_{t}) \right\|^{2} + \sum_{m=1}^{M} 3\mathbb{E} \left\| \nabla g^{m}(\bar{x}_{t}) - \frac{1}{M} \sum_{j=1}^{M} \nabla g^{j}(\bar{x}_{t}) \right\|^{2} \\ &+ \sum_{m=1}^{M} 3\mathbb{E} \left\| \frac{1}{M} \sum_{j=1}^{M} \nabla g^{j}(\bar{x}_{t}) - \frac{1}{M} \sum_{j=1}^{M} \nabla g^{j}(x_{t}^{j}) \right\|^{2} \\ &\leq 3L_{g}^{2} \sum_{m=1}^{M} \mathbb{E} \| x_{t}^{m} - \bar{x}_{t} \|^{2} + 3 \sum_{m=1}^{M} \frac{1}{M} \sum_{j=1}^{M} \mathbb{E} \| \nabla g^{m}(\bar{x}_{t}) - \nabla g^{j}(\bar{x}_{t}) \|^{2} + 3 \sum_{m=1}^{M} \frac{1}{M} \sum_{j=1}^{M} \| \nabla g^{j}(\bar{x}_{t}) - \nabla g^{j}(x_{t}^{j}) \|^{2} \\ &\leq 6L_{g}^{2} \sum_{m=1}^{M} \mathbb{E} \| x_{t}^{m} - \bar{x}_{t} \|^{2} + 3M\delta_{g}^{2}, \end{split}$$

$$\tag{39}$$

where the last inequality is due to the above Assumption 6.

Similarly, we can obtain

$$\sum_{m=1}^{M} \mathbb{E} \left\| g^{m}(x_{t}^{m}) - \frac{1}{M} \sum_{j=1}^{M} g^{j}(x_{t}^{j}) \right\|^{2} \le 6C_{g}^{2} \sum_{m=1}^{M} \mathbb{E} \|x_{t}^{m} - \bar{x}_{t}\|^{2} + 3M\delta_{g}^{2}.$$
(40)

**Lemma 9.** Suppose the iterates  $\{x_t^m\}_{t=1}^T$ , for all  $m \in [M]$  generated from Algorithm 1 satisfy: 

$$\sum_{m=1}^{M} \mathbb{E} \|x_t^m - \bar{x}_t\|^2 \le (q-1) \sum_{l=s_t}^{t-1} \gamma^2 \eta_l^2 \sum_{m=1}^{M} \mathbb{E} \|d_l^m - \bar{d}_l\|^2,$$
(41)

where  $\bar{x}_t = \frac{1}{M} \sum_{m=1}^{M} x_t^m$ ,  $d_t^m = \frac{x_t^m - x_{t+1}^m}{\gamma \eta_t}$  and  $\bar{d}_t = \frac{\bar{x}_t - \bar{x}_{t+1}}{\eta_t \gamma}$ .

Proof. According to the lines 7 and 11 of Algorithm 1, we have

$$x_{t+1}^{m} = x_{t}^{m} - \gamma \eta_{t} A_{t}^{-1} w_{t}^{m} = \arg\min_{x \in \mathbb{R}^{d}} \left\{ \langle w_{t}^{m}, x \rangle + \frac{1}{2\eta_{t}\gamma} (x - x_{t}^{m})^{T} A_{t} (x - x_{t}^{m}) \right\},\\ \bar{x}_{t+1} = \bar{x}_{t} - \gamma \eta_{t} A_{t}^{-1} \bar{w}_{t} = \arg\min_{x \in \mathbb{R}^{d}} \left\{ \langle \bar{w}_{t}, x \rangle + \frac{1}{2\eta_{t}\gamma} (x - \bar{x}_{t})^{T} A_{t} (x - \bar{x}_{t}) \right\},$$

and then we define the gradient mappings as in the above (13):  $d_t^m = \frac{x_t^m - x_{t+1}^m}{\gamma \eta_t} = A_t^{-1} w_t^m$  and  $\bar{d}_t = \frac{\bar{x}_t - \bar{x}_{t+1}}{\eta_t \gamma} = A_t^{-1} \bar{w}_t = \frac{1}{M} \sum_{m=1}^M d_t^m \text{ for any } m \in [M] \text{ and } t \ge 1.$ 

From the line 7 of Algorithm 1, when  $t = s_t = q \lfloor t/q \rfloor + 1$ , we have  $x_t^m = \bar{x}_t = \frac{1}{M} \sum_{m=1}^M x_t^m$  for any  $m \in [M]$ , so the about inequality in the lemma holds trivially.

When  $t \in (s_t, s_t + q)$ , we have

$$x_{t}^{m} = x_{s_{t}}^{m} - \sum_{l=s_{t}}^{t-1} \gamma \eta_{l} d_{l}^{m}, \quad \text{and} \quad \bar{x}_{t} = \bar{x}_{s_{t}} - \sum_{l=s_{t}}^{t-1} \gamma \eta_{l} \bar{d}_{l}$$

Thus we have

$$\sum_{m=1}^{M} \mathbb{E} \|x_{t}^{m} - \bar{x}_{t}\|^{2} = \sum_{m=1}^{M} \mathbb{E} \left\| x_{s_{t}}^{m} - \bar{x}_{s_{t}} - \left( \sum_{l=s_{t}}^{t-1} \gamma \eta_{l} d_{l}^{m} - \sum_{l=s_{t}}^{t-1} \gamma \eta_{l} \bar{d}_{l} \right) \right\|^{2}$$

$$= \sum_{m=1}^{M} \mathbb{E} \left\| \left( \sum_{l=s_{t}}^{t-1} \gamma \eta_{l} d_{l}^{m} - \sum_{l=s_{t}}^{t-1} \gamma \eta_{l} \bar{d}_{l} \right) \right\|^{2} \le (q-1) \sum_{l=s_{t}}^{t-1} \gamma^{2} \eta_{l}^{2} \sum_{m=1}^{M} \mathbb{E} \| d_{l}^{m} - \bar{d}_{l} \|^{2}$$

$$= \sum_{m=1}^{M} \mathbb{E} \left\| \left( \sum_{l=s_{t}}^{t-1} \gamma \eta_{l} d_{l}^{m} - \sum_{l=s_{t}}^{t-1} \gamma \eta_{l} \bar{d}_{l} \right) \right\|^{2} \le (q-1) \sum_{l=s_{t}}^{t-1} \gamma^{2} \eta_{l}^{2} \sum_{m=1}^{M} \mathbb{E} \| d_{l}^{m} - \bar{d}_{l} \|^{2}$$

where the above inequality is due to  $t - s_t \leq q - 1$ . 

**Lemma 10.** Let  $C_{fg}^2 = \max(C_f^2, C_g^2)$ ,  $L_{fg}^2 = L_f^2 C_g^2 + L_g^2$  and  $\eta_t \leq \frac{\rho}{24\gamma q L_{fg} C_{fg}}$  for all  $t \geq 0$ . Further let  $\alpha_{t+1} = c_1 \eta_t^2$ ,  $\beta_{t+1} = c_2 \eta_t^2$  and  $\varrho_{t+1} = c_3 \eta_t^2$ ,  $c_1, c_2, c_3 > 0$  and  $c_1^2 + c_2^2 \leq (\alpha_1 + \alpha_2 + \alpha_3 + \alpha_3 + \alpha_4)$  $\frac{(24)^4 q^2 \gamma^4 L_{fg}^4 C_{fg}^4}{9 o^4}$ . Set  $s_t = q \lfloor t/q \rfloor + 1$  and  $t \in [s_t, s_t + q - 1]$ , we have 

$$\sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \sum_{m=1}^{M} \mathbb{E} \|d_{t}^{m} - \bar{d}_{t}\|^{2}$$

$$\leq \frac{6M}{5} \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \mathbb{E} \|\bar{d}_{t}\|^{2} + \frac{\rho^{2}(c_{1}^{2}+c_{3}^{2})}{120q^{2}\gamma^{4}C_{fg}^{2}L_{fg}^{2}C_{g}^{2}} \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \sum_{m=1}^{M} \mathbb{E} \|h_{t}^{m} - g^{m}(\bar{x}_{t})\|^{2} + \frac{3M\hat{\delta}^{2}}{5\gamma^{2}L_{fg}^{2}} \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t}^{3},$$
(42)

where  $\hat{\delta}^2 = 2c_1^2 L_f^2 \sigma^2 + c_3^2 \sigma^2 + 4c_3^2 \delta_f^2 + 4c_3^2 L_f^2 \delta_a^2 + c_2^2 \sigma^2 + 3c_2^2 \delta_a^2$ 

*Proof.* According to the lines 7 and 11 of Algorithm 1, we have 

1022  
1023 
$$x_{t+1}^m = \arg\min_{x \in \mathbb{R}^d} \Big\{ \langle w_t^m, x \rangle + \frac{1}{2\eta_t \gamma} (x - x_t^m)^T A_t (x - x_t^m) \Big\},$$

1024  
1025 
$$\bar{x}_{t+1} = \arg\min_{x \in \mathbb{R}^d} \left\{ \langle \bar{w}_t, x \rangle + \frac{1}{2\eta_t \gamma} (x - \bar{x}_t)^T A_t (x - \bar{x}_t) \right\},$$

and then we define the gradient mappings as in the above (13):  $d_t^m = \frac{x_t^m - x_{t+1}^m}{\gamma \eta_t} = A_t^{-1} w_t^m$  and  $\bar{d}_t = \frac{\bar{x}_t - \bar{x}_{t+1}}{\eta_t \gamma} = A_t^{-1} \bar{w}_t = \frac{1}{M} \sum_{m=1}^M d_t^m$  for any  $m \in [M]$  and  $t \ge 1$ . Then we have  $\sum_{t=1}^{M} \mathbb{E} \|d_{t}^{m} - \bar{d}_{t}\|^{2} = \sum_{t=1}^{M} \mathbb{E} \|A_{t}^{-1}(w_{t}^{m} - \bar{w}_{t})\|^{2} \le \frac{1}{\rho^{2}} \sum_{t=1}^{M} \mathbb{E} \|w_{t}^{m} - \bar{w}_{t}\|^{2}$ (43) $=\frac{1}{\rho^2}\sum_{t=1}^{M}\mathbb{E}\|(u_t^m)^T - (u_t^m)^T\bar{v}_t + (u_t^m)^T\bar{v}_t - (\bar{u}_t)^T\bar{v}_t + (\bar{u}_t)^T\bar{v}_t - \frac{1}{M}\sum_{t=1}^{M}(u_t^m)^Tv_t^m\|^2$  $\leq \frac{1}{\rho^2} \sum_{t=1}^{M} \left( 3C_g^2 \mathbb{E} \|v_t^m - \bar{v}_t\|^2 + 3C_f^2 \mathbb{E} \|u_t^m - \bar{u}_t\|^2 + 3\|(\bar{u}_t)^T \bar{v}_t - \frac{1}{M} \sum_{t=1}^{M} (u_t^m)^T v_t^m\|^2 \right),$ where the first inequality holds by Assumption 5, i.e.,  $A_t \succeq \rho I_d$  for all  $t \ge 1$ , and the last inequality holds by  $||u_t^m||^2 \le C_g^2$  and  $||\bar{v}_t||^2 \le C_f^2$ . Consider the term  $||(\bar{u}_t)^T \bar{v}_t - \frac{1}{M} \sum_{m=1}^M (u_t^m)^T v_t^m ||^2$ , we have  $\|(\bar{u}_t)^T \bar{v}_t - \frac{1}{M} \sum_{m=1}^M (u_t^m)^T v_t^m \|^2 = \|(\bar{u}_t)^T \bar{v}_t - \frac{1}{M} \sum_{m=1}^M (u_t^m)^T v_t^m \|^2$  $\leq \frac{1}{M} \sum_{t=1}^{M} \|(\bar{u}_t)^T \bar{v}_t - (u_t^m)^T v_t^m\|^2$  $= \frac{1}{M} \sum_{t=1}^{M} \|(\bar{u}_t)^T \bar{v}_t - (\bar{u}_t)^T v_t^m + (\bar{u}_t)^T v_t^m - (u_t^m)^T v_t^m \|^2$  $\leq \frac{1}{M} \sum_{l=1}^{M} \left( 2C_{g}^{2} \| v_{t}^{m} - \bar{v}_{t} \|^{2} + 2C_{f}^{2} \| u_{t}^{m} - \bar{u}_{t} \|^{2} \right).$ (44)By combining the above inequalities (43) and (44), we have

$$\sum_{m=1}^{M} \mathbb{E} \|d_t^m - \bar{d}_t\|^2 \le \frac{9C_g^2}{\rho^2} \sum_{m=1}^{M} \mathbb{E} \|v_t^m - \bar{v}_t\|^2 + \frac{9C_f^2}{\rho^2} \sum_{m=1}^{M} \mathbb{E} \|u_t^m - \bar{u}_t\|^2.$$
(45)

1077 Let  $t = s_t = q \lfloor t/q \rfloor + 1$ . When  $t = s_t$ , we have  $v_t^m = \bar{v}_t$  and  $u_t^m = \bar{u}_t$  for any  $m \in [M]$ , so we have  $\sum_{m=1}^{M} \mathbb{E} \| v_t^m - \bar{v}_t \|^2 = 0$  and  $\sum_{m=1}^{M} \mathbb{E} \| u_t^m - \bar{u}_t \|^2 = 0$ . According to the above inequality (45), 1079 when  $t = s_t$ , we have  $\sum_{m=1}^{M} \mathbb{E} \| d_t^m - \bar{d}_t \|^2 = 0$ . Clearly, the about inequality (42) in the lemma holds trivially.

When 
$$t \in \{s_t, s_t + q\}$$
, we first consider the term  $\sum_{m=1}^{M} \mathbb{E} \| v_t^m - \bar{v}_t \|^2$  as follows:  

$$\int_{m=1}^{M} \mathbb{E} \| v_t^m - v_t \|^2 \qquad (46)$$

$$= \int_{m=1}^{M} \mathbb{E} \| v_t^m - \frac{1}{M} \sum_{m=1}^{M} v_t^m \|^2$$

$$= \int_{m=1}^{M} \mathbb{E} \| \| v_{t-1}^m (h_t^m; \xi_t^m) + (1 - \varrho_t) (v_{t-1}^m - \nabla f^m (h_{t-1}^m; \xi_t^m)) \| - \frac{1}{M} \sum_{m=1}^{M} \Pi_{C_T} \| \nabla f^m (h_t^m; \xi_t^m) + (1 - \varrho_t) (v_{t-1}^m - \nabla f^m (h_{t-1}^m; \xi_t^m)) \| - \frac{1}{M} \sum_{m=1}^{M} (\nabla f^m (h_t^m; \xi_t^m) + (1 - \varrho_t) (v_{t-1}^m - \nabla f^m (h_{t-1}^m; \xi_t^m)) - \frac{1}{M} \sum_{m=1}^{M} (\nabla f^m (h_t^m; \xi_t^m) + (1 - \varrho_t) (v_{t-1}^m - \nabla f^m (h_{t-1}^m; \xi_t^m)) - \frac{1}{M} \sum_{m=1}^{M} (\nabla f^m (h_t^m; \xi_t^m) + (1 - \varrho_t) (v_{t-1}^m - \nabla f^m (h_{t-1}^m; \xi_t^m)) - \frac{1}{M} \sum_{m=1}^{M} (\nabla f^m (h_t^m; \xi_t^m) + (1 - \varrho_t) (v_{t-1}^m - \nabla f^m (h_{t-1}^m; \xi_t^m)) - \frac{1}{M} \sum_{m=1}^{M} \nabla f^m (h_t^m; \xi_t^m) - \frac{1}{M} \sum_{m=1}^{M} \nabla f^m (h_t^m; \xi_t^m) - (1 - \varrho_t) (\nabla f^m (h_{t-1}^m; \xi_t^m) - \frac{1}{M} \sum_{m=1}^{M} \nabla f^m (h_t^m; \xi_t^m)) \|^2$$
Then, we consider the last term of (46):
$$\int_{m=1}^{M} \mathbb{E} \| \nabla f^m (h_t^m; \xi_t^m) - \frac{1}{M} \sum_{m=1}^{M} \nabla f^m (h_{t-1}^m; \xi_t^m) - \frac{1}{M} \sum_{m=1}^{M} (\nabla f^m (h_t^m; \xi_t^m) - \nabla f^m (h_{t-1}^m; \xi_t^m)) \|^2$$

$$= \int_{m=1}^{M} \mathbb{E} \| \nabla f^m (h_t^m; \xi_t^m) - \nabla f^m (h_{t-1}^m; \xi_t^m) - \frac{1}{M} \sum_{m=1}^{M} (\nabla f^m (h_t^m; \xi_t^m) - \nabla f^m (h_{t-1}^m; \xi_t^m)) \|^2$$

$$= \sum_{m=1}^{M} \mathbb{E} \| \nabla f^m (h_t^m; \xi_t^m) - \nabla f^m (h_{t-1}^m; \xi_t^m) - \frac{1}{M} \sum_{m=1}^{M} \nabla f^m (h_{t-1}^m; \xi_t^m) \|^2$$

$$= 2 \sum_{m=1}^{M} \mathbb{E} \| \nabla f^m (h_t^m; \xi_t^m) - \nabla f^m (h_{t-1}^m; \xi_t^m) \|^2$$

$$\leq 2 \sum_{m=1}^{M} \mathbb{E} \| \nabla f^m (h_t^m; \xi_t^m) - \nabla f^m (h_{t-1}^m; \xi_t^m) \|^2$$

$$\leq 2 \sum_{m=1}^{M} \mathbb{E} \| \nabla f^m (h_t^m; \xi_t^m) - \nabla f^m (h_{t-1}^m; \xi_t^m) \|^2$$

$$\leq 2 \sum_{m=1}^{M} \mathbb{E} \| \nabla f^m (h_t^m; \xi_t^m) - \nabla f^m (h_{t-1}^m; \xi_t^m) \|^2$$

$$\leq 2 \sum_{m=1}^{M} \mathbb{E} \| \nabla f^m (h_t^m; \xi_t^m) - \nabla f^m (h_{t-1}^m; \xi_t^m) \|^2$$

$$\leq 2 \sum_{m=1}^{M} \mathbb{E} \| \nabla f^m (h_t^m; \xi_t^m) - \nabla f^m (h_{t-1}^m; \xi_t^m) \|^2$$

$$\leq 2 \sum_{m=1}^{M} \mathbb{E} \| \| \nabla f^m (h_t^m; \xi_t^m) - \nabla f^m (h_{t-1}^m; \xi_t^m) \|^2$$

where the second last inequality is due to Young inequality and the above Lemma 2.

$$\begin{split} & \text{Consider the term } \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m};\xi_{t}^{m}) - \frac{1}{M} \sum_{m=1}^{M} \nabla f^{m}(h_{t-1}^{m};\xi_{t}^{m}) \|^{2}, \text{ we have} \\ & \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m};\xi_{t}^{m}) - \frac{1}{M} \sum_{m=1}^{M} \nabla f^{m}(h_{t-1}^{m};\xi_{t}^{m}) \|^{2} \\ & = \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m};\xi_{t}^{m}) - \nabla f^{m}(h_{t-1}^{m}) - \frac{1}{M} \sum_{m=1}^{M} (\nabla f^{m}(h_{t-1}^{m};\xi_{t}^{m}) - \nabla f^{m}(h_{t-1}^{m})) \\ & + \nabla f^{m}(h_{t-1}^{m}) - \frac{1}{M} \sum_{m=1}^{M} \nabla f^{m}(h_{t-1}^{m}) \|^{2} \\ & \leq 2 \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m};\xi_{t}^{m}) - \nabla f^{m}(h_{t-1}^{m}) \|^{2} \\ & \leq 2 \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m};\xi_{t}^{m}) - \nabla f^{m}(h_{t-1}^{m}) \|^{2} + 2 \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m}) - \frac{1}{M} \sum_{m=1}^{M} \nabla f^{m}(h_{t-1}^{m}) - \frac{1}{M} \sum_{m=1}^{M} \nabla f^{m}(h_{t-1}^{m}) \|^{2} \\ & \leq 2 \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m};\xi_{t}^{m}) - \nabla f^{m}(h_{t-1}^{m}) \|^{2} + 2 \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m}) - \frac{1}{M} \sum_{m=1}^{M} \nabla f^{m}(h_{t-1}^{m}) \|^{2} \\ & \leq 2 M \sigma^{2} + 16 L_{f}^{2} \sum_{m=1}^{M} \mathbb{E} \| h_{t-1}^{m} - g^{m}(\bar{x}_{t-1}) \|^{2} + 8 M \delta_{f}^{2} + 8 M L_{f}^{2} \delta_{g}^{2}, \qquad (48) \end{split}$$
where the last inequality holds by the above Lemma 8.
$$\text{Since } h_{t}^{m} = g^{m}(x_{t}^{m};\zeta_{t}^{m}) - (1 - \alpha_{t})(h_{t-1}^{m} - g^{m}(x_{t-1}^{m};\zeta_{t}^{m})), \text{ we have} \end{aligned}$$

$$= \mathbb{E} \| g^{m}(x_{t}^{m};\zeta_{t}^{m}) - g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) \|^{2} + 2 \Delta_{f}^{2} \mathbb{E} \| h_{t-1}^{m} - g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) \|^{2} \\ & \leq 2 M \sigma^{2} + 16 L_{f}^{2} \sum_{m=1}^{M} \mathbb{E} \| h_{t-1}^{m} - g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) \|^{2} + 2 \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m}) \|^{2} \\ & \leq 2 M \sigma^{2} + 16 L_{f}^{2} \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m}) \|^{2} + 2 \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m}) \|^{2} \\ & \leq 2 M \sigma^{2} + 16 L_{f}^{2} \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m}) \|^{2} + 2 \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m}) \|^{2} \\ & \leq 2 M \sigma^{2} + 16 L_{f}^{2} \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m}) \|^{2} + 2 \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m}) \|^{2} \\ & \leq 2 M \sigma^{2} + 16 L_{f}^{2} \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m}) \|^{2} + 2 \sum_{m=1}^{M} \| \nabla f^{m}(h_{t-1}^{m}) \|^{2} \\ & \leq 2 M \sigma^{2} \| \| \nabla f$$

By combining the above inequalities (46), (47), (48) and (49), we have

 $\sum_{t=1}^{M} \mathbb{E} \|v_t^m - \bar{v}_t\|^2$ (50) $\leq (1+\nu)(1-\varrho_t)^2 \sum_{i=1}^{M} \mathbb{E} \|v_{t-1}^m - \bar{v}_{t-1})\|^2 + (1+\frac{1}{\nu}) \sum_{i=1}^{M} \mathbb{E} \|\nabla f^m(h_t^m; \xi_t^m)\|^2$  $-\frac{1}{M}\sum_{t=1}^{M}\nabla f^{m}(h_{t}^{m};\xi_{t}^{m}) - (1-\varrho_{t})\left(\nabla f^{m}(h_{t-1}^{m};\xi_{t}^{m}) - \frac{1}{M}\sum_{t=1}^{M}\nabla f^{m}(h_{t-1}^{m};\xi_{t}^{m})\right)\Big\|^{2}$  $\leq (1+\nu)(1-\varrho_t)^2 \sum_{i=1}^{M} \mathbb{E} \|v_{t-1}^m - \bar{v}_{t-1})\|^2 + (1+\frac{1}{\nu}) \left(4L_f^2 C_g^2 \sum_{i=1}^{M} \mathbb{E} \|x_t^m - x_{t-1}^m\|^2\right) \|v_{t-1}^m - v_{t-1}^m\|^2$  $+16\alpha_t^2 L_f^2 \sum_{m=1}^M \mathbb{E}\|h_{t-1}^m - g^m(\bar{x}_{t-1})\|^2 + 16\alpha_t^2 L_f^2 C_g^2 \sum_{m=1}^M \mathbb{E}\|\bar{x}_{t-1} - x_{t-1}^m\|^2 + 8M L_f^2 \alpha_t^2 \sigma^2$  $+4M\varrho_t^2\sigma^2+32\varrho_t^2L_f^2\sum_{m=1}^M\mathbb{E}\|h_{t-1}^m-g^m(\bar{x}_{t-1})\|^2+16M\varrho_t^2\delta_f^2+16ML_f^2\varrho_t^2\delta_g^2\right)$  $\leq (1+\nu)(1-\varrho_t)^2 \sum_{m=1}^{M} \mathbb{E}\|v_{t-1}^m - \bar{v}_{t-1}\|^2 + (1+\frac{1}{\nu}) \left(8L_f^2 C_g^2 \eta_{t-1}^2 \gamma^2 \sum_{m=1}^{M} \mathbb{E}\|d_{t-1}^m - \bar{d}_{t-1}\|^2\right)$  $+8L_{f}^{2}C_{g}^{2}\eta_{t-1}^{2}\gamma^{2}\sum_{t}^{M}\mathbb{E}\|\bar{d}_{t-1}\|^{2}+16(q-1)L_{f}^{2}C_{g}^{2}\alpha_{t}^{2}\sum_{t}^{t-1}\gamma^{2}\eta_{l}^{2}\sum_{t}^{M}\mathbb{E}\|d_{l}^{m}-\bar{d}_{l}\|^{2}+8ML_{f}^{2}\alpha_{t}^{2}\sigma^{2}$  $+4M\varrho_t^2\sigma^2+32(\varrho_t^2+\alpha_t^2)L_f^2\sum_{i=1}^M\mathbb{E}\|h_{t-1}^m-g^m(\bar{x}_{t-1})\|^2+16M\varrho_t^2\delta_f^2+16ML_f^2\varrho_t^2\delta_g^2\bigg)$  $\leq (1+\nu)(1-\varrho_t)^2 \sum_{j=1}^{M} \mathbb{E}\|v_{t-1}^m - \bar{v}_{t-1})\|^2 + (1+\frac{1}{\nu}) \left(\frac{72L_f^2 C_g^4 \eta_{t-1}^2 \gamma^2}{\rho^2} \sum_{j=1}^{M} \mathbb{E}\|v_{t-1}^m - \bar{v}_{t-1}\|^2\right)$  $+\frac{72C_f^2L_f^2C_g^2\eta_{t-1}^2\gamma^2}{\rho^2}\sum_{j=1}^M \mathbb{E}\|u_{t-1}^m - \bar{u}_{t-1}\|^2 + 8L_f^2C_g^2\eta_{t-1}^2\gamma^2\sum_{j=1}^M \mathbb{E}\|\bar{d}_{t-1}\|^2$  $+16(q-1)L_{f}^{2}C_{g}^{2}\alpha_{t}^{2}\sum_{l=1}^{t-1}\gamma^{2}\eta_{l}^{2}\left(\frac{9C_{g}^{2}}{\rho^{2}}\sum_{m=1}^{M}\mathbb{E}\|v_{l}^{m}-\bar{v}_{l}\|^{2}+\frac{9C_{f}^{2}}{\rho^{2}}\sum_{m=1}^{M}\mathbb{E}\|u_{l}^{m}-\bar{u}_{l}\|^{2}\right)$  $+32(\varrho_t^2+\alpha_t^2)L_f^2\sum_{t=1}^{M}\mathbb{E}\|h_{t-1}^m-g^m(\bar{x}_{t-1})\|^2+8ML_f^2\alpha_t^2\sigma^2+4M\varrho_t^2\sigma^2+16M\varrho_t^2\delta_f^2+16ML_f^2\varrho_t^2\delta_g^2\bigg),$ (51)

where the second last inequality holds by the above Lemma 9 and the above inequality (45), and the last inequality holds by  $d_{t-1}^m = \frac{x_t^m - x_{t-1}^m}{\eta_t \gamma}$ , and the above inequality (45).

where the second last inequality is due to Young inequality and the above Lemma 2.

$$\begin{aligned} & \text{Consider the term } \sum_{m=1}^{M} \left\| \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) - \frac{1}{M} \sum_{m=1}^{M} \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) \right\|^{2}, \text{ we have} \\ & \sum_{m=1}^{M} \left\| \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) - \frac{1}{M} \sum_{m=1}^{M} \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) \right\|^{2} \\ & = \sum_{m=1}^{M} \left\| \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) - \nabla g^{m}(x_{t-1}^{m}) - \frac{1}{M} \sum_{m=1}^{M} \left( \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) - \nabla g^{m}(x_{t-1}^{m}) \right) \\ & + \nabla g^{m}(x_{t-1}^{m}) - \frac{1}{M} \sum_{m=1}^{M} \nabla g^{m}(x_{t-1}^{m}) \right\|^{2} \\ & \leq 2 \sum_{m=1}^{M} \left\| \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) - \nabla g^{m}(x_{t-1}^{m}) - \frac{1}{M} \sum_{m=1}^{M} \left( \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) - \nabla g^{m}(x_{t-1}^{m}) \right) \right\| \\ & + 2 \sum_{m=1}^{M} \left\| \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) - \nabla g^{m}(x_{t-1}^{m}) \right\|^{2} \\ & \leq 2 \sum_{m=1}^{M} \left\| \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) - \nabla g^{m}(x_{t-1}^{m}) \right\|^{2} \\ & \leq 2 \sum_{m=1}^{M} \left\| \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) - \nabla g^{m}(x_{t-1}^{m}) \right\|^{2} \\ & \leq 2 \sum_{m=1}^{M} \left\| \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) - \nabla g^{m}(x_{t-1}^{m}) \right\|^{2} \\ & \leq 2 \sum_{m=1}^{M} \left\| \nabla g^{m}(x_{t-1}^{m};\zeta_{t}^{m}) - \nabla g^{m}(x_{t-1}^{m}) \right\|^{2} \\ & \leq 2 M \sigma^{2} + 12 L_{g}^{2} \sum_{m=1}^{M} \mathbb{E} \left\| x_{t-1}^{m} - \bar{x}_{t-1} \right\|^{2} + 6M \delta_{g}^{2}, \end{aligned}$$
(55)

where the last inequality holds by the above Lemma 8.

1320 By combining the above inequalities (52), (54) and (55), we have

where the last inequality holds by the above inequality (45).

1352 By summing the above inequalities (50) and (56), we have

 $\sum_{m=1}^{M} \left( \mathbb{E} \| u_{t}^{m} - \bar{u}_{t} \|^{2} + \mathbb{E} \| v_{t}^{m} - \bar{v}_{t} \|^{2} \right)$ 

$$\rho^{2} \qquad \lim_{m=1}^{m=1} (1 + 1)^{2} \sum_{m=1}^{m=1}^{m=1} (1 + 1)^{2} \sum_{m=1}^{m=1}^{m=1} (1 + 1)^{2} \sum_{l=s_{t}}^{l} \gamma^{2} \eta_{l}^{2} \left( \frac{9C_{g}^{2}}{\rho^{2}} \sum_{m=1}^{M} \mathbb{E} \|v_{l}^{m} - \bar{v}_{l}\|^{2} + \frac{9C_{f}^{2}}{\rho^{2}} \sum_{m=1}^{M} \mathbb{E} \|u_{l}^{m} - \bar{u}_{l}\|^{2} \right)$$

$$+ 32(\varrho_{t}^{2} + \alpha_{t}^{2})L_{f}^{2} \sum_{m=1}^{M} \mathbb{E} \|h_{t-1}^{m} - g^{m}(\bar{x}_{t-1})\|^{2} + 8ML_{f}^{2}\alpha_{t}^{2}\sigma^{2} + 4M\varrho_{t}^{2}\sigma^{2} + 16M\varrho_{t}^{2}\delta_{f}^{2} + 16ML_{f}^{2}\varrho_{t}^{2}\delta_{g}^{2} \right)$$

$$\leq \max\left((1 + \nu)(1 - \beta_{t})^{2} + (1 + \frac{1}{\nu})\frac{72C_{f}^{2}(L_{f}^{2}C_{g}^{2} + L_{g}^{2})\eta_{t-1}^{2}\gamma^{2}}{\rho^{2}}, (1 + \nu)(1 - \varrho_{t})^{2} + (1 + \frac{1}{\nu})\frac{72C_{g}^{2}(C_{g}^{2}L_{f}^{2} + L_{g}^{2})\eta_{t-1}^{2}\gamma^{2}}{\rho^{2}} \right)$$

$$\cdot \sum_{m=1}^{M} \left(\mathbb{E} \|u_{t-1}^{m} - \bar{u}_{t-1}\|^{2} + \mathbb{E} \|v_{t-1}^{m} - \bar{v}_{t-1}\|^{2}\right) + 8(1 + \frac{1}{\nu})(L_{f}^{2}C_{g}^{2} + L_{g}^{2})\eta_{t-1}^{2}\gamma^{2}} \sum_{m=1}^{M} \mathbb{E} \|\bar{u}_{t-1}\|^{2} + 24(1 + \frac{1}{\nu})(q - 1)(L_{f}^{2}C_{g}^{2}\alpha_{t}^{2} + L_{g}^{2}\beta_{t}^{2}) \sum_{l=s_{t}}^{t-1} \gamma^{2}\eta_{l}^{2} \left(\frac{9C_{g}^{2}}{\rho^{2}}\sum_{m=1}^{M} \mathbb{E} \|v_{l}^{m} - \bar{v}_{l}\|^{2} + \frac{9C_{f}^{2}}{\rho^{2}}\sum_{m=1}^{M} \mathbb{E} \|u_{l}^{m} - \bar{u}_{l}\|^{2}\right)$$

$$+ 32(1 + \frac{1}{\nu})(\varrho_{t}^{2} + \alpha_{t}^{2})L_{f}^{2} \sum_{m=1}^{M} \mathbb{E} \|h_{t-1}^{m} - g^{m}(\bar{x}_{t-1})\|^{2} + (1 + \frac{1}{\nu})\left(8ML_{f}^{2}\alpha_{t}^{2}\sigma^{2} + 4M\varrho_{t}^{2}\sigma^{2} + 16M\varrho_{t}^{2}\delta_{f}^{2} + 16M\varrho_{t}^{2}\delta_{f}^{2}\right)$$

$$+ 16ML_{f}^{2}\varrho_{t}^{2}\delta_{g}^{2} + 4M\sigma^{2}\beta_{t}^{2} + 12M\delta_{g}^{2}\beta_{t}^{2}\right).$$

$$(59)$$

(58)

(60)

1394 Let  $C_{fg}^2 = \max(C_f^2, C_g^2), L_{fg}^2 = L_f^2 C_g^2 + L_g^2, \nu = \frac{1}{q} \text{ and } \eta_t \leq \frac{\rho}{24\gamma q L_{fg} C_{fg}} \text{ for all } t \geq 0.$  Since 1395  $\beta_t \in (0, 1) \text{ for all } t \geq 0$ , we have

 $\leq (1+\nu)(1-\beta_t)^2 \sum_{i=1}^{M} \mathbb{E}\|u_{t-1}^m - \bar{u}_{t-1}\|^2 + (1+\frac{1}{\nu}) \left(\frac{36C_g^2 L_g^2 \eta_{t-1}^2 \gamma^2}{\rho^2} \sum_{i=1}^{M} \mathbb{E}\|v_{t-1}^m - \bar{v}_{t-1}\|^2\right)$ 

 $+24(q-1)L_{g}^{2}\beta_{t}^{2}\sum_{l=1}^{t-2}\gamma^{2}\eta_{l}^{2}\left(\frac{9C_{g}^{2}}{\rho^{2}}\sum_{l=1}^{M}\mathbb{E}\|v_{l}^{m}-\bar{v}_{l}\|^{2}+\frac{9C_{f}^{2}}{\rho^{2}}\sum_{l=1}^{M}\mathbb{E}\|u_{l}^{m}-\bar{u}_{l}\|^{2}\right)\right)$ 

 $+\frac{72C_f^2L_f^2C_g^2\eta_{t-1}^2\gamma^2}{2}\sum_{m=1}^M \mathbb{E}\|u_{t-1}^m - \bar{u}_{t-1}\|^2 + 8L_f^2C_a^2\eta_{t-1}^2\gamma^2\sum_{m=1}^M \mathbb{E}\|\bar{d}_{t-1}\|^2$ 

 $+\frac{36C_{f}^{2}L_{g}^{2}\eta_{t-1}^{2}\gamma^{2}}{\rho^{2}}\sum_{t}^{M}\mathbb{E}\|u_{t-1}^{m}-\bar{u}_{t-1}\|^{2}+4L_{g}^{2}\eta_{t-1}^{2}\gamma^{2}\sum_{t}^{M}\mathbb{E}\|\bar{d}_{t-1}\|^{2}+4M\sigma^{2}\beta_{t}^{2}+12M\delta_{g}^{2}\beta_{t}^{2}$ 

 $+ (1+\nu)(1-\varrho_t)^2 \sum_{j=1}^{M} \mathbb{E}\|v_{t-1}^m - \bar{v}_{t-1})\|^2 + (1+\frac{1}{\nu}) \left(\frac{72L_f^2 C_g^4 \eta_{t-1}^2 \gamma^2}{\rho^2} \sum_{j=1}^{M} \mathbb{E}\|v_{t-1}^m - \bar{v}_{t-1}\|^2\right)$ 

- 1399  $(1+\nu)(1-\beta_t)^2 + (1+\frac{1}{\nu})\frac{72C_f^2(L_f^2C_g^2 + L_g^2)\eta_{t-1}^2\gamma^2}{\rho^2}$

$$\begin{aligned} & \begin{array}{l} \mathbf{1400} \\ & \mathbf{1401} \\ & \mathbf{1401} \\ & \begin{array}{l} \mathbf{1402} \\ & \mathbf{1403} \\ \end{array} & \\ & \begin{array}{l} \leq 1 + \frac{1}{q} + (1+q) \frac{72C_f^2(L_f^2C_g^2 + L_g^2)\gamma^2}{\rho^2} \frac{\rho^2}{576\gamma^2 q^2 L_{fg}^2 C_{fg}^2} \\ & \\ & \begin{array}{l} \leq 1 + \frac{1}{q} + \frac{1+q}{8q^2} \leq 1 + \frac{5}{4q}. \\ \end{array} \end{aligned}$$

Similarly, since  $\varrho_t \in (0,1)$  for all  $t \ge 0$ , we have  $(1+\nu)(1-\varrho_t)^2 + (1+\frac{1}{\nu})\frac{72C_g^2(C_g^2L_f^2 + L_g^2)\eta_{t-1}^2\gamma^2}{\sigma^2} \le 1$  $1 + \frac{5}{4a}$ . Based on the above inequality (58) and the parameters, then we have  $\sum_{t=1}^{M} \left( \mathbb{E} \left\| u_t^m - \bar{u}_t \right\|^2 + \mathbb{E} \left\| v_t^m - \bar{v}_t \right\|^2 \right)$ (61) $\leq \left(1 + \frac{5}{4q}\right) \sum_{i=1}^{M} \left(\mathbb{E} \left\| u_{t-1}^{m} - \bar{u}_{t-1} \right\|^{2} + \mathbb{E} \left\| v_{t-1}^{m} - \bar{v}_{t-1} \right\|^{2}\right) + 8(q+1)L_{fg}^{2}\eta_{t-1}^{2}\gamma^{2} \sum_{i=1}^{M} \mathbb{E} \|\bar{d}_{t-1}\|^{2}$  $+216(q^{2}-1)\frac{C_{fg}^{2}L_{fg}^{2}\gamma^{2}}{\rho^{2}}\left(\alpha_{t}^{2}+\beta_{t}^{2}\right)\sum_{l=1}^{t-1}\eta_{l}^{2}\sum_{l=1}^{M}\left(\mathbb{E}\|v_{l}^{m}-\bar{v}_{l}\|^{2}+\mathbb{E}\|u_{l}^{m}-\bar{u}_{l}\|^{2}\right)$  $+ 32(1+q)(\varrho_t^2 + \alpha_t^2)L_f^2 \sum_{i=1}^M \mathbb{E} \|h_{t-1}^m - g^m(\bar{x}_{t-1})\|^2$  $+4M(q+1)\Big(2L_{f}^{2}\alpha_{t}^{2}\sigma^{2}+\varrho_{t}^{2}\sigma^{2}+4\varrho_{t}^{2}\delta_{f}^{2}+4L_{f}^{2}\varrho_{t}^{2}\delta_{q}^{2}+\sigma^{2}\beta_{t}^{2}+3\delta_{q}^{2}\beta_{t}^{2}\Big)$  $\leq \left(1 + \frac{5}{4q}\right) \sum_{i=1}^{M} \left(\mathbb{E} \left\| u_{t-1}^{m} - \bar{u}_{t-1} \right\|^{2} + \mathbb{E} \left\| v_{t-1}^{m} - \bar{v}_{t-1} \right\|^{2}\right) + \frac{\rho^{2}}{36qC_{t,i}^{2}} \sum_{i=1}^{M} \mathbb{E} \left\| \bar{d}_{t-1} \right\|^{2}$  $+\frac{3(c_1^2+c_2^2)}{8}\eta_{t-1}^2\sum_{l=1}^{t-2}\eta_l^2\sum_{l=1}^{M}\left(\mathbb{E}\|v_l^m-\bar{v}_l\|^2+\mathbb{E}\|u_l^m-\bar{u}_l\|^2\right)+\frac{\rho^2(c_1^2+c_3^2)}{9q\gamma^2C_{t-2}^2C_{t-2}^2}\eta_{t-1}^2\sum_{l=1}^{M}\mathbb{E}\|h_{t-1}^m-g^m(\bar{x}_{t-1})\|^2$  $+\frac{M\rho}{3\gamma L_{fg}C_{fg}}\Big(2c_1^2L_f^2\sigma^2+c_3^2\sigma^2+4c_3^2\delta_f^2+4c_3^2L_f^2\delta_g^2+c_2^2\sigma^2+3c_2^2\delta_g^2\Big)\eta_{t-1}^3,$ (62)where the first inequality holds by the above inequality (52) and  $\nu = \frac{1}{a}$ , and the last inequality holds by  $\alpha_t = c_1 \eta_{t-1}^2$ ,  $\beta_t = c_2 \eta_{t-1}^2$ ,  $\varrho_t = c_3 \eta_{t-1}^2$  and  $\eta_t \leq \frac{\rho}{24\gamma q L_{fg} C_{fg}}$  for all  $t \geq 0$ , and  $\frac{L_f^2}{L_{fg}^2} \leq \frac{1}{C_g^2}$ .

According to the above inequality (61), we have  $\sum_{m=1}^{M} \left( \mathbb{E} \left\| u_t^m - \bar{u}_t \right\|^2 + \mathbb{E} \left\| v_t^m - \bar{v}_t \right\|^2 \right)$  $\leq \frac{\rho^2}{36qC_t^2} \sum_{s=1}^{t-1} \left(1 + \frac{5}{4q}\right)^{t-1-s} \sum_{s=1}^M \mathbb{E}\|\bar{d}_s\|^2$  $+\frac{3(c_1^2+c_2^2)}{8}\sum_{l=1}^{t-1}\left(1+\frac{5}{4q}\right)^{t-1-s}\eta_s^2\sum_{l=1}^{s-2}\eta_l^2\sum_{l=1}^{M}\left(\mathbb{E}\|v_l^m-\bar{v}_l\|^2+\mathbb{E}\|u_l^m-\bar{u}_l\|^2\right)$  $+\frac{\rho^2(c_1^2+c_3^2)}{9a\gamma^2 C_{t-}^2 C_a^2} \sum_{c_{s-1}}^{t-1} \left(1+\frac{5}{4a}\right)^{t-1-s} \eta_s^2 \sum_{c_{s-1}}^M \mathbb{E}\|h_{s-1}^m - g^m(\bar{x}_{s-1})\|^2$  $+\frac{M\rho}{3\gamma L_{fg}C_{fg}} \left(2c_1^2 L_f^2 \sigma^2 + c_3^2 \sigma^2 + 4c_3^2 \delta_f^2 + 4c_3^2 L_f^2 \delta_g^2 + c_2^2 \sigma^2 + 3c_2^2 \delta_g^2\right) \sum_{s=1}^{t-1} \left(1 + \frac{5}{4q}\right)^{t-1-s} \eta_s^3$  $\leq \frac{\rho^2}{9qC_{4s}^2} \sum_{s=1}^{t-1} \sum_{s=1}^M \mathbb{E}\|\bar{d}_s\|^2 + \frac{3(c_1^2 + c_2^2)}{2} \sum_{s=1}^{t-1} \eta_s^2 \sum_{s=1}^{s-2} \eta_l^2 \sum_{s=1}^M \left(\mathbb{E}\|v_l^m - \bar{v}_l\|^2 + \mathbb{E}\|u_l^m - \bar{u}_l\|^2\right)$  $+\frac{4\rho^2(c_1^2+c_3^2)}{9q\gamma^2C_{\ell_s}^2C_{\epsilon_s}^2}\sum_{i=1}^{t-1}\eta_s^2\sum_{i=1}^M\mathbb{E}\|h_{s-1}^m-g^m(\bar{x}_{s-1})\|^2$  $+\frac{4M\rho}{3\gamma L_{fg}C_{fg}} \left(2c_1^2 L_f^2 \sigma^2 + c_3^2 \sigma^2 + 4c_3^2 \delta_f^2 + 4c_3^2 L_f^2 \delta_g^2 + c_2^2 \sigma^2 + 3c_2^2 \delta_g^2\right) \sum_{s=1}^{t-1} \eta_s^3$  $\leq \frac{M\rho^2}{9qC_{f_a}^2} \sum_{s=1}^{t-1} \mathbb{E}\|\bar{d}_s\|^2 + \frac{\rho^2(c_1^2 + c_2^2)}{24 * 16\gamma^2 qL_{f_a}^2 C_{f_a}^2} \sum_{s=1}^{t-1} \eta_s^2 \sum_{s=1}^{M} \left(\mathbb{E}\|v_s^m - \bar{v}_s\|^2 + \mathbb{E}\|u_s^m - \bar{u}_s\|^2\right)$  $+\frac{4\rho^2(c_1^2+c_3^2)}{9q\gamma^2C_{t_a}^2C_a^2}\sum_{s=1}^{t-1}\eta_s^2\sum_{s=1}^M \mathbb{E}\|h_{s-1}^m-g^m(\bar{x}_{s-1})\|^2$  $+\frac{4M\rho}{3\gamma L_{fg}C_{fg}} \left(2c_1^2 L_f^2 \sigma^2 + c_3^2 \sigma^2 + 4c_3^2 \delta_f^2 + 4c_3^2 L_f^2 \delta_g^2 + c_2^2 \sigma^2 + 3c_2^2 \delta_g^2\right) \sum_{i=1}^{t-1} \eta_s^3,$ (63)

where the second inequality holds by  $\left(1+\frac{5}{4a}\right)^{t-1-s} \leq \left(1+\frac{5}{4a}\right)^q \leq e^{5/4} \leq 4$  and the last inequality holds by  $\eta_t \leq \frac{\rho}{24\gamma q L_{f_a} C_{f_a}}$  for all  $t \geq 0$ .

By multiplying both sides of (63) by  $\eta_t$  and summing over  $t = s_t$  to  $s_t + q - 1$ , we have

$$\begin{split} & \overset{s_{t}+q-1}{1501} \qquad \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \sum_{m=1}^{M} \left( \mathbb{E} \left\| u_{t}^{m} - \bar{u}_{t} \right\|^{2} + \mathbb{E} \left\| v_{t}^{m} - \bar{v}_{t} \right\|^{2} \right) \\ & \overset{1503}{1503} \qquad \qquad \leq \frac{M\rho^{2}}{9C_{fg}^{2}} \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \mathbb{E} \left\| \bar{d}_{t} \right\|^{2} + \frac{\rho^{4}(c_{1}^{2}+c_{2}^{2})}{24^{3}*16\gamma^{4}q^{2}L_{fg}^{4}C_{fg}^{4}} \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \sum_{m=1}^{M} \left( \mathbb{E} \| v_{t}^{m} - \bar{v}_{t} \|^{2} + \mathbb{E} \| u_{t}^{m} - \bar{u}_{t} \|^{2} \right) \\ & \overset{1506}{1506} \qquad + \frac{\rho^{4}(c_{1}^{2}+c_{3}^{2})}{24*54q^{2}\gamma^{4}C_{fg}^{4}L_{fg}^{2}C_{g}^{2}} \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \sum_{m=1}^{M} \mathbb{E} \| h_{t}^{m} - g^{m}(\bar{x}_{t}) \|^{2} \\ & \overset{1509}{1510} \qquad + \frac{M\rho^{2}}{18\gamma^{2}L_{fg}^{2}C_{fg}^{2}} \left( 2c_{1}^{2}L_{fg}^{2}\sigma^{2} + c_{3}^{2}\sigma^{2} + 4c_{3}^{2}\delta_{f}^{2} + 4c_{3}^{2}L_{fg}^{2}\delta_{g}^{2} + c_{2}^{2}\sigma^{2} + 3c_{2}^{2}\delta_{g}^{2} \right) \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t}^{3}, \quad (64) \end{aligned}$$

Given 
$$c_1^2 + c_2^2 \leq \frac{(24)^4 q^2 \gamma^4 L_{fg}^4 C_{fg}^4}{9\rho^4}$$
, we have  $\frac{60}{72} \leq 1 - \frac{\rho^4 (c_1^2 + c_2^2)}{24^3 * 16\gamma^4 q^2 L_{fg}^4 C_{fg}^4}$ , we have  
 $\sum_{t=s_t}^{s_t+q-1} \eta_t \sum_{m=1}^M \left( \mathbb{E} \left\| u_t^m - \bar{u}_t \right\|^2 + \mathbb{E} \left\| v_t^m - \bar{v}_t \right\|^2 \right)$   
 $\leq \frac{2M\rho^2}{15C_{fg}^2} \sum_{t=s_t}^{s_t+q-1} \eta_t \mathbb{E} \|\bar{d}_t\|^2 + \frac{\rho^4 (c_1^2 + c_3^2)}{1080q^2 \gamma^4 C_{fg}^4 L_{fg}^2 C_g^2} \sum_{t=s_t}^{s_t+q-1} \eta_t \sum_{m=1}^M \mathbb{E} \| h_t^m - g^m(\bar{x}_t) \|^2$   
 $+ \frac{M\rho^2}{15\gamma^2 L_{fg}^2 C_{fg}^2} \left( 2c_1^2 L_f^2 \sigma^2 + c_3^2 \sigma^2 + 4c_3^2 \delta_f^2 + 4c_3^2 L_f^2 \delta_g^2 + c_2^2 \sigma^2 + 3c_2^2 \delta_g^2 \right) \sum_{t=s_t}^{s_t+q-1} \eta_t^3.$  (65)  
1523

According to the above inequality (45) and  $C_{fg}^2 = \max(C_f^2, C_g^2)$ , we have

$$\sum_{m=1}^{1526} \sum_{m=1}^{M} \mathbb{E} \|d_t^m - \bar{d}_t\|^2 \le \frac{9C_g^2}{\rho^2} \sum_{m=1}^{M} \mathbb{E} \|v_t^m - \bar{v}_t\|^2 + \frac{9C_f^2}{\rho^2} \sum_{m=1}^{M} \mathbb{E} \|u_t^m - \bar{u}_t\|^2 \le \frac{9C_{fg}^2}{\rho^2} \sum_{m=1}^{M} \left( \mathbb{E} \|u_t^m - \bar{u}_t\|^2 + \mathbb{E} \|v_t^m - \bar{v}_t\|^2 \right)$$

$$(66)$$

Thus we have

$$\sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \sum_{m=1}^{M} \mathbb{E} \|d_{t}^{m} - \bar{d}_{t}\|^{2}$$

$$\leq \frac{9C_{fg}^{2}}{\rho^{2}} \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \sum_{m=1}^{M} \left( \mathbb{E} \|u_{t}^{m} - \bar{u}_{t}\|^{2} + \mathbb{E} \|v_{t}^{m} - \bar{v}_{t}\|^{2} \right)$$

$$\leq \frac{6M}{5} \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \mathbb{E} \|\bar{d}_{t}\|^{2} + \frac{\rho^{2}(c_{1}^{2} + c_{3}^{2})}{120q^{2}\gamma^{4}C_{fg}^{2}L_{fg}^{2}C_{g}^{2}} \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \sum_{m=1}^{M} \mathbb{E} \|h_{t}^{m} - g^{m}(\bar{x}_{t})\|^{2}$$

$$+ \frac{3M}{5\gamma^{2}L_{fg}^{2}} \left( 2c_{1}^{2}L_{f}^{2}\sigma^{2} + c_{3}^{2}\sigma^{2} + 4c_{3}^{2}\delta_{f}^{2} + 4c_{3}^{2}L_{f}^{2}\delta_{g}^{2} + c_{2}^{2}\sigma^{2} + 3c_{2}^{2}\delta_{g}^{2} \right) \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t}^{3}. \quad (67)$$

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\|\nabla F(\bar{x}_t)\| \le \left(\frac{\sqrt{2G}n^{1/6}}{T^{1/2}} + \frac{\sqrt{2G}}{T^{1/3}}\right)\sqrt{\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\|A_t\|^2},\tag{68}$$

where  $C_{fg}^2 = \max(C_f^2, C_g^2)$ ,  $L_{fg}^2 = L_f^2 C_g^2 + L_g^2$ ,  $G = \frac{4(F(\bar{x}_1) - F^*)}{k\rho\gamma} + \frac{12n^{1/3}\sigma^2}{qk^2\rho^2} + 4k^2 \Big(\frac{\hat{\delta}^2}{4\gamma^2 L_{fg}^2} + \frac{\hat{\delta}^2}{qk^2\rho^2} + \frac{12n^{1/3}\sigma^2}{qk^2\rho^2} + \frac{12n^{1/3}\sigma^2}{qk^2} +$  $\frac{\left(c_{1}^{2}+c_{2}^{2}+c_{3}^{2}\right)\sigma^{2}}{_{3\rho\gamma qL_{fg}C_{fg}}}\right)\ln(n+T) \text{ and } \hat{\delta}^{2} = 2c_{1}^{2}L_{f}^{2}\sigma^{2} + c_{3}^{2}\sigma^{2} + 4c_{3}^{2}\delta_{f}^{2} + 4c_{3}^{2}L_{f}^{2}\delta_{g}^{2} + c_{2}^{2}\sigma^{2} + 3c_{2}^{2}\delta_{g}^{2}.$ 

*Proof.* Since  $\eta_t = \frac{k}{(n+t)^{1/3}}$  on t is decreasing and  $n \ge k^3$ , we have  $\eta_t \le \eta_0 = \frac{k}{n^{1/3}} \le 1$  and  $\gamma \leq \frac{n^{1/3}\rho}{2Lk} \leq \frac{\rho}{2L\eta_0} \leq \frac{\rho}{2L\eta_t}$  for any  $t \geq 0$ . Since  $\eta_t \leq \frac{\rho}{24\gamma qL_{fg}C_{fg}}$  for all  $t \geq 0$ , we have 1566  $\frac{k}{n^{1/3}} = \eta_0 \leq \eta_t \leq \frac{\rho}{24\gamma q L_{fg} C_{fg}}$ , then we have  $n \geq \frac{(24k\gamma q L_{fg} C_{fg})^3}{\rho^3}$ . Due to  $0 < \eta_t \leq 1$  and 1567  $n \geq (c_1k)^3$ , we have  $\alpha_{t+1} = c_1\eta_t^2 \leq c_1\eta_t \leq \frac{c_1k}{n^{1/3}} \leq 1$ . Similarly, due to  $n \geq (c_2k)^3$  and 1568  $n \geq (c_3 k)^3$ , we have  $\beta_{t+1} \leq 1$  and  $\rho_{t+1} \leq 1$ . 1569 1570 According to Lemma 7, for any  $m \in [M]$ , we have 1571  $\frac{1}{n}\mathbb{E}\|h_{t+1}^m - g^m(x_{t+1}^m)\|^2 - \frac{1}{n}\mathbb{E}\|h_t^m - g^m(x_t^m)\|^2$ 1572 (69)1573  $\leq \left(\frac{1-\alpha_{t+1}}{\eta_t} - \frac{1}{\eta_{t-1}}\right) \mathbb{E} \|h_t^m - g^m(x_t^m)\|^2 + 2C_g^2 \mathbb{E} \|x_{t+1}^m - x_t^m\|^2 + 2\alpha_{t+1}^2 \sigma^2$ 1574 1575  $= \left(\frac{1}{n_t} - \frac{1}{n_{t-1}} - c_1 \eta_t\right) \mathbb{E} \|h_t^m - g^m(x_t^m)\|^2 + 2C_g^2 \mathbb{E} \|x_{t+1}^m - x_t^m\|^2 + 2\alpha_{t+1}^2 \sigma^2,$ 1576 1577 where the second equality is due to  $\alpha_{t+1} = c_1 \eta_t^2$ . Similarly, since  $\beta_{t+1} = c_2 \eta_t^2$ , we have 1579  $\frac{1}{\eta_t} \mathbb{E} \| u_{t+1}^m - \nabla g^m(x_{t+1}^m) \|^2 - \frac{1}{\eta_{t-1}} \mathbb{E} \| u_t^m - \nabla g^m(x_t^m) \|^2$ 1580 (70)1581  $\leq \Big(\frac{1-\beta_{t+1}}{m} - \frac{1}{m-1}\Big)\mathbb{E}\|u_t^m - \nabla g^m(x_t^m)\|^2 + 2L_g^2\mathbb{E}\|x_{t+1}^m - x_t^m\|^2 + 2\beta_{t+1}^2\sigma^2$ 1582 1583 1584  $= \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} - c_2\eta_t\right) \mathbb{E} \|u_t^m - \nabla g^m(x_t^m)\|^2 + 2L_g^2 \mathbb{E} \|x_{t+1}^m - x_t^m\|^2 + 2\beta_{t+1}^2 \sigma^2.$ 1585 1586 And we have 1587  $\frac{1}{n} \mathbb{E} \|v_{t+1}^m - \nabla f^m(h_{t+1}^m)\|^2 - \frac{1}{n} \mathbb{E} \|v_t^m - \nabla f^m(h_t^m)\|^2$ (71)1589  $\leq \big(\frac{1-\varrho_{t+1}}{2} - \frac{1}{m-1}\big)\mathbb{E}\|v_t^m - \nabla f^m(h_t^m)\|^2 + 4L_f^2 C_g^2 \mathbb{E}\|x_{t+1}^m - x_t^m\|^2 + 2\varrho_{t+1}^2 \sigma^2$ 1590 1591  $+8\alpha_{t+1}^{2}L_{t}^{2}\mathbb{E}\|h_{t}^{m}-q^{m}(x_{t}^{m})\|^{2}+8L_{t}^{2}\alpha_{t+1}^{2}\sigma^{2}$ 1592 1593  $= \left(\frac{1}{m} - \frac{1}{m-1} - c_3 \eta_t\right) \mathbb{E} \|v_t^m - \nabla f^m(h_t^m)\|^2 + 4L_f^2 C_g^2 \mathbb{E} \|x_{t+1}^m - x_t^m\|^2 + 2\varrho_{t+1}^2 \sigma^2$ 1594  $+8\alpha_{t+1}^2 L_t^2 \mathbb{E} \|h_t^m - q^m(x_t^m)\|^2 + 8L_t^2 \alpha_{t+1}^2 \sigma^2,$ 1596 1597 where the second equality is due to  $\rho_{t+1} = c_2 \eta_t^2$ . 1598 By  $\eta_t = \frac{k}{(n+t)^{1/3}}$ , we have

$$\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} = \frac{1}{k} \left( (n+t)^{\frac{1}{3}} - (n+t-1)^{\frac{1}{3}} \right) \le \frac{1}{3k(n+t-1)^{2/3}} \le \frac{1}{3k\left(n/2+t\right)^{2/3}} \\ \le \frac{2^{2/3}}{3k(n+t)^{2/3}} = \frac{2^{2/3}}{3k^3} \frac{k^2}{(n+t)^{2/3}} = \frac{2^{2/3}}{3k^3} \eta_t^2 \le \frac{2}{3k^3} \eta_t,$$
(72)

where the first inequality holds by the concavity of function  $f(x) = x^{1/3}$ , *i.e.*,  $(x + y)^{1/3} \le x^{1/3} + \frac{y}{3x^{2/3}}$ ; the second inequality is due to  $n \ge 2$ , and the last inequality is due to  $0 < \eta_t \le 1$ .

Let 
$$c_1 \ge \frac{2}{3k^3} + B$$
, for any  $m \in [M]$ , we have

1604

$$\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\end{array}}{1} & \frac{1}{\eta_t} \mathbb{E} \|h_{t+1}^m - g^m(x_{t+1}^m)\|^2 - \frac{1}{\eta_{t-1}} \mathbb{E} \|h_t^m - g^m(x_t^m)\|^2 & (73) \\ \end{array} \\
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\begin{array}{ll}
\end{array}}{1612} & \leq -B\eta_t \mathbb{E} \|h_t^m - g^m(x_t^m)\|^2 + 2C_g^2 \mathbb{E} \|x_{t+1}^m - x_t^m\|^2 + 2\alpha_{t+1}^2 \sigma^2 \\ \end{array} \\
\begin{array}{ll}
\end{array}}{1613} & = -B\eta_t \mathbb{E} \|h_t^m - g^m(x_t^m)\|^2 + 2C_g^2 \eta_t^2 \gamma^2 \mathbb{E} \|d_t^m - \bar{d}_t + \bar{d}_t\|^2 + 2\alpha_{t+1}^2 \sigma^2 \\ \end{array} \\
\begin{array}{ll}
\end{array}}{1614} & = -B\eta_t \mathbb{E} \|h_t^m - g^m(x_t^m)\|^2 + 4C_g^2 \eta_t^2 \gamma^2 \mathbb{E} \|d_t^m - \bar{d}_t\|^2 + 4C_g^2 \eta_t^2 \gamma^2 \mathbb{E} \|\bar{d}_t\|^2 + 2\alpha_{t+1}^2 \sigma^2 \\ \end{array} \\
\end{array} \\
\begin{array}{ll}
\end{array}}{1615} & \leq -B\eta_t \mathbb{E} \|h_t^m - g^m(\bar{x}_t)\|^2 + 4C_g^2 \eta_t^2 \gamma^2 \mathbb{E} \|d_t^m - \bar{d}_t\|^2 + 4C_g^2 \eta_t^2 \gamma^2 \mathbb{E} \|\bar{d}_t\|^2 + 2\alpha_{t+1}^2 \sigma^2 \\ \end{array} \\
\end{array} \\
\begin{array}{ll}
\end{array}}{1616} \\
\end{array} \\
\begin{array}{ll}
\end{array}}{1617} & \leq -\frac{B}{2} \eta_t \|h_t^m - g^m(\bar{x}_t)\|^2 + BC_g^2 \eta_t \|x_t^m - \bar{x}_t\|^2 + 4C_g^2 \eta_t^2 \gamma^2 \mathbb{E} \|d_t^m - \bar{d}_t\|^2 + 4C_g^2 \eta_t^2 \gamma^2 \mathbb{E} \|\bar{d}_t\|^2 + 2\alpha_{t+1}^2 \\ \end{array} \\
\end{array} \\
\end{array} \\
\begin{array}{ll}
\end{array}}{1618} \\
\end{array} \\$$
where the last inequality holds by  $-\|h_t^m - g^m(x_t^m)\|^2 \leq -\frac{1}{2}\|h_t^m - g^m(\bar{x}_t)\|^2 + \|g^m(x_t^m) - g^m(\bar{x}_t)\|^2 + C_g^2\|x_t^m - \bar{x}_t\|^2. \end{array}$ 

 $\sigma^2$ ,

Let  $c_2 \geq \frac{2}{3L^3} + 5C_f^2$ , for any  $m \in [M]$ , we have  $\frac{1}{n} \mathbb{E} \|u_{t+1}^m - \nabla g^m(x_{t+1}^m)\|^2 - \frac{1}{n} \mathbb{E} \|u_t^m - \nabla g^m(x_t^m)\|^2$ (74) $< -5C_{f}^{2}\eta_{t}\mathbb{E}\|u_{t}^{m} - \nabla q^{m}(x_{t}^{m})\|^{2} + 2L_{a}^{2}\mathbb{E}\|x_{t+1}^{m} - x_{t}^{m}\|^{2} + 2\beta_{t+1}^{2}\sigma^{2}$  $= -5C_t^2 \eta_t \mathbb{E} \|u_t^m - \nabla q^m (x_t^m)\|^2 + 2L_s^2 \eta_t^2 \gamma^2 \mathbb{E} \|d_t^m - \bar{d}_t + \bar{d}_t\|^2 + 2\beta_{t+1}^2 \sigma^2$  $\leq -5C_{f}^{2}\eta_{t}\mathbb{E}\|u_{t}^{m}-\nabla g^{m}(x_{t}^{m})\|^{2}+4L_{a}^{2}\eta_{t}^{2}\gamma^{2}\mathbb{E}\|d_{t}^{m}-\bar{d}_{t}\|^{2}+4L_{a}^{2}\eta_{t}^{2}\gamma^{2}\mathbb{E}\|\bar{d}_{t}\|^{2}+2\beta_{t}^{2}\beta_{t}^{2}\sigma^{2}$  $\leq -\frac{5C_{f}^{2}\eta_{t}}{2}\mathbb{E}\|u_{t}^{m}-\nabla g^{m}(\bar{x}_{t})\|^{2}+5C_{f}^{2}L_{q}^{2}\eta_{t}\|x_{t}^{m}-\bar{x}_{t}\|^{2}+4L_{q}^{2}\eta_{t}^{2}\gamma^{2}\mathbb{E}\|d_{t}^{m}-\bar{d}_{t}\|^{2}+4L_{q}^{2}\eta_{t}^{2}\gamma^{2}\mathbb{E}\|\bar{d}_{t}\|^{2}+2\beta_{t+1}^{2}\sigma^{2},$ where the last inequality holds by  $-\|u_t^m - \nabla g^m(x_t^m)\|^2 \le -\frac{1}{2}\|u_t^m - \nabla g^m(\bar{x}_t)\|^2 + \|\nabla g(x_t^m) - \nabla g^m(\bar{x}_t)\|^2$  $\nabla g^m(\bar{x}_t)\|^2 \le -\frac{1}{2} \|u_t^m - \nabla g^m(\bar{x}_t)\|^2 + L_a^2 \|x_t^m - \bar{x}_t\|^2.$ Let  $c_3 \geq \frac{2}{3k^3} + 5C_a^2$ , for any  $m \in [M]$ , we have  $\frac{1}{\eta_t} \mathbb{E} \| v_{t+1}^m - \nabla f^m(h_{t+1}^m) \|^2 - \frac{1}{\eta_{t-1}} \mathbb{E} \| v_t^m - \nabla f^m(h_t^m) \|^2$ (75) $\leq -5C_g^2\eta_t\mathbb{E}\|v_t^m - \nabla f^m(h_t^m)\|^2 + 4L_f^2C_g^2\mathbb{E}\|x_{t+1}^m - x_t^m\|^2 + 2\varrho_{t+1}^2\sigma^2 + 8\alpha_{t+1}^2L_f^2\mathbb{E}\|h_t^m - g^m(x_t^m)\|^2 + 8L_f^2\alpha_{t+1}^2\sigma^2 + 2\rho_{t+1}^2\sigma^2 + 2\rho_{t+1}^$  $= -5C_a^2\eta_t \mathbb{E} \|v_t^m - \nabla f^m(h_t^m)\|^2 + 4L_f^2 C_a^2 \eta_t^2 \gamma^2 \mathbb{E} \|d_t^m - \bar{d}_t + \bar{d}_t\|^2 + 2\varrho_{t+1}^2 \sigma^2 \eta_t^2 \nabla f^m(h_t^m)\|^2 + 4L_f^2 C_a^2 \eta_t^2 \gamma^2 \mathbb{E} \|d_t^m - \bar{d}_t + \bar{d}_t\|^2 + 2\varrho_{t+1}^2 \sigma^2 \eta_t^2 \nabla f^m(h_t^m)\|^2 + 4L_f^2 C_a^2 \eta_t^2 \gamma^2 \mathbb{E} \|d_t^m - \bar{d}_t + \bar{d}_t\|^2 + 2\varrho_{t+1}^2 \sigma^2 \eta_t^2 \nabla f^m(h_t^m)\|^2 + 4L_f^2 C_a^2 \eta_t^2 \gamma^2 \mathbb{E} \|d_t^m - \bar{d}_t + \bar{d}_t\|^2 + 2\varrho_{t+1}^2 \sigma^2 \eta_t^2 \nabla f^m(h_t^m)\|^2 + 4L_f^2 C_a^2 \eta_t^2 \gamma^2 \mathbb{E} \|d_t^m - \bar{d}_t + \bar{d}_t\|^2 + 2\varrho_{t+1}^2 \sigma^2 \eta_t^2 \nabla f^m(h_t^m)\|^2 + 4L_f^2 \nabla f^m(h$  $+8\alpha_{t+1}^2 L_f^2 \mathbb{E} \|h_t^m - q^m (x_t^m)\|^2 + 8L_f^2 \alpha_{t+1}^2 \sigma^2$  $\leq -5C_{g}^{2}\eta_{t}\mathbb{E}\|v_{t}^{m}-\nabla f^{m}(h_{t}^{m})\|^{2}+8L_{f}^{2}C_{q}^{2}\eta_{t}^{2}\gamma^{2}\mathbb{E}\|d_{t}^{m}-\bar{d}_{t}\|^{2}+8L_{f}^{2}C_{q}^{2}\eta_{t}^{2}\gamma^{2}\mathbb{E}\|\bar{d}_{t}\|^{2}+2\varrho_{t+1}^{2}\sigma^{2}$  $+8\alpha_{t+1}^2 L_t^2 \mathbb{E} \|h_t^m - q^m(x_t^m)\|^2 + 8L_t^2 \alpha_{t+1}^2 \sigma^2$  $\leq -5C_g^2\eta_t \mathbb{E}\|v_t^m - \nabla f^m(h_t^m)\|^2 + 8L_f^2 C_g^2\eta_t^2 \gamma^2 \mathbb{E}\|d_t^m - \bar{d}_t\|^2 + 8L_f^2 C_g^2\eta_t^2 \gamma^2 \mathbb{E}\|\bar{d}_t\|^2 + 2\varrho_{t+1}^2 \sigma^2 ||v_t^m||^2 + 2\rho_{t+1}^2 \sigma^2 ||v_t^m||^2 +$  $+\frac{c_2^2 L_f^2}{864q^3 \gamma^3 L_{fa}^3 C_{fa}^3} \eta_t \mathbb{E} \|h_t^m - g^m(\bar{x}_t)\|^2 + \frac{c_2^2 C_g^2 L_f^2}{864q^3 \gamma^3 L_{fa}^3 C_{fa}^3} \eta_t \mathbb{E} \|x_t^m - \bar{x}_t\|^2 + 8L_f^2 \alpha_{t+1}^2 \sigma^2,$ where the last inequality holds by Assumption ,  $\alpha_{t+1} = c_2 \eta_t^2$  and  $\eta_t \leq \frac{\rho}{24q\gamma L_{fg}C_{fg}}$  for all  $t \geq 0$ . According to Lemma 4, we have  $F(\bar{x}_{t+1}) - F(\bar{x}_t) \le \frac{1}{M} \sum_{j=1}^{M} \left( \frac{2C_f^2 \eta_t \gamma}{\rho} \| u_t^m - \nabla g^m(\bar{x}_t) \|^2 + \frac{4C_g^2 \eta_t \gamma}{\rho} \| v_t^m - \nabla f^m(h_t^m) \|^2 \right)$  $+\frac{4C_g^2 L_f^2 \eta_t \gamma}{\rho} \|h_t^m - g^m(\bar{x}_t)\|^2 \Big) - \frac{\rho}{2n_t \gamma} \|\bar{x}_{t+1} - \bar{x}_t\|^2$  $= \frac{1}{M} \sum_{a}^{M} \left( \frac{2C_{f}^{2} \eta_{t} \gamma}{\rho} \| u_{t}^{m} - \nabla g^{m}(\bar{x}_{t}) \|^{2} + \frac{4C_{g}^{2} \eta_{t} \gamma}{\rho} \| v_{t}^{m} - \nabla f^{m}(h_{t}^{m}) \|^{2} \right)$  $+\frac{4C_g^2 L_f^2 \eta_t \gamma}{2} \|h_t^m - g^m(\bar{x}_t)\|^2 \Big) - \frac{\rho \eta_t \gamma}{2} \|\bar{d}_t\|^2.$ (76) $\sum_{t=1}^{s_t+q-1} \eta_t \sum_{t=1}^M \mathbb{E} \|d_t^m - \bar{d}_t\|^2$  $\leq \frac{6M}{5} \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \mathbb{E} \|\bar{d}_{t}\|^{2} + \frac{\rho^{2}(c_{1}^{2}+c_{3}^{2})}{120q^{2}\gamma^{4}C_{fa}^{2}L_{fa}^{2}C_{q}^{2}} \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t} \sum_{m=1}^{M} \mathbb{E} \|h_{t}^{m} - g^{m}(\bar{x}_{t})\|^{2} + \frac{3M\hat{\delta}^{2}}{5\gamma^{2}L_{fg}^{2}} \sum_{t=s_{t}}^{s_{t}+q-1} \eta_{t}^{3},$ Next, we define a *potential* function, for any  $t \ge 1$  $\Omega_t = \mathbb{E}\Big[F(\bar{x}_t) + \frac{\gamma}{\rho m_{t-1}} \frac{1}{M} \sum_{i=1}^{M} \Big( \|h_t^m - g^m(x_t^m)\|^2 + \|u_t^m - \nabla g^m(x_t^m)\|^2 + \|v_t^m - \nabla f^m(h_t^m)\|^2 \Big) \Big].$ 

<sup>674</sup> Tl	en we have
75	
76	
77	
78	
79	
30	
31	
32	
33	
34	
<sup>85</sup> O	a = 0
36	+1 $2t$
87 = 88	$F(\bar{x}_{t+1}) - F(\bar{x}_t) + \frac{\gamma}{M\rho} \sum_{m=1}^{M} \left( \frac{1}{\eta_t} \mathbb{E} \ h_{t+1}^m - g^m(x_{t+1}^m)\ ^2 - \frac{1}{\eta_{t-1}} \mathbb{E} \ h_t^m - g^m(x_t^m)\ ^2 + \frac{1}{\eta_t} \ u_{t+1}^m - \nabla g^m(x_{t+1}^m)\ ^2 \right)$
)0 )1	$-\frac{1}{\eta_{t-1}}\ u_t^m - \nabla g^m(x_t^m)\ ^2 + \frac{1}{\eta_t}\ v_{t+1}^m - \nabla f^m(h_{t+1}^m)\ ^2 - \frac{1}{\eta_{t-1}}\ v_t^m - \nabla f^m(h_t^m)\ ^2\right)$
2 3 ≤ 4	$\frac{1}{M}\sum_{m=1}^{M} \left(\frac{2C_{f}^{2}\eta_{t}\gamma}{\rho}\ u_{t}^{m} - \nabla g^{m}(\bar{x}_{t})\ ^{2} + \frac{4C_{g}^{2}\eta_{t}\gamma}{\rho}\ v_{t}^{m} - \nabla f^{m}(h_{t}^{m})\ ^{2} + \frac{4C_{g}^{2}L_{f}^{2}\eta_{t}\gamma}{\rho}\ h_{t}^{m} - g^{m}(\bar{x}_{t})\ ^{2}\right) - \frac{\rho\eta_{t}\gamma}{2}\ \bar{d}_{t}\ ^{2}$
5 6 7	$+\frac{\gamma}{M\rho}\sum_{m=1}^{M}\left(-\frac{B}{2}\eta_{t}\ h_{t}^{m}-g^{m}(\bar{x}_{t})\ ^{2}+BC_{g}^{2}\eta_{t}\ x_{t}^{m}-\bar{x}_{t}\ ^{2}+4C_{g}^{2}\eta_{t}^{2}\gamma^{2}\mathbb{E}\ d_{t}^{m}-\bar{d}_{t}\ ^{2}+4C_{g}^{2}\eta_{t}^{2}\gamma^{2}\mathbb{E}\ \bar{d}_{t}\ ^{2}+2\alpha_{t+1}^{2}\sigma^{2}$
8	$-\frac{5C_{f}^{2}\eta_{t}}{2}\mathbb{E}\ u_{t}^{m}-\nabla g^{m}(\bar{x}_{t})\ ^{2}+5C_{f}^{2}L_{g}^{2}\eta_{t}\ x_{t}^{m}-\bar{x}_{t}\ ^{2}+4L_{g}^{2}\eta_{t}^{2}\gamma^{2}\mathbb{E}\ d_{t}^{m}-\bar{d}_{t}\ ^{2}+4L_{g}^{2}\eta_{t}^{2}\gamma^{2}\mathbb{E}\ \bar{d}_{t}\ ^{2}+2\beta_{t+1}^{2}\sigma^{2}$
)0 )1 )2	$-5C_{g}^{2}\eta_{t}\mathbb{E}\ v_{t}^{m}-\nabla f^{m}(h_{t}^{m})\ ^{2}+8L_{f}^{2}C_{g}^{2}\eta_{t}^{2}\gamma^{2}\mathbb{E}\ d_{t}^{m}-\bar{d}_{t}\ ^{2}+8L_{f}^{2}C_{g}^{2}\eta_{t}^{2}\gamma^{2}\mathbb{E}\ \bar{d}_{t}\ ^{2}+2\varrho_{t+1}^{2}\sigma^{2}$ $c_{2}^{2}L_{t}^{2}$
)3 )4	$+\frac{2}{864q^{3}\gamma^{3}L_{fg}^{3}C_{fg}^{3}}\eta_{t}\mathbb{E}\ h_{t}^{m}-g^{m}(\bar{x}_{t})\ ^{2}+\frac{2}{864q^{3}\gamma^{3}L_{fg}^{3}C_{fg}^{3}}\eta_{t}\mathbb{E}\ x_{t}^{m}-\bar{x}_{t}\ ^{2}+8L_{f}^{2}\alpha_{t+1}^{2}\sigma^{2}\Big)$
)5 ≤ )6	$\frac{1}{M}\sum_{m=1}^{\infty} \left( -\frac{\nabla f^{\gamma}}{2\rho} \eta_t \  u_t^m - \nabla g^m(\bar{x}_t) \ ^2 - \frac{\nabla g^{\gamma}}{\rho} \eta_t \  v_t^m - \nabla f^m(h_t^m) \ ^2 \right)$
8 9	$-\frac{\gamma}{\rho} \left(\frac{B}{2} - 4C_g^2 L_f^2 - \frac{c_2^2 C_g^2 L_f^2}{864q^3 \gamma^3 L_{fg}^3 C_{fg}^3} \right) \eta_t \ h_t^m - g^m(\bar{x}_t)\ ^2 \right) - \left(\frac{\rho \gamma \eta_t}{2} - \frac{4C_g^2 \eta_t^2 \gamma^3}{\rho} - \frac{4L_g^2 \eta_t^2 \gamma^3}{\rho} - \frac{8L_f^2 C_g^2 \eta_t^2 \gamma^3}{\rho} \right) \ \bar{d}_t\ ^2$
0 1 2	$+\frac{\gamma}{M\rho}\Big(BC_g^2+5C_f^2L_g^2+\frac{c_2^2C_g^2L_f^2}{864q^3\gamma^3L_{fg}^3C_{fg}^3}\Big)\eta_t(q-1)\sum_{l=s_t}^{t-1}\gamma^2\eta_l^2\sum_{m=1}^M\mathbb{E}\ d_l^m-\bar{d}_l\ ^2$
3 4	$+\frac{\gamma}{M_{o}}\left(4C_{g}^{2}\eta_{t}^{2}\gamma^{2}+4L_{g}^{2}\eta_{t}^{2}\gamma^{2}+8L_{f}^{2}C_{g}^{2}\eta_{t}^{2}\gamma^{2}\right)\sum^{M}\mathbb{E}\ d_{t}^{m}-\bar{d}_{t}\ ^{2}+\frac{2\sigma^{2}\gamma}{c}\left(\alpha_{t+1}^{2}+\beta_{t+1}^{2}+\varrho_{t+1}^{2}\right),$
5	$\mu \rho \sim \rho \sim \rho $ (70)
6	(78)
7	
8	
9	
U	
<u> </u>	
2 1	
•	
6	
27 wl in	here the first inequality holds by the above inequalities (73), (74), (75) and (76), and the last equality is due to Lemma 9.

Let  $s_t = q \lfloor t/q \rfloor + 1$ , summing the above inequality (78) over  $t = s_t$  to  $s_t + q - 1$ , we have  $\sum_{t=1}^{s_t+q-1} \left(\Omega_{t+1} - \Omega_t\right)$  $\leq \sum_{t=1}^{s_t+q-1} \frac{1}{M} \sum_{t=1}^{M} \left( -\frac{C_f^2 \gamma}{2\rho} \eta_t \| u_t^m - \nabla g^m(\bar{x}_t) \|^2 - \frac{C_g^2 \gamma}{\rho} \eta_t \| v_t^m - \nabla f^m(h_t^m) \|^2 \right)$  $-\frac{\gamma}{\rho} \Big(\frac{B}{2} - 4C_g^2 L_f^2 - \frac{c_2^2 C_g^2 L_f^2}{864 q^3 \gamma^3 L_{s_-}^3 C_s^3} \Big) \eta_t \|h_t^m - g^m(\bar{x}_t)\|^2 \Big)$  $-\sum_{p=1}^{s_t+q-1} \Big(\frac{\rho \gamma \eta_t}{2} - \frac{4C_g^2 \eta_t^2 \gamma^3}{\rho} - \frac{4L_g^2 \eta_t^2 \gamma^3}{\rho} - \frac{8L_f^2 C_g^2 \eta_t^2 \gamma^3}{\rho} \Big) \|\bar{d}_t\|^2$  $+\frac{\gamma}{M\rho}\Big(BC_g^2+5C_f^2L_g^2+\frac{c_2^2C_g^2L_f^2}{864q^3\gamma^3L_{so}^3C_{to}^3}\Big)\sum_{l=1}^{s_t+q-1}\eta_t(q-1)\sum_{l=1}^{t-1}\gamma^2\eta_l^2\sum_{l=1}^M\mathbb{E}\|d_l^m-\bar{d}_l\|^2$  $+\sum_{k=1}^{s_{t}+q-1}\frac{\gamma}{M\rho}\Big(4C_{g}^{2}\eta_{t}^{2}\gamma^{2}+4L_{g}^{2}\eta_{t}^{2}\gamma^{2}+8L_{f}^{2}C_{g}^{2}\eta_{t}^{2}\gamma^{2}\Big)\sum_{k=1}^{M}\mathbb{E}\|d_{t}^{m}-\bar{d}_{t}\|^{2}+\sum_{k=1}^{s_{t}+q-1}\frac{2\sigma^{2}\gamma}{\rho}\big(\alpha_{t+1}^{2}+\beta_{t+1}^{2}+\varrho_{t+1}^{2}\big)$  $\leq \sum_{t=1}^{s_t+q-1} \frac{1}{M} \sum_{t=1}^{M} \left( -\frac{C_f^2 \gamma}{2\rho} \eta_t \| u_t^m - \nabla g^m(\bar{x}_t) \|^2 - \frac{C_g^2 \gamma}{\rho} \eta_t \| v_t^m - \nabla f^m(h_t^m) \|^2 \right)$  $-\frac{\gamma}{\rho} \left(\frac{B}{2} - 4C_g^2 L_f^2 - \frac{c_2^2 C_g^2 L_f^2}{864q^3 \gamma^3 L_{t_a}^3 C_{t_a}^3} \eta_t \|h_t^m - g^m(\bar{x}_t)\|^2 \right)$  $-\sum_{q=1}^{s_t+q-1} \Big(\frac{\rho\gamma\eta_t}{2} - \frac{4C_g^2\eta_t^2\gamma^3}{\rho} - \frac{4L_g^2\eta_t^2\gamma^3}{\rho} - \frac{8L_f^2C_g^2\eta_t^2\gamma^3}{\rho}\Big) \|\bar{d}_t\|^2$  $+\frac{\gamma}{M\rho}\Big(BC_g^2+5C_f^2L_g^2+\frac{c_2^2C_g^2L_f^2}{864q^3\gamma^3L_{ta}^3C_{ta}^3}\Big)\frac{\rho^2}{(24)^2L_{ta}^2C_{ta}^2}\sum_{l=1}^{s_t+q-1}\eta_t\sum_{l=1}^M\mathbb{E}\|d_t^m-\bar{d}_t\|^2$  $+\frac{\gamma}{M\rho}\frac{\gamma\rho}{6aL_{fa}C_{fa}}\left(C_{g}^{2}+L_{g}^{2}+2L_{f}^{2}C_{g}^{2}\right)\sum_{t=1}^{s_{t}+q-1}\eta_{t}\sum_{t=1}^{M}\mathbb{E}\|d_{t}^{m}-\bar{d}_{t}\|^{2}$  $+\frac{\sigma^2}{12qL_{fq}C_{fq}}\left(c_1^2+c_2^2+c_3^2\right)\sum_{t=1}^{s_t+q-1}\eta_t^3,$ (79)

where the second inequality is due to  $\eta_t \leq \frac{\rho}{24q\gamma L_{fg}C_{fg}}$  for all  $t \geq 0$ .

1775 Let  $\gamma^2 \ge \frac{\rho^2 \sqrt{c_1^2 + c_3^2}}{24 \sqrt{30} q L_{fg}^2 C_{fg}^2}$ , we have 

 $\frac{\rho^2 C_g^2}{(24)^2 L_{fg}^2 C_{fg}^2} \frac{\rho^2 (c_1^2 + c_3^2)}{120q^2 \gamma^4 C_{fg}^2 L_{fg}^2 C_g^2} \le \frac{1}{4}.$ (80)

$$\begin{aligned} & \text{Set } \Theta = \left(5C_{f}^{2}L_{g}^{2} + \underset{i \neq q}{\frac{2}{3}}C_{f_{g}}^{2}C_{f_{g}}^{2}} + \underset{i \neq q}{\frac{2}{3}}C_{f_{g}}^{2}C_{f_{g}}^{2}} \left(C_{g}^{2} + L_{g}^{2} + 2L_{f}^{2}C_{g}^{2}\right). \text{ Based on the} \\ & \text{above Lemma 10, then we have} \\ & \text{Sinter-1} \left(\Omega_{t+1} - \Omega_{t}\right) \\ & \leq \sum_{t=n_{t}}^{n_{t}+1} \frac{1}{M} \sum_{m=1}^{M} \left(-\frac{C_{f}^{2}\gamma}{2\rho} \eta_{t} \mathbb{E} ||u_{t}^{m} - \nabla g^{m}(x_{t})||^{2} - \frac{C_{q}^{2}\gamma}{\rho} \eta_{t} \mathbb{E} ||v_{t}^{m} - \nabla f^{m}(h_{t}^{m})||^{2} \\ & \quad = \sum_{t=n_{t}}^{n_{t}+1} \frac{1}{M} \sum_{m=1}^{M} \left(-\frac{C_{f}^{2}\gamma}{864q^{3}\gamma^{2}} L_{f}^{2}C_{f}^{2}\right) \eta_{t} \mathbb{E} ||h_{t}^{m} - g^{m}(x_{t})||^{2} \\ & \quad = \sum_{t=n_{t}}^{n_{t}+1} \frac{1}{M} \sum_{m=1}^{M} \left(-\frac{C_{f}^{2}\gamma}{864q^{3}\gamma^{2}} L_{f}^{2}C_{f}^{2}\right) \eta_{t} \mathbb{E} ||h_{t}^{m} - g^{m}(x_{t})||^{2} \\ & \quad = \sum_{t=n_{t}}^{n_{t}+1} \frac{1}{M} \sum_{m=1}^{M} \left(\frac{\rho Q_{f}}{2}L_{f}^{2} - \frac{q^{2}C_{f}^{2}L_{f}^{2}}{120q^{2}\gamma^{4}} L_{f}^{2}q^{2}r^{2}\gamma^{3}}{\rho} - \frac{8L_{f}^{2}C_{f}^{2}q^{2}r^{2}\gamma}{\rho}\right) \mathbb{E} ||d_{t}||^{2} \\ & \quad + \frac{\gamma}{M\rho} \left(BC_{g}^{2} + 5C_{g}^{2}L_{g}^{2} + \frac{q^{2}C_{f}^{2}L_{f}^{2}}{120q^{2}\gamma^{4}C_{f}^{2}} L_{f}^{2}C_{f}^{2}}\right) \frac{\rho^{2}}{(24)^{2}L_{f}^{2}C_{f}^{2}} \\ & \quad \cdot \left(\frac{6M}{5} \sum_{t=n_{t}}^{t+r-1} \eta_{t}\mathbb{E} ||d_{t}||^{2} + \frac{\rho^{2}(c_{f}^{2} + c_{g}^{2})}{120q^{2}\gamma^{4}C_{f}^{2}} L_{f}^{2}C_{g}^{2}}\right) \frac{\rho^{2}}{t_{t=n_{t}}}} \\ & \quad \cdot \left(\frac{6M}{5} \sum_{t=n_{t}}^{t+r-1} \eta_{t}\mathbb{E} ||d_{t}||^{2} + \frac{\rho^{2}(c_{f}^{2} + c_{g}^{2})}{120q^{2}\gamma^{4}C_{f}^{2}} L_{f}^{2}C_{g}^{2}}\right) \frac{\rho^{2}}{t_{t=n_{t}}}} \\ & \quad \cdot \left(\frac{6M}{5} \sum_{t=n_{t}}^{t+r-1} \eta_{t}\mathbb{E} ||d_{t}||^{2} + \frac{\rho^{2}(c_{f}^{2} + c_{g}^{2})}{120q^{2}\gamma^{4}C_{f}^{2}} L_{f}^{2}} C_{g}^{2}}\right) \frac{\rho^{2}}{t_{t=n_{t}}}} \\ & \quad \cdot \left(\frac{6M}{5} \sum_{t=n_{t}}^{t+r+1} \eta_{t}\mathbb{E} ||d_{t}||^{2} + \frac{\rho^{2}(c_{f}^{2} + c_{g}^{2})}{120q^{2}\gamma^{4}C_{f}^{2}} L_{f}^{2}} C_{g}^{2}}\right) \\ & \quad \cdot \left(\frac{6M}{5} \sum_{t=n_{t}}^{t+r+1} \eta_{t}\mathbb{E} ||d_{t}||^{2} + \frac{\rho^{2}(c_{f}^{2} + c_{g}^{2})}{120q^{2}\gamma^{4}C_{f}^{2}} L_{f}^{2}} L_{g}^{2}} \eta_{t}^{2}}\right) \\ & \quad \cdot \left(\frac{6M}{5} \sum_{t=n_{t}}^{t+r+1} \eta_{t}\mathbb{E} ||d_{t}||^{2} + \frac{\rho^{2}(c_{f}^{2} + c_{g}^{2})}{120q^{2}\gamma^{4}} L_{f}^{2}}$$

where the second inequality holds by the above inequality (80), and the last inequality holds by  $B \geq 20C_g^2 L_f^2 + \frac{c_2^2 C_g^2 L_f^2}{216q^3 \gamma^3 L_f^3 g C_{fg}^3} + \frac{\Theta \rho^2 (c_1^2 + c_3^2)}{30q^2 \gamma^4 C_{fg}^2 L_f^2 g C_g^2}, \ \gamma \leq \frac{3\rho q L_{fg} C_{fg}}{4(C_g^2 + L_g^2 + 2L_f^2 C_g^2)}$ (i.e., the following inequality (82) and  $\Theta + \frac{B C_g^2 \rho^2}{(24)^2 L_{fg}^2 C_{fg}^2} \leq \frac{5\rho^2}{48}$ .

$$\begin{aligned} & \text{Since } \eta_{i} \leq \frac{3}{44p^{2}L_{12}C_{13}} \text{ and } \gamma \leq \frac{3}{4Q^{2}L_{22}^{2}L_{22}^{2}L_{22}^{2}} \frac{24q^{2}L_{12}G_{12}}{\rho} \leq \frac{\rho^{2}}{32\eta_{1}(C_{2}^{2} + L_{2}^{2} + 2L_{1}^{2}C_{2}^{2})}. \quad (82) \end{aligned} \\ & \text{Summing the above inequality (81) from  $t = 1$  to  $T$ , then we have  
$$\frac{\gamma^{2}}{t=1} \frac{\rho^{2}}{11} \frac{1}{M} \sum_{n=1}^{M} \left( -\frac{C_{2}^{2}\gamma}{2\rho} \eta_{1} \mathbb{E} \|u_{1}^{n} - \nabla g^{m}(\tilde{x}_{1})\|^{2} - \frac{C_{n}^{2}\gamma}{\rho} \eta_{1} \mathbb{E} \|v_{1}^{n} - \nabla f^{m}(h_{1}^{n})\|^{2} - \frac{\gamma^{2}C_{2}^{2}L_{2}^{2}}{\rho} \eta_{1} \mathbb{E} \|h_{1}^{n} - g^{m}(\tilde{x}_{1})\|^{2} \right) \\ & \leq \sum_{t=1}^{r} \frac{1}{M} \sum_{m=1}^{M} \left( -\frac{C_{2}^{2}\gamma}{2\rho} \eta_{1} \mathbb{E} \|u_{1}^{n} - \nabla g^{m}(\tilde{x}_{1})\|^{2} - \frac{C_{n}^{2}\gamma}{\rho} \eta_{1} \mathbb{E} \|v_{1}^{n} - \nabla f^{m}(h_{1}^{n})\|^{2} - \frac{\gamma^{2}C_{2}^{2}L_{2}^{2}}{\rho} \eta_{1} \mathbb{E} \|h_{1}^{n} - g^{m}(\tilde{x}_{1})\|^{2} \right) \\ & -\sum_{t=1}^{r} \frac{1}{PT} \sum_{m} |u_{1}^{n}|^{2} + \frac{\rho^{32}}{10PL_{2}^{2}} \sum_{t=1}^{T} \eta_{1}^{3} + \frac{\sigma^{2}}{12qL_{2}^{2}G_{1}^{2}} \left(c_{1}^{2} + c_{2}^{2} + c_{3}^{2}\right) \sum_{t=1}^{T} \eta_{1}^{3}. \qquad (83) \\ & \text{Since } h_{1}^{m} = \frac{1}{q} \sum_{s=1}^{r} \frac{1}{g} \sum_{m=1}^{r} \eta_{m}^{2} + \frac{\sigma^{2}}{12qL_{2}^{2}G_{1}^{2}} \left(c_{1}^{2} + c_{2}^{2} + c_{3}^{2}\right) \sum_{t=1}^{T} \eta_{1}^{3}. \qquad (83) \\ & \Omega_{1} = \mathbb{E} \left[F(\tilde{x}_{1}) + \frac{\gamma}{\eta_{p}} \frac{1}{M} \sum_{m=1}^{M} \left[\|h_{1}^{m} - g^{m}(x_{1}^{m})|^{2} + \||u_{1}^{m} - \nabla g^{m}(x_{1}^{m})|^{2} + \|v_{1}^{n} - \nabla f^{m}(h_{1}^{m})\|^{2}\right)\right] \\ & \leq F(\tilde{x}_{1}) + \frac{3\gamma\sigma^{2}}{q\rho p_{0}}. \qquad (84) \\ & \text{where the ats inequality bolds by Assumption 2. \\ & \text{Since } \eta_{1} = \frac{1}{q} \sum_{t=1}^{T} \mathbb{E} \left[\left[\frac{1}{M} \sum_{m=1}^{M} \frac{1}{p^{2}} \left(2C_{1}^{2}\|u_{1}^{m} - \nabla g^{m}(\tilde{x}_{1})\|^{2} + U_{1}^{2}g_{0}U_{1}^{2}} \sum_{t=1}^{T} \eta_{s}^{3} \\ & \leq \frac{4}{T\rho\gamma\eta m} \sum_{t=1}^{T} (\Omega_{1} - \Omega_{1+1}) + \frac{4\tau^{2}}{4\tau^{2}} \sum_{t=1}^{T} \sum_{t=1}^{T} \eta_{s}^{2} + \frac{(c_{1}^{2} + c_{2}^{2} + c_{3}^{2})\sigma^{2}}{3T\rho\eta m} T_{2} \int_{t=1}^{T} \eta_{s}^{3} \\ & \leq \frac{4}{T\rho\gamma\eta m} \left[F(\tilde{x}_{1}) + \frac{3\gamma\sigma^{2}}{q\rho p_{0}} - F^{*}\right) + \left(\frac{\Lambda^{2}}{4T\gamma^{2}\eta m} L_{1}^{2}} + \frac{(c_{1}^{2} + c_{2}^{2} + c_{3}^{2})\sigma^{2}}{3T\rho\eta m} T_{1} \int_{t=1}^{T} \eta_{s}^{3} \\ & \leq \frac{4}{T\rho\gamma\eta m} \left[F(\tilde{x}_{1}) +$$$$

where the first inequality holds by Lemma 5, and the last inequality holds by (85). 

Since 
$$\bar{d}_t = \frac{\bar{x}_t - \bar{x}_{t+1}}{\eta_t \gamma} = \frac{\eta_t \gamma A_t^{-1} \bar{w}_t}{\eta_t \gamma} = A_t^{-1} \bar{w}_t$$
, and let  $\mathcal{G}_t = \frac{1}{\rho} \| \bar{w}_t - \nabla F(\bar{x}_t) \| + \| \bar{d}_t \|$ , we have  

$$\mathcal{G}_t = \frac{1}{\rho} \| \bar{w}_t - \nabla F(\bar{x}_t) \| + \| \bar{d}_t \| = \frac{1}{\rho} \| \bar{w}_t - \nabla F(\bar{x}_t) \| + \| A_t^{-1} \bar{w}_t \|$$

$$\geq \| A_t^{-1} \bar{w}_t \| + \frac{1}{\rho} \| \bar{w}_t - \nabla F(\bar{x}_t) \|$$

$$= \frac{1}{\| A_t \|} \| A_t^{-1} \bar{w}_t \| + \frac{1}{\rho} \| \bar{w}_t - \nabla F(\bar{x}_t) \|$$

$$\geq \frac{1}{\| A_t \|} \| \bar{w}_t \| + \frac{1}{\rho} \| \bar{w}_t - \nabla F(\bar{x}_t) \|$$

$$\geq \frac{1}{\| A_t \|} \| \bar{w}_t \| + \frac{1}{\rho} \| \bar{w}_t - \nabla F(\bar{x}_t) \|$$

$$\geq \frac{1}{\| A_t \|} \| \bar{w}_t \| + \frac{1}{\rho} \| \bar{w}_t - \nabla F(\bar{x}_t) \|$$

$$(87)$$

where the inequality (i) holds by  $||A_t|| \ge \rho$  for all  $t \ge 1$  due to Assumption 6. Then we have 

$$\nabla F(\bar{x}_t) \| \le \|A_t\| \mathcal{G}_t.$$
(88)

According to Cauchy-Schwarz inequality, we have 

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\|\nabla F(\bar{x}_t)\| \le \frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\left[\mathcal{G}_t\|A_t\|\right] \le \sqrt{\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[\mathcal{G}_t^2]}\sqrt{\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\|A_t\|^2}.$$
(89)

According to the above inequality (86), we have 

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}[\mathcal{G}_{t}^{2}] \leq \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}\Big[\frac{2}{\rho^{2}} \|\bar{w}_{t} - \nabla F(\bar{x}_{t})\|^{2} + 2\|\bar{d}_{t}\|^{2}\Big] \\
\leq \frac{2G}{T} (n+T)^{1/3}.$$
(90)

Combining the above inequalities (89) with (90), we have 

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\|\nabla F(\bar{x}_t)\| \leq \sqrt{\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[\mathcal{G}_t^2]}\sqrt{\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\|A_t\|^2}$$

$$\leq \left(\frac{\sqrt{2G}n^{1/6}}{T^{1/2}} + \frac{\sqrt{2G}}{T^{1/3}}\right) \sqrt{\frac{1}{T} \sum_{t=1}^{T} \mathbb{E} \|A_t\|^2}.$$
(91)

### В ADDITIONAL EXPERIMENTS

### **B**.1 TASK-DISTRIBUTED META LEARNING

In this subsection, we evaluate the effectiveness of our proposed algorithms for personalized federat-ed learning, which can be described a task-distributed meta learning Huang et al. (2021a). From the above (2), the task-distributed meta learning problem can be rewritten as a distributed composition optimization problem, defined as

1942  
1943 
$$\min_{x \in \mathbb{R}^d} \frac{1}{M} \sum_{m=1}^M \exp\left(f^m \left(x - \eta \nabla f^m(x)\right) / \lambda\right), \tag{92}$$

1958 1959

1976

1978 1979



1954 Figure 3: We evaluated the test accuracy (%) of different FCO methods for solving the distributed 1955 problem (92) under various settings using a heterogeneous CIFAR-10 dataset. We varied the hyperparameter  $\chi$  to control the percentage of samples from the dominant class, with values of 0.3, 0.5, 0.7, and 0.9 from left to right. 1957



Figure 4: We computed the cross entropy on the test sets of various FCO methods for solving the 1975 distributed problem (92) using a heterogeneous CIFAR-10 dataset under different settings. The results were obtained for  $\chi = 0.7$  on the left and  $\chi = 0.9$  on the right.

where  $\lambda > 0$  is a regularization parameter. The Problem (92) is also a special case of the above 1981 Problem (1).

1982 In the experiments, we consider a multi-class classification task over the CIFAR10 Krizhevsky et al. (2009) dataset, using a 7-layer CNN. We create a heterogeneous training (validation) dataset con-1984 sisting of 10 clients and 1 server, where each client has images from a dominant class and a small percentage of images from other classes. Specifically, the *m*-th client owns  $\chi$  percentage of images 1986 from the *m*-th class and  $(1 - \chi)/9$  for the other classes. For  $\chi > 0.1$ , the images of each client 1987 are dominated by a different class, which is referred to as the dominant class. In our experiments, 1988 each client has a different dominant class, for example, one client has 60% samples from the air-1989 plane class and the remaining 40% samples from other classes. The hyper-parameter  $\chi$  controls the percentage of samples from the dominant class over each client. The dominant class is different for 1990 different clients.

1992 In the experiments, we utilize a grid search approach to determine the optimal hyper-parameters for 1993 all methods, and the search space is described in subsection B.2. As each client's data distribution is constructed to be heterogeneous, tuning a personalized model for each client can offer additional benefits. For selecting the learning rate, we typically set the learning rate to 0.05 at inner loops and the learning rate to 0.1 at outer loops. Additionally, we randomly select five clients to participate in training per epoch, and we set the asynchronization step q to 5 if it is not specified. The total number 1997 of training iteration steps is set to 600.

1998 From Figures 3 and 4, it is apparent that the convergence difficulty increases as the hyperparame-1999 ter  $\chi$  increases. Among all of the heterogeneous ratios, our MFGCD and AdaMFGCD algorithms 2000 outperform other baselines, particularly when the data distribution is significantly heterogeneous. The other composition federated optimization methods such as FEDNEST Tarzanagh et al. (2022), 2002 Local-MOML Wang et al. (2021), and Local-SCGDM Gao et al. (2022) exhibit inferior performance in both test accuracy and loss under these circumstances. Meanwhile, although ComFedL Huang 2003 et al. (2021a) obtains a better performance, the loss (accuracy) curve is quite noisy. Our MFCGD 2004 and AdaMFCGD algorithms adaptively adjust the weight of clients based on their performance on the task (training loss). In other words, if a client's data distribution is challenging to learn (i.e., 2006 higher training loss), the algorithm increases its learning rate, while for clients with simpler distri-2007 butions, the learning rate is reduced. Overall, the results indicate that our MFCGD and AdaMFCGD 2008 approaches can also enhance personalized FL. 2009



Figure 5: Comparing the accuracy(%) (left) and cross-entropy loss (right) on different regularization parameter  $\lambda$  on our AdaMFCGD algorithm

2029 Figure 5 demonstrates the robustness of our AdaMFCGD algorithm by varying the regularization 2030 parameter  $\lambda$ . The results show that our AdaMFCGD algorithm achieves good test accuracy and 2031 low test loss for different values of  $\lambda$ . Additionally, decreasing  $\lambda$  leads to faster convergence of the algorithm, particularly when running multiple local epochs. Figure 6 illustrates the robustness 2032 of our algorithm when varying the asynchronization step q. It is noteworthy that our AdaMFCGD achieves its optimal performance when q = 1, as it enables the momentum-based variance reduction technique to fully demonstrate its potential by calculating an adaptive matrix in each iteration. In 2035 comparison to varying q in Robust FL, q shows a less significant influence in Task-Distributed Meta 2036 Learning due to the heterogeneity of the data, which can result in significant changes in the gradient across iterations, further impacting the stability of the convergence curve. In Task-Distributed Meta 2038 Learning, asynchronization step q also relatively controls the degree of heterogeneity.

Layer Type	Output Size	Kernel Size	Stride	Activation
Input	28 x 28 x 1	-	-	-
Convolution	24 x 24 x 6	5 x 5	1	ReLU
Max Pooling	12 x 12 x 6	2 x 2	2	-
Convolution	8 x 8 x 16	5 x 5	1	ReLU
Max Pooling	4 x 4 x 16	2 x 2	2	-
Convolution	120	5 x 5	1	ReLU
Dense	360	-	-	ReLU
Output	10	-	-	Softmax

2049 2050

2010

2026

Table 2: Structure of a 4-layer CNN for MNIST



Figure 6: Comparing the accuracy(%) (left) and cross-entropy loss (right) on different synchronization step q on our AdaMFCGD algorithm.

Layer Type	Output Size	Kernel Size	Stride	Activation
Input	32 x 32 x 3	-	-	-
Convolution	30 x 30 x 96	3 x 3	1	ReLU
Convolution	14 x 14 x 96	3 x 3	2	ReLU
Convolution	14 x 14 x 196	1 x 1	1	ReLU
Convolution	14 x 14 x 10	1 x 1	1	ReLU
Flatten	1,960	-	-	-
Dense	1000	-	-	ReLU
Dense	1000	-	-	ReLU
Output	10	-	-	Softmax

Table 3: Structure of a 7-layer CNN for CIFAR-10

2083 B.2 IMPLEMENTATION DETAILS

In this subsection, we provide the specific backbone networks of the above two tasks, which are described in Table 2 and Table 3, respectively.

In the above **Robust Federated Learning** experiments, we conduct a search for the learning rate within the range [0.001, 0.01, 0.05, 0.1, 0.2, 0.5, 1]. We observe that for most methods, a learning rate greater than 0.5 caused divergence. In the case of ComFedL Huang et al. (2021a), we conduct a search for the regularization parameter within the range [0.1, 0.5, 1.5, 2] and 0.5 is the best. For FEDNEST Tarzanagh et al. (2022), we directly use the hyperparameters as it reports. We also conduct a search for the hyperparameter of Local-SCGDM Gao et al. (2022) within the range [0.1, 0.5, 1, 1.5, 2] and 1 is the best.

2094 In the above Task-Distributed Meta Learning experiments: for the learning rate (both inner and outer if two types of learning rates are needed), we search from [0.001, 0.01, 0.05, 0.1, 2095 0.2, 0.5, 1]. For our method, we search the regularization parameter from [0.1, 0.5, 1, 5]; For 2096 FEDNEST Tarzanagh et al. (2022), we directly use the hyperparameters as it reports. We also con-2097 duct a search for the hyperparameter of Local-SCGDM Gao et al. (2022) within the range [0.1, 0.5, 2098 1, 1.5, 2] and 0.5 is the best. For Local-MOML Wang et al. (2021), we search its weight parameter 2099  $\beta$  within the range [0.1, 0.3, 0.5, 0.7, 0.9] with the fixed inner and outer learning rates for a fair 2100 comparison and  $\beta = 0.7$  is the best. 2101

2102

2067

2068

2081 2082

2084

2103

2104