# BIDIRECTIONAL COMMUNICATION-EFFICIENT NON CONVEX ADAPTIVE FEDERATED LEARNING

Anonymous authors

Paper under double-blind review

# ABSTRACT

Within the framework of federated learning, we introduce two novel strategies: New Lazy Aggregation (NLA) and Accelerated Aggregation (AA). The NLA strategy reduces communication and computational costs through adaptive gradient skipping, while the AA strategy accelerates computation and decreases communication costs via adaptive gradient accumulation. Building upon these innovative strategies and compression techniques, we propose two new algorithms: FedBN-LACA and FedBACA, aimed at minimizing bidirectional communication costs. We provide theoretical guarantees for client participation (either full or partial) in these algorithms under non-convex settings and heterogeneous data. In the context of non-convex optimization with full client participation, our proposed FedBN-LACA and FedBACA algorithms achieve the same convergence rate of  $\mathcal{O}(1/T)$  as their non-tight counterparts. Extensive experimental results demonstrate that our protocols facilitate effective training in non-convex environments and exhibit robustness across a wide range of devices, partial participation, and imbalanced data.

024 025 026

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

# 1 INTRODUCTION

027 028

029 Federated Learning (FL) is a promising machine learning (ML) framework that enables collaborative model training without sharing data. This approach ensures that data privacy is protected, as data does not need to leave the owner's possession McMahan et al. (2017). FL models involve edge clients 031 or devices, such as smartphones and personal computers (PCs), training local ML models using their own data without sharing. Instead of sending raw data to a central server (e.g. cloud server), 033 the edge client updates the parameters of the local server model (e.g. smartphone, PC updates). 034 The server then generates a shared global ML model by aggregating local parameter updates. In FL, traditional distributed stochastic gradient descent (SGD) is unsuitable for challenging federated learning scenarios, where data does not follow independent homogeneous distributions and only 037 a small fraction of clients participate in each communication round. To address this issue, many 038 federated optimization methods use local client updates. To reduce significantly the amount of communication required to train the model, local clients update their models multiple times before communicating with the server. One of the most popular FL methods is the FedAvg algorithm 040 proposed by McMahan et al. (2017), in which the global model is updated by averaging over multiple 041 local SGD update steps. 042

Although FedAvg algorithm can be trained without sharing data and achieve good results, it in practice still presents challenges in FL. 1) Lack of adaptivity. As SGD-based updates may not be suitable for stochastic gradient noise with heavy-tailed distributions, which usually occur when training large models Devlin et al. (2019); Brown et al. (2020); 2) Unaffordable communication. Repeated synchronization of the uplink and downlink between the client and the server leads to significant communication overhead, but some of these parameters are passed non-essential. Here, the uplink represents a transmission from clients to servers and the downlink represents a transmission from server to clients.

Many works have researched the aforementioned issues. 1) For adaptivity in the federated learning framework, Reddi et al. (2020) proposed the FedAdam algorithm and its variants Tong et al. (2020);
Wang et al. (2022a); Wu et al. (2023), which integrate the adaptive gradient approach. 2) Recently, to reduce communication costs, three approaches have been developed: (i) Local approach Huang

054 et al. (2023); Li & Wang (2019); Li et al. (2019b); Mishchenko et al. (2022); Karimireddy et al. 055 (2019a). The local devices implement multiple rounds of optimization before sending the information 056 to the server. This reduces the rounds of global communication, thereby having an overall lower 057 communication cost. (ii) Compression methods Reisizadeh et al. (2020); Chen et al. (2021); Richtárik 058 et al. (2023); Beznosikov et al. (2023). During each round of communication, the local devices send compressed information to the server, reducing communication costs through some compression mechanisms. (iii) Lazy aggregation algorithmChen et al. (2018b); Sun et al. (2019); Ghadikolaei et al. 060 (2021); Mishchenko et al. (2022). Some parameters may be similar to those of the previous round, so 061 they do not need to be transmitted (i.e., the parameters of the previous round are used in the current 062 round), which also effectively reduces the communication cost. 063

064 However, although the lazy aggregation algorithm provides a strategy to reduce communication costs, it is more difficult to implement in practice, especially in joint learning where only some clients 065 are involved in the training. Although these algorithms have achieved great success in reducing 066 communication costs, they are unidirectional in the sense that they are uplink communication, not 067 downlink communication. A natural thought is: How do we design a simple and more efficient 068 strategy to improve communication efficiency based on bidirectional communication? In this 069 paper, we initially introduce two innovative updated strategies, namely NLA and AA. These strategies are subsequently applied to the FedAMS and FedCAMS algorithms as outlined in Wang et al. 071 (2022a). Furthermore, we advance the development of bidirectional communication-efficient adaptive 072 algorithms within a non-convex framework and in the context of heterogeneous data. 073

### Main contributions 074

075 • We present two novel communication strategies for joint learning in federated learning (FL), 076 referred to as NLA and AA, which improve upon the existing LAG algorithmChen et al. (2018b); 077 Sun et al. (2019); Mishchenko et al. (2022). In contrast to conventional approaches, NLA and AA do not require parameters from multiple iterations; rather, they utilize only the parameters from the preceding iteration along with those from the current iteration. This characteristic enhances their 079 adaptability and simplifies their implementation. 080

081 • We propose two novel and efficient adaptive optimization techniques for communication: 082 FedNLACA and FedACA. Both methods enhance communication efficiency and adaptability within 083 the context of joint learning. The FedNLACA algorithm notably decreases communication costs by employing real error feedback, a compression strategy, and our proposed NLA strategy, all while 084 preserving high accuracy. Similarly, the FedACA algorithm effectively reduces communication 085 costs and sustains high accuracy through the use of real error feedback, a compression strategy, 086 and our proposed AA strategy. It is important to highlight that FedNLACA and FedACA attains a 087 convergence rate of  $\mathcal{O}(1/T)$ . 088

• We have developed a novel cross-device compatible adaptive joint optimization method, referred to 089 as FedBNLACA and FedBACA, which employs two strategic approaches to facilitate a reduction in bidirectional communication costs for both uplink and downlink transmissions. A convergence analy-091 sis was performed under conditions of general non-convexity and data heterogeneity, demonstrating 092 that both FedBNLACA and FedBACA can also attain a convergence rate of  $\mathcal{O}(1/T)$ . 093

• Extensive experiments on various benchmarks showed that our proposed algorithms are well adaptive in training real-world machine learning models.

096

094

#### 2 **RELATED WORK**

098

Adaptive Gradient Methods: Adaptive gradients (Zeiler, 2012; Duchi et al., 2011; Kingma & 100 Ba, 2014; Reddi et al., 2018) are a series of algorithms that effectively reduce the relatively slow 101 convergence and over-sensitivity to parameters of gradient descent methods in the face of heavy-tailed 102 stochastic gradients, and are heavily used to train large networks. 103

104 Federated Learning: FedAvg was introduced by McMahan et al. (2017) as the inaugural algorithm 105 for federated learning (FL). By employing periodic model averaging, this approach significantly mitigates communication overheads. Initial studies focused on the analysis of FL algorithms within a 106 homogeneous data framework, while more recent investigations have expanded the scope of federated 107 learning to encompass heterogeneous data environments (non-iid) and non-convex models (Li & 108 Wang, 2019; Li et al., 2019a; Sahu et al., 2018; Yang et al., 2019; Wang et al., 2022a)). (Li et al., 109 2019a; Sahu et al., 2018) proposed FedProx and FedDANE algorithms for federated optimization 110 against heterogeneity. Wang et al. (2022a)gave theoretical results of data heterogeneous federated 111 learning under the conditions of adaptive methods. In this paper, we follow the ideas of Wang et al. 112 (2022a). Numerous studies have been conducted that build upon the FedAvg framework, including notable contributions such as FedNova Wang et al. (2020), and SCAFFOLD Karimireddy et al. 113 (2019a), as well as various other adaptations of FedAvg (Yang et al., 2021b; Wang et al., 2022a). 114 In a recent advancement, Reddi et al. (2020) introduced several adaptive federated optimization 115 techniques, including FedAdagrad, FedYogi, and FedAdam, aimed at addressing the convergence 116 challenges associated with FedAvg. Additionally, Chen et al. (2020) presented Local AMSGrad, 117 while Tong et al. (2020) introduced a suite of federated adaptive gradient methods that incorporate 118 calibration mechanisms. 119

Effective methods of communication: Many methods to reduce communication costs have been 120 proposed in federated learning, three main ideas are 1) multiple rounds of local iteration methods 121 (Elgabli et al., 2022; McMahan et al., 2017; Li et al., 2019a; Sahu et al., 2018), 2) compression and 122 error feedback (Reisizadeh et al., 2020; Elgabli et al., 2022; Haddadpour et al., 2020; Wang et al., 123 2022a), 3) lazy aggregation (Sun et al., 2020; Mao et al., 2022; Chen et al., 2018b; Sun et al., 2019). 124 Most of the above methods are only unidirectional (uplink), and although (Sun et al., 2020; Wang 125 et al., 2022b) considered bidirectional algorithms, they did not consider the heterogeneity problem in 126 FL, the involvement of some of the servers in the training as well as the related theoretical guarantee, 127 and the adaptive problem. Moreover, their strategies to reduce the communication cost are difficult to 128 implement under heterogeneous partial clients participation. Therefore, how to give simple strategies 129 to reduce communication cost, and how to design adaptively efficient algorithms for bidirectional communication with partial clients participation is the focus of this paper. 130

131 Notation: For vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ ,  $\sqrt{\mathbf{x}}, \mathbf{x}^2, \mathbf{x}/\mathbf{y}$  denote the *n* element-wise square root, square, and 132 division of the vectors. For vector  $\mathbf{x}$  and matrix  $A, \|\cdot\|$  denotes the  $\ell_2$  norm of vector/matrix, i.e., 133  $\|\mathbf{x}\| = \|\mathbf{x}\|_2$  and  $\|A\| = \|A\|_2$ . In algorithms, *t* denotes the *t*-th iteration. The *m* is the number of all 134 clients/devices.

# 3 PRELIMINARY

136

137 138

139

140 141

142

151 152 This paper aims to study the federated learning non-convex optimization problem, which is formulated as follows:

$$\min_{\theta \in \mathbb{R}^d} f(\theta) = \frac{1}{m} \sum_{i=1}^m F_i(\theta),$$

143 where  $F_i(\theta) = \mathbb{E}_{\xi_i \sim \mathcal{D}_i} F_i(\theta, \xi_i)$ ,  $F_i(\theta)$  denotes the local non-convex loss,  $\mathcal{D}_i$  represents the data 144 distribution on *i* clients. *m* represents the number of all clients, *d* denotes the dimension of the model 145 parameters. In the non i.i.d setting, distributions  $\mathcal{D}_i$  and  $\mathcal{D}_j$  can vary from each other, i.e.,  $\mathcal{D}_i \neq \mathcal{D}_j$ , 146  $\forall i \neq j$ . In the stochastic setting, one can only get unbiased estimates of  $F_i(\theta)$ , i.e., the stochastic 147 gradient  $\mathbf{g}_t^i = \nabla F_i(\theta, \xi_i)$ .

148 Assumption 3.1. (Smoothness). There exists an L such that each loss function on the *i*-th worker 149  $F_i(\theta)$  satisfies the following equation,  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ , 150

$$|F_i(\mathbf{x}) - F_i(\mathbf{y}) - \langle \nabla F_i(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle| \le \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

The smoothness of  $F_i$  also means that the *L*-gradient Lipschitz condition, i.e.,  $\|\nabla F_i(\mathbf{x}) - \nabla F_i(\mathbf{y})\| \le L\|\mathbf{x} - \mathbf{y}\|$ . This assumption is widely used Yang et al. (2021a); Reddi et al. (2018).

Assumption 3.2. (Bounded Gradient). Each loss function on the *i*-th worker  $F_i(\theta)$  has *G*-bounded stochastic gradient on  $\ell_2$ , i.e., for all  $\xi$ ,  $E || \nabla f_i(\theta, \xi) || \le G$ . In addition, we also assume that  $||\theta|| \le$ *H*. The assumption of bounded gradient is usually adopted in adaptive gradient methods Reddi et al. (2018); Chen et al. (2018a).

**Assumption 3.3.** (Bounded Variance). The bounded local variance, i.e. for all  $\theta$ ,  $i \in [m]$ ,  $\mathbb{E}[\|\nabla f_i(\theta,\xi) - \nabla F_i(\theta)\|^2] \le \sigma_l^2$ ; and global variance constraint, i.e.  $\frac{1}{m} \sum_{i=1}^m \|\nabla F_i(\theta) - \nabla f(\theta)\|^2 \le \sigma_a^2$ , where  $\sigma_l^2$  and  $\sigma_a^2$  are some positive constants. The assumption of bounded variance also is usually adopted in adaptive gradient methods Yang et al. (2021a); Reddi et al. (2020). The bounded local variance represents the randomness of stochastic gradients, while the bounded global variance represents the heterogeneity of data between clients. It is important to note that these variances are bounded. The value of  $\sigma_g = 0$  indicates the i.i.d setting, where datasets from each client have the same distribution.

Assumption 3.4. (Biased Compressor). Consider a biased operator  $C : \mathbb{R}^d \to \mathbb{R}^d$ : for  $\forall \theta \in \mathbb{R}^d$ , there exists constant  $0 \le q \le 1$  such that

$$\|\mathcal{C}(\theta) - \theta\| \le q \|\theta\|, \forall \theta \in \mathbb{R}^d.$$

Note that q = 0 means no compression to  $\theta$ . Here are two examples: scaled-sign compressor and top-k compressor.

**Top-***k* Shi et al. (2019); Stich et al. (2018): For  $1 \le k \le d$  and  $\forall \theta \in \mathbb{R}^d$ , the coordinate of  $\theta$  is ordered by the magnitude  $|\theta_{(1)}| \le |\theta_{(2)}| \le \cdots \le |\theta_{(d)}|$ . Denote  $h_1, h_2, \dots, h_d$  as standard unit basis vectors in  $\mathbb{R}^d$ . The compressor  $\mathcal{C}_{top} : \mathbb{R}^d \to \mathbb{R}^d$  is defined as:  $\mathcal{C}_{top}(\theta) = \sum_{i=d-k+1}^d \theta_{(i)} h_{(i)}$ .

177 Define the compression ratio as r = k/d. It can be shown that  $||C_{top}(\theta) - \theta||^2 \le (1 - k/d) ||\theta||^2 = (1 - r) ||\theta||^2$ , and thus we have  $q = \sqrt{1 - r}$ .

179 Scaled sign Karimireddy et al. (2019b): For  $1 \le k \le d$  and  $\forall \theta \in \mathbb{R}^d$ , the compressor  $\mathcal{C}_{sign} : \mathbb{R}^d \to \mathbb{R}^d$  is defined as

$$C_{\text{sign}}(\theta) = \|\theta\|_1 \cdot \text{sign}(\theta)/d$$

For scaled sign compressor, when  $\|C_{\text{sign}}(\theta) - \theta\|^2 = (1 - \|\theta\|_1^2/d\|\theta\|^2)\|\theta\|^2$ , thus  $q = \sqrt{1 - \|\theta\|_1^2/d\|\theta\|^2}$ .

## 3.1 TWO STRATEGIES

Prior to the presentation of my strategy, it is essential to examine the LAG algorithmChen et al. (2018b); Sun et al. (2019); Mishchenko et al. (2022):

$$\left\|\nabla F_m(\boldsymbol{\theta}_m^{t-1}) - \nabla F_m(\boldsymbol{\theta}^t)\right\|^2 \le \frac{1}{\alpha^2 m^2} \sum_{r=1}^R \xi_R \left\|\boldsymbol{\theta}^{t+1-R} - \boldsymbol{\theta}^{t-R}\right\|^2,\tag{1}$$

204 205 206

210 211

212

182

183

184 185 186

187

188

189 190

169 170

 $L_m^2 \left\| \boldsymbol{\theta}_m^{t-1} - \boldsymbol{\theta}^t \right\|^2 \le \frac{1}{\alpha^2 m^2} \sum_{r=1}^R \xi_R \left\| \boldsymbol{\theta}^{t+1-R} - \boldsymbol{\theta}^{t-R} \right\|^2.$ (2)

While the aforementioned concept is innovative, it necessitates the establishment of parameters for the initial R iterations. Determining the appropriate value for R poses a challenge, and furthermore, the inclusion of parameters from all R iterations in the computation may lead to complications in data storage. Therefore, we propose two new aggregation strategies.

Here we introduce the meanings of some parameter representations. The  $S_t$  denotes the sum of all participating training clients at the *t*-th iteration,  $C, D, \alpha$  are some postive constants, *m* represents the number of all clients, here we first introduce the meanings of some parameter representations, R represents the number of iterations selected,  $L_m$  represents the *L*-gradient Lipschitz condition of *m*.

**NLA Strategy** (New Lazy Aggregation). For any x, let  $\rho_t = x_t - x_{t-1}$ . If

$$\|\rho_t\| \le \frac{C}{\alpha S_t} \|\mathbf{x}_{t-1}\| : \mathbf{x}_t \leftarrow \mathbf{x}_{t-1}, else : \mathbf{x}_t \leftarrow \mathbf{x}_t.$$
(3)

Example 1: Let  $q_t^i = \Delta_t^i - \Delta_{t-1}^i$ , if  $||q_t^i|| \le \frac{C}{\alpha S_t} ||\Delta_{t-1}^i||$ ,  $i \in M_t : \tilde{\Delta}_t^i \leftarrow \Delta_{t-1}^i$ , else:  $\tilde{\Delta}_t^i \leftarrow \Delta_t^i$ , the specific parameters are given in Algorithm 1.

**Example 2:** Let 
$$C(q_t^i) = \widehat{\Delta}_t^i - \widehat{\Delta}_{t-1}^i$$
, if  $\|C(q_t^i)\| \leq \frac{C}{\alpha S_t} \|\widehat{\Delta}_{t-1}^i\|, i \in M_t : \widehat{\widehat{\Delta}}_t^i \leftarrow \widehat{\Delta}_{t-1}^i$ , else:

 $\Delta_t \leftarrow \Delta_t^i$ , the specific parameters are given in Algorithm 3.

Example 3: Let 
$$C(Q_t^i) = \hat{\theta}_t^i - \hat{\theta}_{t-1}^i$$
, if  $\|C(Q_t^i)\| \leq \frac{C}{\alpha S_t} \|\hat{\theta}_{t-1}^i\|, i \in M_t : \hat{\theta}_t^i \leftarrow \hat{\theta}_{t-1}^i$ , else:  
 $\hat{\theta}_t^i \leftarrow \hat{\theta}_t^i$ , the specific parameters are given in Algorithm 2.

216 *Remark* 3.1. The NLA Strategy presented herein represents a modification of the lazy aggregation 217 method as described in previous works Chen et al. (2018b); Sun et al. (2019). This approach is 218 characterized by its operational simplicity, necessitating only a comparison with the parameters from 219 the preceding iteration. As an innovative lazy aggregation technique, the NLA Strategy effectively 220 diminishes communication costs by minimizing the number of communication parameters required. (The NLA strategy assesses whether the difference between the parameters from the (t - 1)th and 221 t-th iterations falls within a minimal threshold, indicating that these parameters are closely aligned. 222 If this proximity is confirmed, the parameter from the (t-1)th iteration is retained in place of the 223 parameter from the *t*-th iteration.) 224

AA Strategy (Accelerated Aggregation). For any x, let  $\rho_t = x_t - x_{t-1}$ . if

228 229 230

231 232

233 234

235 236

237 238

248 249

250 251

252

253 254

255

256

257

258

259

260

261

262

225

226

227

**Example 4:** Let  $q_t^i = \Delta_t^i - \Delta_{t-1}^i$ , if  $\|q_t^i\| \leq \frac{D}{\alpha S_t} \|\Delta_{t-1}^i\|, i \in M_t : \tilde{\Delta}_t^i \leftarrow \Delta_{t-1}^i + \Delta_t^i$ , else:  $\tilde{\Delta}_t^i \leftarrow \Delta_t^i$ , the specific parameters are given in Algorithm 1

 $\|\rho_t\| \leq \frac{D}{\alpha S_t} \|\mathbf{x}_{t-1}\| : \mathbf{x}_t \leftarrow \mathbf{x}_{t-1} + \mathbf{x}_t, else : \mathbf{x}_t \leftarrow \mathbf{x}_t.$ 

(4)

**Example 5:** Let  $C(q_t^i) = \widehat{\Delta}_t^i - \widehat{\Delta}_{t-1}^i$ , if  $\|C(q_t^i)\| \le \frac{D}{\alpha S_t} \|\widehat{\Delta}_{t-1}^i\|, i \in M_t : \widehat{\widehat{\Delta}}_t^i \leftarrow \widehat{\Delta}_{t-1}^i + \widehat{\Delta}_t^i$ , else:

 $\widehat{\hat{\Delta}}_t^{\iota} \leftarrow \widehat{\Delta}_t^i$ , the specific parameters are given in Algorithm 3.

**Example 6:** Let  $C(Q_t^i) = \widehat{\theta}_t^i - \widehat{\theta}_{t-1}^i$ , if  $\|C(Q_t^i)\| \leq \frac{D}{\alpha S_t} \|\widehat{\theta}_{t-1}^i\|, i \in M_t : \widehat{\widehat{\theta}}_t^i \leftarrow \widehat{\theta}_{t-1}^i + \widehat{\theta}_t^i$ , else:

 $\hat{\theta}_t \leftarrow \hat{\theta}_t^i$ , the specific parameters are given in Algorithm 2.

239 *Remark* 3.2. The AA strategy is a novel acceleration method designed to reduce communication 240 costs by speeding up the process. The proposal of this strategy is based on a fundamental motivation: 241 when the parameters of the (t-1)th iteration are very close to the parameters of the tth iteration, it is 242 possible to achieve an update of both steps at once (i.e., by adding the parameters of the (t-1)th 243 and tth iterations). This core idea is consistent with the principles of the NLA algorithm. Through 244 adaptive accelerated iterative descent, the AA strategy can effectively reduce communication costs. It 245 is worth noting that even without the AA strategy, the iterative process of conventional algorithms 246 can still achieve results similar to those of the AA strategy, but a more in-depth analysis of the AA 247 strategy will be reserved for future research.

#### **METHODS** 4

FEDERATED NEW LAZY AGGREGATION AMSGRAD AND FEDERATED ACCELERATION 4.1 AMSGRAD

In this section, we present two new frameworks of the adaptive algorithm: Federated New Lazy Aggregation AMSGrad (FedNLAA) and Federated Accelerated AMSGrad (FedAA). In FedNLAA and FedAA algorithms.  $\theta_t$  is the t-th iteration t of the global model parameters. At iteration t, the participating client i in the selected subset  $S_t$  (of size n) receives the model  $\theta_t$  from the server, i.e.,  $\theta_{t,0}^i = \theta_t$ . Then, the client performs K steps of local SGD updating with local learning rate  $\eta_l$  to get the local model  $\theta_{t,K}^i$ , judges whether the client *i* model difference  $\Delta_t^i = \theta_{t,K}^i - \theta_t$  satisfies the strategy NLA (in FedNLAA) or AA (in FedAA), and then send the judged model difference  $\Delta_t^i$  to the server. The server updates the global model difference  $\Delta_t$  by simply averaging the local model differences  $\widehat{\Delta}_{t}^{i}$ . Algorithm 1 gives the detailed procedure for FedNLAA and FedAA.

263 264 265

4.2 CONVERGENCE ANALYSIS FOR FEDNLAA AND FEDAA

**Full Participation:** All clients participate in training, i.e.,  $|S_t| = m, \forall t \in [t]$ .

267 268 269

266

**Theorem 4.1.** Under Assumptions 3.1-3.3, if learning rate  $\eta_1$  satisfies the following condition:  $\eta_1 \leq \eta_2$  $\min\left\{\frac{1}{8KL}, \frac{1}{KL}\right\}$ rithm 1 under the full

$$\frac{\epsilon}{\sqrt{\beta_2 K^2 G^2 + \epsilon} [(3 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G]} \right\}, \text{ then FedLAA in Algor}$$

#### 270 Algorithm 1 FedNLAA and FedAA 271 **Input:** Initial value $\theta_1$ , local step size $\eta_l$ , global step size $\eta$ , constants $\beta_1$ , $\beta_2$ and $\epsilon$ , $\Delta_0^i =$ 272 0 273 1: $\mathbf{m}_0 \leftarrow 0, \mathbf{v}_0 \leftarrow 0$ 274 2: for t = 1 to T do 275 Server randomly selects a subset of clients $S_t$ and transmits $\theta_t$ to the subset of clients $S_t$ . 3: 276 4: $\theta_{t,0}^i \leftarrow \theta_t$ 277 for each client $i \in S_t$ in parallel do 5: 278 for k = 0, ..., K - 1 do 6: Compute local stochastic gradient: 279 7: $\begin{aligned} \mathbf{g}_{t,k}^{i} &= \nabla F_{i}(\theta_{t,k}^{i};\xi_{t,k}^{i}), \\ \text{Update } \theta_{t,k+1}^{i} &= \theta_{t,k}^{i} - \eta_{l}\mathbf{g}_{t,k}^{i}. \end{aligned}$ 281 8: 9: end for Compute $\Delta_t^i = \theta_{t,K}^i - \theta_t$ , $q_t^i = \Delta_t^i - \Delta_{t-1}^i$ , 10: 284 Judges: If $q_t^i$ satisfies NLA (Example 1) or AA (Example 4). 11: 285 Outputs: $\Delta_t^i$ . 12: 286 13: end for Server aggregates: $\tilde{\Delta}_t = \frac{1}{|S_t|} \sum_{i \in S_t} \tilde{\Delta}_t^i$ , 287 14: 288 Update: $\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \tilde{\Delta}_t$ , 15: 289 Update: $\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \tilde{\Delta}_t^2$ , $[\mathbf{\hat{v}}_t = \max(\mathbf{\hat{v}}_{t-1}, \mathbf{v}_t, \epsilon) \text{ and } \theta_{t+1} = \theta_t + \eta \frac{\mathbf{m}}{\sqrt{\mathbf{v}_t}}]$ , 16: 290 291 $[\widehat{\mathbf{v}}_t = \max(\widehat{\mathbf{v}}_{t-1}, \mathbf{v}_t) \text{ and } \theta_{t+1} = \theta_t + \eta \frac{\mathbf{m}}{\sqrt{\mathbf{v}_t + \epsilon}}].$ 292 293 17: end for 294

participation has

295 296

304 305 306

307 308 309

323

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \leq 4\sqrt{\beta_2 \eta_l^2 K^2 G^2 + \epsilon} \cdot \left[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\right],$$
where  $\Xi = \frac{C_1 G^2 d}{\sqrt{\epsilon}} + \frac{2C_1^2 \eta \eta_l KL G^2 d}{\epsilon}, \Omega = \frac{5\eta_l^2 K^2 L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_g^2) + (3 + C_1^2)\eta^2 L + 2\sqrt{2(1 - \beta_2)}\eta G] (\frac{2\eta_l}{m\eta\epsilon} \sigma_l^2 + \frac{2KC^2 \eta_l G^2}{\alpha^2 \eta m^2 \epsilon}) + \frac{\sqrt{2}GC}{\alpha m\epsilon}, and C_1 = \frac{\beta_1}{1 - \beta_1}.$ 
Theorem 4.2. Under Assumptions 3.1-3.3, if learning rate  $\eta_l$  satisfies the following condition:  $\eta_l < 1$ 

**I neorem 4.2.** Under Assumptions 3.1-3.3, if learning rate  $\eta_1$  satisfies the following condition:  $\eta_l \leq \min\left\{\frac{1}{8KL}, \frac{\epsilon}{K\sqrt{\beta_2 K^2 G^2 + \epsilon}[(3+C_1^2)\eta L + 2\sqrt{2(1-\beta_2)}G]}\right\}$ , then FedAA in Algorithm 1 under the full participation has

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 4\sqrt{\beta_2 \eta_l^2 K^2 G^2 + \epsilon} \cdot \left[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\right]$$

 $\begin{array}{rcl} \textbf{310} \\ \textbf{311} \\ \textbf{312} \\ \textbf{312} \\ \textbf{313} \end{array} & \textbf{where} \ \Xi &= \ \frac{C_1 G^2 d}{\sqrt{\epsilon}} \ + \ \frac{2C_1^2 \eta \eta_l K L G^2 d}{\epsilon}, \Omega \\ \textbf{312} \\ 2\sqrt{2(1-\beta_2)} \eta G](\frac{2\eta_l}{m\eta\epsilon} \sigma_l^2 + \frac{2K \eta_l G^2}{\eta\epsilon}) + \frac{\sqrt{2}G}{\epsilon}, \text{ and } C_1 = \frac{\beta_1}{1-\beta_1}. \end{array}$ 

*Remark* 4.1. When the parameters C = D,  $\frac{C}{\alpha m} = 1$ , the result of Theorem 4.1 becomes the result of Theorem 4.2. The upper bound for  $\min_{t \in [T]} \mathbb{E}[|\nabla f(\theta_t)||^2]$  contains three parts: The first two terms decrease as T increases, and this term tends to zero as t tends to infinity. The last term relates to the local stochastic variance  $\sigma_l$  and global variance  $\sigma_g$ . In the i.i.d setting, where the global variance is zero and each worker has the same data distribution, i.e.,  $\sigma_g = 0$ , the variance term  $\Omega$  will be smaller.

**Corollary 4.1.** Suppose choose local learning rate  $\eta_l = \Theta(\frac{1}{\sqrt{T}K})$  and global learning rate  $\eta = \Theta(\sqrt{Km})$ , when *T* is sufficiently large, i.e.,  $T \ge Km$ , the convergence rate for FedNLAA and FedAA in Algorithm 2 under full participation has

$$\min_{t \in [T]} \mathbb{E} \left[ \|\nabla f(\theta_t)\|^2 \right] = \mathcal{O} \left( \frac{1}{\sqrt{TKm}} \right)$$

*Remark* 4.2. Corollary 4.1 suggests that with sufficient large T, when  $T = \mathcal{O}(Km)$ , FedNLAA and FedAA achieve a convergence rate of  $\mathcal{O}(\frac{1}{T})$ , which matches the result for general federated non-convex optimization methods such as FedAMS Wang et al. (2022a) and FedAdam Reddi et al. (2020).

**Partial Participation:** We assume that only n of m workers participate the local updating and communicate with the central server on each step t, i.e.,  $|S_t| = n, \forall t \in [1, T]$ . The partial participation includes the randomness of sampling, and the coefficient varies for different sampling methods. Here we consider the random sampling without replacement. At the *t*-th iteration, we randomly sample a subset  $S_t$  contains n workers for local updating, for any two workers  $i, j \in S_t$ , the probability of being sampled to participate in the model update are  $\mathbb{P}\{i \in S_t\} = \frac{n}{m}$  and  $\mathbb{P}\{i, j \in S_t\} = \frac{n(n-1)}{m(m-1)}$ .

**Theorem 4.3.** Under Assumptions 3.1-3.3, if  $\eta_l$  satisfies:  $\eta_l \leq \min\left\{\frac{1}{8KL},\right\}$ 

 $\left.\frac{n(m-1)\epsilon}{48m(n-1)K\sqrt{\beta_2K^2G^2+\epsilon\cdot[3\eta L+C_1^2\eta L+2\sqrt{2(1-\beta_2)}G]}}\right\}, \ then \ FedLAA \ in \ Algorithm \ 1 \ under \ partial$ participation has

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 8\sqrt{\beta_2 \eta_l^2 K^2 G^2 + \epsilon} \left[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\right]$$

where  $\Xi = \frac{C_1 G^2 d}{\sqrt{\epsilon}} + \frac{2G_1^2 \eta \eta K L G^2 d}{\epsilon}, \quad \Omega = \frac{5\eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_q^2) + [(3 + C_l^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{n\eta\epsilon}\sigma_l^2 + \frac{2\eta_l C^2 K^2 G^2}{\alpha^2 \eta n^2 \epsilon}) + [(3 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G]\frac{\eta_l(m-n)}{2n(m-1)\epsilon} [15K^2 L^2 \eta_l^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + [(3 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G]\frac{\eta_l}{2n(m-1)\epsilon} [15K^2 L^2 \eta_l^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + [(3 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) + C_1^2 (\sigma_l^2 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{2n(m-1)\epsilon}) +$  $6K\sigma_g^2) + 3K\sigma_g^2]\frac{1}{\eta\eta_l K} + \frac{\sqrt{2}GC}{\alpha n\epsilon} \text{ and } C_1 = \frac{\beta_1}{1-\beta_1}$ 

**Theorem 4.4.** Under Assumptions 3.1-3.3, if  $\eta_l$  satisfies:  $\eta_l \leq \min\left\{\frac{1}{8KL},\right\}$ 

 $\frac{n(m-1)\epsilon}{48m(n-1)K\sqrt{\beta_2K^2G^2+\epsilon \cdot [3\eta L+C_1^2\eta L+2\sqrt{2(1-\beta_2)}G]}}\right\}, then FedAA in Algorithm 1 under partial$ participation ha

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 8\sqrt{\beta_2 \eta_l^2 K^2 G^2 + \epsilon} \left[\frac{f_0 - f_*}{\eta \eta_l K T} + \frac{\Xi}{T} + \Omega\right],$$

where  $\Xi = \frac{C_1 G^2 d}{\sqrt{\epsilon}} + \frac{2G_1^2 \eta \eta K L G^2 d}{\epsilon}, \quad \Omega = \frac{5\eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_q^2) + [(3 + C_l^2)\eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{n\eta\epsilon}\sigma_l^2 + \frac{2\eta_l K^2 G^2}{\eta\epsilon}) + [(3 + C_1^2)\eta L + 2\sqrt{2(1 - \beta_2)}G]\frac{\eta_l(m-n)}{2n(m-1)\epsilon} [15K^2 L^2 \eta_l^2 (\sigma_l^2 + 6K\sigma_g^2) + 3K\sigma_g^2]\frac{1}{\eta\eta_l K} + \frac{\sqrt{2}G}{\epsilon} \text{ and } C_1 = \frac{\beta_1}{1 - \beta_1}.$ 

*Remark* 4.3. When the parameters C = D and  $\frac{C}{\alpha n} = 1$ , result of Theorem 4.3 becomes the one of Theorem 4.4. The upper bound for  $\min_{t \in [T]} \mathbb{E}[\|\nabla f(\theta_t)\|^2]$  contains three terms: The first two terms decrease as T increases, and this term tends to zero as t tends to infinity. The last term relates to the local stochastic variance  $\sigma_l$  and global variance  $\sigma_g$ . In the i.i.d setting, where the global variance is zero and each worker has the same data distribution, i.e.,  $\sigma_g = 0$ , the variance term  $\Omega$  will be smaller. 

**Corollary 4.2.** Suppose choose local learning rate  $\eta_l = \Theta(\frac{1}{\sqrt{T}K})$  and global learning rate  $\eta =$  $\Theta(\sqrt{Kn})$ , the convergence rate for FedNLAA and FedAA in Algorithm 1 under partial participation without replacement sampling is 

$$\min_{f \in [T]} \mathbb{E} \left[ \|\nabla f(\theta_t)\|^2 \right] = \mathcal{O} \left( \frac{\sqrt{K}}{\sqrt{Tn}} \right)$$

Remark 4.4. Note that Corollary 4.2 suggests that Theorems 4.3 and 4.4 directly relate to the global variance  $\sigma_a^2$ . Such convergence rate is consistent with the partial participation result of FedAvg in the non-i.i.d case in Yang et al. (2021a). It is shown that the global variance has more influence on the convergence behavior in partial participation cases. This is especially true for highly non-i.i.d cases where  $\sigma_q$  is large. The effect of the number of local updates, K, is complex. In partial participation settings, the larger value of K results in a slower convergence, while full participation suggests the opposite. A similar slowdown was also seen in Wang et al. (2022a).

378	Algo	rithm 2 FedBNLACA and FedBACA.
379	Inpu	<b>it:</b> initial value $\theta_1, \theta_0 = 0$ , local step size $\eta_l$ , global step size $\eta$ , constants $\beta_1, \beta_2$ and $\epsilon$ , for each
300	clien	t $i \in S_t$ , $\Delta_0^i = 0$ , compressor $C(\cdot)$ .
301	1: 1	$\mathbf{m}_0 \leftarrow 0, \mathbf{v}_0 \leftarrow 0, \mathbf{e}_1^i = 0, \mathbf{E}_1^i = 0.$
302	2: 1	for $t = 1$ to $T$ do
383	3:	On the server: Server randomly selects a subset of clients $S_t$
384	4:	for each client $j \in S_t$ in parallel <b>do</b>
385	5:	$\theta_t^i = \theta_t$
386	6:	Compress: $\theta_t^i = \mathcal{C}(\theta_t^i + \mathbf{E}_t^i), \mathcal{C}(Q_t^i) = \theta_t^i - \theta_{t-1}^i,$
387	7:	Judge: If $\mathcal{C}(Q_t^i)$ satisfies NLA (Example 3) or AA (Example 6), output: $\widehat{\theta}_t^i$
389	0	
390	8:	Update: $\mathbf{E}_{t+1}^{\circ} = \theta_t^{\circ} + \mathbf{E}_t^{\circ} - \theta_t$ ,
391	9: 10:	end for each client $i \notin S$ in parallel do
392	10.	maintain stale compression error $\mathbf{F}^{j} = \mathbf{F}^{j}$
393	11:	and for
394	12.	
395	13:	On the clients: $\theta_{t,0}^i = \theta_t$
396	14:	for each client $i \in S_t$ in parallel do
397	15:	<b>for</b> for $k = 0,, K - 1$ <b>do</b> Compute level SCD: $\pi^i = \nabla F(0^i + t^i)$
398	10:	Compute local SOD: $\mathbf{g}_{t,k}^{i} = \nabla F_{i}(\sigma_{t,k}^{i}; \boldsymbol{\xi}_{t,k}^{i}),$
399	17:	$ heta_{t,k+1}^{*}= heta_{t,k}^{*}-\eta_{l}\mathbf{g}_{t,k}^{*}.$
400	18:	end for $\gamma^i$
401	19:	$\Delta_t^i =  heta_{t,K}^i - \widehat{ heta}_t$
402	20:	Following the same way as in Algorithm 3 (Line 11-13)
403	21:	end for
404	22:	Following the same way as in Algorithm 3 (Line 14-18)
405	23:	end for

In the following, we will give Algorithms 2 and 3. Due to space constraints, we will only give Algorithm 2 here, and Algorithm 3 is moved into the appendix A. Algorithm 2, which is the most important one in this paper, is an effective adaptive federated learning algorithm for non-convex heterogeneity with bidirectional communication.

# 4.3 FEDBNLACA AND FEDBACA ALGORITHMS

Algorithm 3 gives only a one-way algorithm to reduce the cost of communication (unplink). In the section, we propose two bidirectional communication algorithms (uplink and downlink) with efficiently adaptive non-convex optimization: Federated Bidirectional New Lazy Aggregation Compression
AMSGrad (FedBNLACA) and Federated Bidirectional Accelerated Compression AMSGrad (FedBACA). The detailed procedure is given in Algorithm 2.

Next, we show the convergence analysis for FedBNLACA and FedBACA. Due to the space limit, we only show the full participation setting and leave the partial participation setting in Appendix E.2.
We show the full participation setting and leave the partial participation setting in Appendix E.2.

**Theorem 4.5.** Under Assumptions 3.1-3.3, if the local learning rate  $\eta_l$  satisfies:  $\eta_l \leq \min\left\{\frac{1}{8KL},\right\}$ 

$$\frac{\epsilon}{KC_{\beta,q}[3\eta L+2C_2\eta L+2\sqrt{2(1-\beta_2)}G]} \bigg\}, \text{ where } C_{\beta,q} = \sqrt{4\beta_2(1+q^2)^3(1-q^2)^{-2}K^2G^2+\epsilon}, \text{ then}$$

FedBNLACA in Algorithm 2 under partial participation has

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 4\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 + \epsilon} \Big[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\Big],$$

429  
430 where 
$$\Xi = \frac{C_1 G^2 d}{\sqrt{\epsilon}} + \frac{2C_1^2 \eta \eta K L G^2 d}{\epsilon}$$
,  $\Omega = \left[G + \frac{L \eta \eta_l K G}{\sqrt{\epsilon}} + \frac{L \eta \eta_l C_1 K G d}{\epsilon}\right] \frac{\eta (\gamma + \frac{C}{\alpha m})H}{(1-\beta)\sqrt{\epsilon}} + \frac{431}{(1-\beta)^2 \epsilon} + \frac{2L \eta^2 (\gamma^2 + \frac{C^2}{\alpha^2 m^2})H^2}{(1-\beta)\sqrt{\epsilon}} \frac{5\eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_g^2) + \left[(3 + 2C_2)\eta L + \frac{C}{\alpha^2 m^2}\right] + \frac{2L \eta^2 (\gamma^2 + \frac{C^2}{\alpha^2 m^2})H^2}{(1-\beta)\sqrt{\epsilon}} + \frac{C}{\alpha^2 m^2} \frac{\delta \eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_g^2) + \left[(3 + 2C_2)\eta L + \frac{C}{\alpha^2 m^2}\right] + \frac{C}{\alpha^2 m^2} \frac{\delta \eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_g^2) + \frac{C}{\alpha^2 m^2} \frac{\delta \eta^2 K L^2}{(1-\beta)^2 \epsilon} + \frac{C}{\alpha^2 m^2} \frac{\delta \eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_g^2) + \frac{C}{\alpha^2 m^2} \frac{\delta \eta^2 K L^2}{\sqrt{2\epsilon}} \frac{\delta \eta^2 K L^2}{\sqrt{2\epsilon}} + \frac{C}{\alpha^2 m^2} \frac{\delta \eta^2 K L^2}{\sqrt{2\epsilon}} + \frac{C}{\alpha^2 m^2} \frac{\delta \eta^2 K L^2}{\sqrt{2\epsilon}} + \frac{C}{\alpha^2 m^2} \frac{\delta \eta^2 K L^2}{\sqrt{2\epsilon}} \frac{\delta \eta^2 K L^2}{\sqrt{2\epsilon}} + \frac{C}{\alpha^2 m^2} \frac{\delta \eta^2 K L^2$ 

**Theorem 4.6.** Under Assumptions 3.1-3.3, if the local learning rate  $\eta_l$  satisfies:  $\eta_l \leq \min\left\{\frac{1}{8KL}\right\}$ 

 $\frac{\epsilon}{KC_{\beta,q}[3\eta L+2C_2\eta L+2\sqrt{2(1-\beta_2)G}]} \bigg\}, \text{ where } C_{\beta,q} = \sqrt{4\beta_2(1+q^2)^3(1-q^2)^{-2}K^2G^2 + \epsilon}, \text{ then } FedBACA \text{ in Algorithm 2 under partial participation has}$ 

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 4\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 + \epsilon \Big[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\Big]},$$

$$\begin{array}{lll} \mbox{where } \Xi &=& \frac{C_1 G^2 d}{\sqrt{\epsilon}} + \frac{2C_1^2 \eta \eta K L G^2 d}{\epsilon}, \ \Omega &=& \left[G + \frac{L \eta \eta_l K G}{\sqrt{\epsilon}} + \frac{L \eta \eta_l C_1 K G d}{\epsilon}\right] \cdot \frac{\eta (\gamma + \frac{C}{\alpha m}) H}{(1 - \beta) \sqrt{\epsilon}} + \\ & \frac{2L \eta^2 (\gamma^2 + \frac{C^2}{\alpha^2 m^2}) H^2}{(1 - \beta)^2 \epsilon} &+& \frac{2L \eta^2 (\gamma^2 + \frac{C^2}{\alpha^2 m^2}) H^2}{(1 - \beta) \sqrt{\epsilon}} \frac{5 \eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6 K \sigma_g^2) + \left[(3 + 2C_2) \eta L + 2\sqrt{2(1 - \beta_2)} G\right] \frac{\eta_l}{2m \eta \epsilon} \sigma_l^2, \ C_1 &=& \frac{\beta_1}{1 - \beta_1} + \sqrt{\frac{12q^2}{(1 - q^2)^2} + \frac{(1 - q^2)^2 C^2}{\alpha^2 m^2 q^2}} + \frac{2T L^2 \eta^2 (\gamma^2 + \frac{C^2}{\alpha^2 m^2}) H^2}{(1 - \beta)^2 \epsilon} (\frac{C^2}{\alpha^2 m^2} + 1) \\ & and \ C_2 &=& \frac{\beta_1^2}{(1 - \beta_1)^2} + \frac{4(q + \gamma + \frac{\lambda C}{\alpha m})^2}{(1 - q^2)^2}. \end{array} \end{array}$$

*Remark* 4.5. When the parameters C = D and  $\frac{C}{\alpha m} = 1$ , result of Theorem 4.5 becomes the one of Theorem 4.6. The upper bound for  $\min_{t \in [T]} \mathbb{E}[||\nabla f(\theta_t)||^2]$  contains three terms: The first two terms decrease as T increases, and this term tends to zero as t tends to infinity. The last term relates to the local stochastic variance  $\sigma_l$  and global variance  $\sigma_q$ . In the i.i.d setting, where the global variance is zero and each worker has the same data distribution, i.e.,  $\sigma_g = 0$ , the variance term  $\Omega$  will be smaller. 

**Corollary 4.3.** Suppose choose local learning rate  $\eta_l = \Theta(\frac{1}{\sqrt{T}K})$  and global learning rate  $\eta =$  $\Theta(\sqrt{Km})$ , when T is sufficiently large, i.e.,  $T = \mathcal{O}(Km)$ , the convergence rate for FedNLAA and FedAA in Algorithm I under full participation has

$$\min_{t \in [T]} \mathbb{E} \left[ \|\nabla f(\theta_t)\|^2 \right] = \mathcal{O} \left( \frac{1}{T} \right)$$

Remark 4.6. Corollary 4.3 suggests that with sufficient large T, FedBNLACA and FedBACA achieve a convergence rate of  $\mathcal{O}(\frac{1}{T})$ , which matches the result for general federated non-convex optimization methods such as FedAMS Wang et al. (2022a) and FedAdam Reddi et al. (2020).

#### EXPERIMENTS

We compare our proposed algorithms with several state-of-the-art algorithms (FedAvgMcMahan et al. (2017), FedAMS, FedCAMS(Wang et al. (2022a)), Fedadam (Reddi et al. (2020)). We use MNIST and Fashion-MNIST datasets, and models by MLP and CNN, respectively. A total of 100 clients for all federated training experiments are used. Set the partial participation rate to 0.5, i.e. in each round, the server selects 50 clients out of 100 to participate in communication and model updating. In each round, the client completes 3 local epochs with batch size 32. In experiments, we respectively sample Independent Identical Distribution (I.I.D.) and non-I.I.D. client data from the dataset. Choose compression rate in Topk to be 1/8 and 1/128. For parameters C and D, the values are not too large. From theoretical analysis, the larger the values of C and D, the larger the errors. In addition, C and D are too small, and basically does not help to reduce communication cost. We also verify this result. We suggest that C and D in the vicinity of  $\frac{\sqrt{\alpha \log m}}{m}$ 

Figures 1-2 represent the relationship between the accuracy of prediction and communication Bits when the model is CNN and the I.I.D. client data is sampled from Fashion-MNIST dataset. From Figure 1, we can find that (i) the proposed algorithm FedNLAA not only communicates fewer Bits than the other three state-of-the-art algorithms, but also has higher accuracy; (ii) our proposed algorithms (FedNLACA and FedBNLACA) require only few communication Bits to achieve good accuracy, especially the bidirectional compression algorithm FedBNLACA requires even fewer Bits. These shows that our proposed algorithms are communication efficient. From Figure 2, we can find that (i) the proposed algorithm FedAA algorithm can converge more quickly than the other three



Figure 1: NLA strategy, on Fashion-MNIST via Figure 2: AA strategy, on Fashion-MNIST via CNN model, and the I.I.D. client sampling. CNN model, and the I.I.D. client sampling.



Figure 3: NLA strategy, on Fashion-MNIST via Figure 4: AA strategy, on Fashion-MNIST via CNN model, and the non-I.I.D. client sampling. CNN model, and the non-I.I.D. client sampling.

algorithms, thus disguising the reduction of communication cost; (ii) FedACA and FedBACA can
 also converge quickly and with higher accuracy than the other three algorithms, especially FedBACA.

Figures 3-4 represent the relationship between the accuracy of prediction and communication Bits when the model is CNN and the non-I.I.D. client data is sampled from Fashion-MNIST dataset..From Figure 3, we can find that (i) the proposed algorithm FedNLAA not only communicates fewer Bits than the other three state-of-the-art algorithms, but also has higher accuracy; (ii) FedAvg performs the worst; (iii) our proposed algorithms (FedNLACA and FedBNLACA) require only few communication Bits to achieve good accuracy, especially the bidirectional compression algorithm FedBNLACA requires even fewer Bits. These also show that our proposed algorithms are communication efficient. From Figure 4, we can find that (i) the proposed FedAA algorithm can converge more quickly than the other three algorithms, thus disguising the reduction of communication cost; (ii) our algorithms (FedACA and FedBACA) can also converge quickly and with higher accuracy than the other three algorithms, especially FedBACA. 

6 CONCLUSION

We propose two novel strategies: NLA and AA in the framework of federated learning. They are simple to operate and effective in reducing the communication cost. The NLA strategy achieves communication cost reduction by reducing the amount of information passed and AA strategy reduces the communication cost by accelerating computation. By combining our proposed strategies with compression techniques, we design FedNLAA and FedAA algorithms, which not only achieve communication cost reduction in one-way, but also extend them to bidirectional algorithms (FedBNLACA and FedBACA), which achieve communication cost reduction in bidirectional as well.

# 540 REFERENCES

549

550

551 552

553

554

555

567

568

569

570

571

572 573

574

575

576

580

581

582

542	Dan Alistarh, Torsten Hoefler, Mikael Johansson, Sarit Khirirat, Nikola Konstantinov, and Cédric
543	Renggli. The convergence of sparsified gradient methods. arXiv preprint, arXiv:1809.10505, 2018.
544	

- Yoshua Bengio and Yann LeCun. Scaling learning algorithms towards AI. In *Large Scale Kernel Machines*. MIT Press, 2007.
- Aleksandr Beznosikov, Samuel Horvath, Peter Richtárik, and M. H. Safaryan. On biased compression
   for distributed learning. *Journal of Machine Learning Research*, 24(276):1–50, 2023.
  - T. B. Brown, B. Mann, N. Ryder, and etal. Language models are few-shot learners. *arXiv preprint*, arXiv:2005.14165, 2020.
  - Jinghui Chen, Yuan Cao, and Quanquan Gu. Closing the generalization gap of adaptive gradient methods in training deep neural networks. In *International Joint Conference on Artificial Intelligence*, 2018a.
- M. Chen, N. Shlezinger, H. V. Poor, and etal. Communication-efficient federated learning. *Proceed*ings of the National Academy of Sciences, 118, 2021.
- Tianyi Chen, Georgios B. Giannakis, Tao Sun, and Wotao Yin. Lag: Lazily aggregated gradient for communication-efficient distributed learning. In *Neural Information Processing Systems*, 2018b.
- Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong. On the convergence of a class of adam-type
   algorithms for non-convex optimization. *International Conference on Learning Representations*,
   2018c.
- Xiangyi Chen, Xiaoyun Li, and P. Li. Toward communication efficient adaptive gradient method.
   *Proceedings of the 2020 ACM-IMS on Foundations of Data Science Conference*, 2020.
  - J. Devlin, M.W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*, 2019.
  - J. C. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12:2121–2159, 2011.
  - A. Elgabli, C. Ben Issaid, A. S. Bedi, and etal. Fednew: A communication-efficient and privacypreserving newton-type method for federated learning. In *International Conference on Machine Learning*, volume 162, pp. 5861–5877, 2022.
- Hossein Shokri Ghadikolaei, Sebastian U. Stich, and Martin Jaggi. Lena: Communication-efficient
   distributed learning with self-triggered gradient uploads. In *International Conference on Artificial Intelligence and Statistics*, 2021.
  - Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- Farzin Haddadpour, Mohammad Mahdi Kamani, Aryan Mokhtari, and Mehrdad Mahdavi. Fed erated learning with compression: Unified analysis and sharp guarantees. *arXiv preprint*,
   arXiv:2007.01154, 2020.
- Geoffrey E. Hinton, Simon Osindero, and Yee Whye Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18:1527–1554, 2006.
- M. Huang, D. Zhang, and K. Ji. Achieving linear speedup in non-iid federated bilevel learning.
   *Proceedings of the 40th International Conference on Machine Learning*, 2023.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Sebastian U. Stich, and
   Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In International Conference on Machine Learning, 2019a.

594 Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian U. Stich, and Martin Jaggi. Error feedback 595 fixes signsgd and other gradient compression schemes. Proceedings of the 36th International 596 Conference on Machine Learning, 97:3252–3261, 2019b. 597 D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. International Conference on 598 Learning Representations, 2014. 600 D. Li and J. Wang. Fedmd: Heterogenous federated learning via model distillation. arXiv preprint, 601 arXiv:1910.03581, 2019. 602 T. Li, A. K. Sahu, M. Zaheer, and etal. Feddane: A federated newton-type method. Asilomar 603 Conference on Signals, Systems and Computers, 2019a. 604 605 X. Li, K. Huang, W.Yang, and etal. On the convergence of fedavg on non-iid data. arXiv preprint, 606 arXiv:1907.02189, 2019b. 607 Yuzhu Mao, Zihao Zhao, Guangfeng Yan, Yang Liu, Tian Lan, Linqi Song, and Wenbo Ding. 608 Communication-efficient federated learning with adaptive quantization. ACM Transactions on 609 Intelligent Systems and Technology (TIST), 13:1–26, 2022. 610 611 H. B. McMahan, E. Moore, D. Ramage, and etal. Communication-efficient learning of deep networks 612 from decentralized data. In International Conference on Artificial Intelligence and Statistics, pp. 613 1273-1282, 2017. 614 Konstantin Mishchenko, Grigory Malinovsky, Sebastian U. Stich, and Peter Richt'arik. Proxskip: 615 Yes! local gradient steps provably lead to communication acceleration! finally! In International 616 Conference on Machine Learning, 2022. 617 618 S. J. Reddi, S. Kale, and S. Kumar. On the convergence of adam and beyond. International Conference 619 on Learning Representations, 2018. 620 S.J. Reddi, Zachary B. Charles, M. Zaheer, Z. Garrett, and etal. Adaptive federated optimization. 621 arXiv preprint, arXiv:2003.00295, 2020. 622 623 Amirhossein Reisizadeh, Aryan Mokhtari, Hamed Hassani, Ali Jadbabaie, and Ramtin Pedarsani. 624 Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization. In International Conference on Artificial Intelligence and Statistics, pp. 2021–2031, 625 2020. 626 627 Richtárik, Igor Sokolov, and Ilyas Fatkhullin. Ef21: A new, simpler, theoretically better, and. 628 Advances in Neural Information Processing Systems, 34:4384–4396, 2023. 629 A. K. Sahu, T. Li, M. Sanjabi, and etal. Federated optimization in heterogeneous networks. *Conference* 630 on Machine Learning and Systems, 2018. 631 632 Shaohuai Shi, Xiaowen Chu, Ka Chun Cheung, and S. See. Understanding top-k sparsification in 633 distributed deep learning. arXiv preprint, arXiv:1911.08772, 2019. 634 Sebastian U. Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified sgd with memory. Neural 635 Information Processing Systems, 2018. 636 637 J. Sun, T. Chen, G. B. Giannakis, and etal. Lazily aggregated quantized gradient innovation for 638 communication-efficient federated learning. IEEE Transactions on Pattern Analysis and Machine 639 Intelligence, 44:2031-2044, 2020. 640 Jun Sun, Tianyi Chen, Georgios B. Giannakis, and Zaiyue Yang. Communication-efficient distributed 641 learning via lazily aggregated quantized gradients. In Neural Information Processing Systems, 642 2019. 643 644 Qianqian Tong, Guannan Liang, and Jinbo Bi. Effective federated adaptive gradient methods with 645 non-iid decentralized data. arXiv preprint, arXiv:2009.06557, 2020. 646 Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H. Vincent Poor. Tackling the objective 647 inconsistency problem in heterogeneous federated optimization. ArXiv, abs/2007.07481, 2020.

648 649 650	Y.J. Wang, L.Li, and J.H Chen. Communication-efficient adaptive federated learning. <i>Proceedings of the 39th International Conference on Machine Learning(ICML 2022)</i> , 2022a.
651 652 653	Yujia Wang, Lu Lin, and Jinghui Chen. Communication-compressed adaptive gradient method for distributed nonconvex optimization. <i>In International Conference on Artificial Intelligence and Statistics</i> , pp. 6292–6320, 2022b.
654 655	X.D. Wu, F.H. Huang, Z.M. Hu, and H. Huang. Faster adaptive federated learning. Association for the Advancement of Artificial Intelligence(AAAI 2023), 2023.
657 658	Haibo Yang, Minghong Fang, and Jia Liu. Achieving linear speedup with partial worker participation in non-iid federated learning. <i>International Conference on Learning Representations</i> , 2021a.
659 660	Haibo Yang, Minghong Fang, and Jia Liu. Achieving linear speedup with partial worker participation in non-iid federated learning. <i>ICLR</i> , 2021b.
662 663 664	Zhaohui Yang, Mingzhe Chen, Walid Saad, Choong Seon Hong, and Mohammad R. Shikh-Bahaei. Energy efficient federated learning over wireless communication networks. <i>IEEE Transactions on Wireless Communications</i> , 20:1935–1949, 2019.
665 666 667	M. D. Zeiler. Adadelta: An adaptive learning rate method. <i>arXiv preprint</i> , arXiv:1212.5701, 2012.
668 669	
671 672	
673 674	
675 676	
678 679	
680 681	
682 683	
685 686	
687 688	
689 690 691	
692 693	
694 695 696	
697 698	
699 700 701	

# Appendix

# A FEDNLACA AND FEDACA ALGORITHMS

## A.1 PROCEDURE OF FEDNLACAAND FEDACA

In order to reduce communication costs, we propose **Fed**erated New Lazy Aggregation Compression AMSGrad (FedNLACA) and **Fed**erated Accelerated Compression AMSGrad (FedACA). These two algorithms combine two of our proposed strategies(strategy 1 (FedNLACA), strategy 2 (FedACA)) and compression techniques. The detailed procedure is given in Algorithm 3.

# A.2 CONVERGENCE ANALYSIS FOR FEDNLACAAND FEDACA

715 In the case of full participation: 

 **Theorem A.1.** Under Assumption 3.1-3.3, if the local learning rate  $\eta_l$  satisfies:  $\eta \leq \min\left\{\frac{1}{8KL}, \frac{\epsilon}{KC_{\beta,q}[3\eta L+2C_2\eta L+2\sqrt{2(1-\beta_2)}G]}\right\}$ , where  $C_{\beta,q} = \sqrt{4\beta_2(1+q^2)^3(1-q^2)^{-2}K^2G^2 + \epsilon}$ , then the iterates of FedBNLACA in Algorithm 2 under partial

 $\sqrt{4\beta_2(1+q^2)^3(1-q^2)^{-2}K^2G^2} + \epsilon$ , then the iterates of FedBNLACA in Algorithm 2 under partial participation scheme satisfy

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 4\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 + \epsilon \Big[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\Big]},$$

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 4\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 + \epsilon} \Big[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\Big],$$

where  $\Xi = \frac{C_1 G^2 d}{\sqrt{\epsilon}} + \frac{2C_1^2 \eta \eta K L G^2 d}{\epsilon}, \Omega = \frac{5\eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K \sigma_g^2) + [(3 + 2C_2)\eta L + 2\sqrt{2(1-\beta_2)}G] \frac{\eta_l}{2m\eta\epsilon} \sigma_l^2, C_1 = \frac{\beta_1}{1-\beta_1} + \sqrt{\frac{12q^2}{(1-q^2)^2} + \frac{(1-q^2)^2}{q^2}} and C_2 = \frac{\beta_1^2}{(1-\beta_1)^2} + \frac{4(q+\gamma+\lambda)^2}{(1-q^2)^2}.$ 

**Remark** A.1. When the parameters C = D,  $\frac{C}{\alpha m} = 1$ , the result of Theorem A.1 becomes the result of Theorem A.2. The upper bound for  $\min_{t \in [T]} \mathbb{E}[||\nabla f(\theta_t)||^2]$  contains three terms: The first two terms decrease as T increases, and this term tends to zero as t tends to infinity. The last term relates to the local stochastic variance  $\sigma_l$  and global variance  $\sigma_g$ . In the i.i.d setting, where the global variance is zero and each worker has the same data distribution, i.e.,  $\sigma_g = 0$ , the variance term  $\Omega$  will be smaller.

**Corollary A.1.** Suppose we choose local learning rate  $\eta_l = \Theta(\frac{1}{\sqrt{TK}})$  and the global learning rate  $\eta = \Theta(\sqrt{Km})$ , when T is sufficient large, i.e.,  $T \ge Km$ , the convergence rate for FedNLACA and FedACA in Algorithm I under full participation scheme satisfies

$$\min_{t \in [T]} \mathbb{E} \left[ \|\nabla f(\theta_t)\|^2 \right] = \mathcal{O} \left( \frac{1}{\sqrt{TKm}} \right)$$

751 *Remark* A.2. Corollary A.1 suggests that with sufficient large *T*, FedCAMS achieves the desired 752  $\mathcal{O}(\frac{1}{\sqrt{TKm}})$  convergence rate which matches the result for its uncompressed counterpart FedNLAA 753 and FedAA. In addition, when  $T = \mathcal{O}(Km)$ ,  $\min_{t \in [T]} \mathbb{E}[||\nabla f(\theta_t)||^2] = \mathcal{O}(\frac{1}{T})$ . This suggests that 754 FedNLACA and FedACA can indeed achieve better communication efficiency without sacrificing 755 much on the accuracy.

756 Algorithm 3 FedNLACA and FedACA 757 **Input:** initial value  $\theta_1$ , local step size  $\eta_l$ , global step size  $\eta$ ,constant  $\beta_1,\beta_2,\epsilon$ , for each client  $i \in S_t$ , 758  $\Delta_0^i = 0$ , compressor  $C(\cdot)$ 759 1:  $\mathbf{m}_0 \leftarrow 0, \mathbf{v}_0 \leftarrow 0, \mathbf{e}_1^i = 0$ 760 2: for t = 1 to T do 761 Randomly select a subset of clients  $S_t$  and the server transmits  $\theta_t$  to the subset of clients  $S_t$ 3: 762 4:  $\theta_{t,0}^i = \theta_t$ 763 5: for each client  $i \in S_t$  in parallel do for for k = 0, ..., K - 1 do 764 6: Compute local stochastic gradient:  $\mathbf{g}_{t,k}^i = \nabla F_i(\theta_{t,k}^i; \xi_{t,k}^i)$ 7: 765  $\theta^i_{t,k+1} = \theta^i_{t,k} - \eta_l \mathbf{g}^i_{t,k}$ 766 8: 767 9: end for 768 10:  $\Delta_t^i = \theta_{t,K}^i - \theta_t$ Compress  $\widehat{\Delta}_t^i = \mathcal{C}(\Delta_t^i + \mathbf{e}_t^i), \mathcal{C}(q_t^i) = \widehat{\Delta}_t^i - \widehat{\Delta}_{t-1}^i$ , 769 11: 770 Judge: If  $C(q_t^i)$  satisfies NLA (Example 2) or AA (Example 5), 12: 771 then outputs the result  $\widehat{\hat{\Delta}}_t^i$  of the judgement and passes it to the server and update  $\mathbf{e}_{t+1}^i =$ 13: 772  $\Delta_t^i + \mathbf{e}_t^i - \widehat{\widehat{\Delta}}_t^i$ 773 774 14: end for 775 15: for each client  $j \notin S_t$  in parallel do **do** 776 client j maintains the stale compression error  $\mathbf{e}_{t+1}^{j} = \mathbf{e}_{t}^{j}$ 16: 777 17: end for 778 Server aggregates local update:  $\widehat{\Delta}_t = \frac{1}{|S_t|} \sum_{i \in S_t} \widehat{\widehat{\Delta}}_t^i$ 18: 779 Server updates  $\mathbf{x}_{t+1}$  using  $\hat{\Delta}_t$  in the same way as in Algorithm 1 (Line 14-16) 19: 780 20: end for 781

In the case of partial participation

782 783

784

785

786

787

788

798

799 800

801

**Theorem A.3.** Under Assumption 3.1-3.3, if the local learning rate  $\eta_l$  satisfies:  $\eta_l \leq \min\left\{\frac{1}{8KL}, \frac{\epsilon}{KC_{\beta,q}[3\eta L+2C_2\eta L+2\sqrt{2(1-\beta_2)}G]}\right\}$ , where  $C_{\beta,q} = \sqrt{4\beta_2(1+q^2)^3(1-q^2)^{-2}K^2G^2 + \epsilon}$ , then the iterates of FedBNLACA in Algorithm 2 under partial participation scheme satisfy

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 8\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 + \epsilon} \left[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\right]$$

 $\begin{array}{l} \text{, where } \Xi &= \frac{C_1 G^3 d}{\sqrt{\epsilon}} + \frac{2C_1^2 \eta \eta_l K L G^2 d}{\epsilon}, \Omega &= \frac{C_1 \eta \eta_l K L G^2}{\epsilon} + \frac{5\eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K \sigma_g^2) + [\eta L + \sqrt{2(1 - \beta_2)}G] \frac{\eta_l}{\eta n \epsilon} \sigma_l^2 + [\eta L + \sqrt{2(1 - \beta_2)}G] \frac{\eta_l (m - n)}{n(m - 1)\epsilon} [15K^2 L^2 \eta_l^2 (\sigma_l^2 + 6K \sigma_g^2) + 3K \sigma_g^2] \text{ and } C_1 &= \frac{\beta_1}{1 - \beta_1} + \frac{m}{n} \sqrt{\frac{12q^2}{(1 - q^2)^2} + \frac{(1 - q^2)^2 C^2}{\alpha^2 n^2 q^2}}. \end{array}$ 

**Theorem A.4.** Under Assumption 3.1-3.3, if the local learning rate  $\eta_l$  satisfies:  $\eta_l \leq \min\left\{\frac{1}{8KL}, \frac{\epsilon}{KC_{\beta,q}[3\eta L + 2C_2\eta L + 2\sqrt{2(1-\beta_2)}G]}\right\}$ , where  $C_{\beta,q} = \sqrt{4\beta_2(1+q^2)^3(1-q^2)^{-2}K^2G^2 + \epsilon}$ , then the iterates of FedBACA in Algorithm 2 under partial participation scheme satisfy

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 8\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 + \epsilon} \left[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\right]$$

 $\begin{array}{ll} \text{806} \\ \text{807} \\ \text{808} \\ \text{808} \\ \text{808} \\ \text{809} \end{array} & \text{where } \Xi = \frac{C_1 G^3 d}{\sqrt{\epsilon}} + \frac{2C_1^2 \eta \eta_l K L G^2 d}{\epsilon}, \Omega = \frac{C_1 \eta \eta_l K L G^2}{\epsilon} + \frac{5\eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K \sigma_g^2) + [\eta L + \sqrt{2(1 - \beta_2)}G] \frac{\eta_l (m - n)}{n(m - 1)\epsilon} [15K^2 L^2 \eta_l^2 (\sigma_l^2 + 6K \sigma_g^2) + 3K \sigma_g^2] \text{ and } \\ \text{809} \\ C_1 = \frac{\beta_1}{1 - \beta_1} + \frac{m}{n} \sqrt{\frac{12q^2}{(1 - q^2)^2} + \frac{(1 - q^2)^2}{q^2}}. \end{array}$ 

*Remark* A.3. When the parameters C = D,  $\frac{C}{\alpha n} = 1$ , the result of Theorem A.3 becomes the result of Theorem A.4. The upper bound for  $\min_{t \in [T]} \mathbb{E}[||\nabla f(\theta_t)||^2]$  contains three terms: The first two terms decrease as T increases, and this term tends to zero as t tends to infinity. The last term relates to the local stochastic variance  $\sigma_l$  and global variance  $\sigma_q$ . In the i.i.d setting, where the global variance is zero and each worker has the same data distribution, i.e.,  $\sigma_q = 0$ , the variance term  $\Omega$  will be smaller.

**Corollary A.2.** Suppose we choose local learning rate  $\eta_l = \Theta(\frac{1}{\sqrt{T}K})$  and the global learning rate  $\eta = \Theta(\sqrt{Kn})$ , the convergence rate for FedNLACA, FedACA in Algorithm 1 under partial participation scheme without replacement sampling is 

$$\min_{t \in [T]} \mathbb{E} \left[ \|\nabla f(\theta_t)\|^2 \right] = \mathcal{O} \left( \frac{\sqrt{K}}{\sqrt{Tn}} \right)$$

Remark A.4. Note that Corollary A.2 suggests that Theorem A.3, A.4 directly relates to the global variance  $\sigma_a^2$ . Such convergence rate is consistent with the partial participation result of FedAvg in the non i.i.d case in Yang et al. (2021a). It is shown that the global variance has more influence on the convergence behaviour in partial participation cases. This is especially true for highly non i.i.d. cases where  $\sigma_q$  is large. The effect of the number of local updates, K, is complex. In partial participation settings, larger values of K result in slower convergence, while full participation suggests the opposite. A similar slowdown was also seen in Wang et al. (2022a).

#### В ALL EXPERIMENTS



Figure 5: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is MLP, the data is MNIST dataset, the dataset obeys independent identical distribution and the data is partially used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 6: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is MLP, the data is MNIST dataset, the dataset obeys independent identical distribution and the all data is used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 7: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is MLP, the data is MNIST dataset, the dataset does not obeys independent identical distribution and the data is partially used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 8: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is MLP, the data is MNIST dataset, the dataset does not obeys independent identical distribution and all the data is used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 9: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is MLP, the data is Fashion-MNIST dataset, the dataset obeys independent identical distribution and the data is partially used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 10: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is MLP, the data is Fashion-MNIST dataset, the dataset obeys independent identical distribution and the all data is used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 11: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is MLP, the data is Fashion-MNIST dataset, the dataset does not obeys independent identical distribution and the data is partially used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 12: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is MLP, the data is Fashion-MNIST dataset, the dataset does not obeys independent identical distribution and all the data is used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 13: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is CNN, the data is MNIST dataset, the dataset obeys independent identical distribution and the data is partially used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 14: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is CNN, the data is MNIST dataset, the dataset obeys independent identical distribution and the all data is used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 15: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is CNN, the data is MNIST dataset, the dataset does not obeys independent identical distribution and the data is partially used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.

- 1132
- 1133



Figure 16: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is CNN, the data is MNIST dataset, the dataset does not obeys independent identical distribution and all the data is used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 17: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is CNN, the data is Fashion-MNIST dataset, the dataset obeys independent identical distribution and the data is partially used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 18: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is CNN, the data is Fashion-MNIST dataset, the dataset obeys independent identical distribution and the all data is used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.



Figure 19: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is CNN, the data is Fashion-MNIST dataset, the dataset does not obeys independent identical distribution and the data is partially used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.

- 1240
- 1241



Figure 20: The above figure represents the relationship between the accuracy of the prediction and the communication Bits when the model used is CNN, the data is Fashion-MNIST dataset, the dataset does not obeys independent identical distribution and all the data is used. The left side represents the method using NLA (New Lazy Aggregation) strategy and the right side represents the method using AA (Accelerated Aggregation) strategy.

1264 C PROOF IN SECTION 4.1

 C.1 PROOF OF THEOREM 4.1

Similar to previous work in the field of adaptive methods Chen et al. (2018c); Wang et al. (2022a), we introduce a Lyapunov sequence  $z_t$ : assume  $\theta_0 = \theta_1$ , for each  $t \ge 1$ , there is the following equation

$$\mathbf{z}_{t} = \theta_{t} + \frac{\beta_{1}}{1 - \beta_{1}} (\theta_{t} - \theta_{t-1}) = \frac{1}{1 - \beta_{1}} \theta_{t} - \frac{\beta_{1}}{1 - \beta_{1}} \theta_{t-1}.$$
 (C.1)

1280 For the difference of sequence z, there is the following equation 1281

$$\begin{aligned} \mathbf{z}_{t+1} - \mathbf{z}_t &= \frac{1}{1 - \beta_1} (\theta_{t+1} - \theta_t) - \frac{\beta_1}{1 - \beta_1} (\theta_t - \theta_{t-1}) \\ &= \frac{1}{1 - \beta_1} (\eta \widehat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t) - \frac{\beta_1}{1 - \beta_1} \eta \widehat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}_{t-1} \\ &= \frac{1}{1 - \beta_1} \eta \widehat{\mathbf{V}}_t^{-1/2} \Big[ \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \widetilde{\Delta}_t \Big] - \frac{\beta_1}{1 - \beta_1} \eta \widehat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}_{t-1} \\ &= \eta \widehat{\mathbf{V}}_t^{-1/2} \Delta_t - \eta \frac{\beta_1}{1 - \beta_1} \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_t^{-1/2} \right) \mathbf{m}_{t-1}. \end{aligned}$$

Since f is L-smooth, taking conditional expectation at time t, we get

$$\mathbb{E}[f(\mathbf{z}_{t+1})] - f(\mathbf{z}_{t}) \\
\leq \mathbb{E}[\langle \nabla f(\mathbf{z}_{t}), \mathbf{z}_{t+1} - \mathbf{z}_{t} \rangle] + \frac{L}{2} \mathbb{E}[\|\mathbf{z}_{t+1} - \mathbf{z}_{t}\|^{2}] \\
\leq \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \widetilde{\Delta}_{t} \right\rangle \right] - \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \eta \frac{\beta_{1}}{1 - \beta_{1}} \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) \mathbf{m}_{t-1} \right\rangle \right] \\
+ \frac{\eta^{2} L}{2} \mathbb{E}\left[\left\| \widehat{\mathbf{V}}_{t}^{-1/2} \widetilde{\Delta}_{t} - \frac{\beta_{1}}{1 - \beta_{1}} \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) \mathbf{m}_{t-1} \right\|^{2} \right] \\
= \underbrace{\mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \widetilde{\Delta}_{t} \right\rangle}_{I_{1}} - \eta \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\beta_{1}}{1 - \beta_{1}} \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) \mathbf{m}_{t-1} \right\rangle \right]}_{I_{2}} \\
+ \underbrace{\frac{\eta^{2} L}{2} \mathbb{E}\left[\left\| \widehat{\mathbf{V}}_{t}^{-1/2} \widetilde{\Delta}_{t} - \frac{\beta_{1}}{1 - \beta_{1}} \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) \mathbf{m}_{t-1} \right\|^{2} \right]}_{I_{3}} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \widetilde{\Delta}_{t} \right\rangle \right]}_{I_{4}}_{I_{4}} \tag{C.2}$$

Recall the notation  $\widehat{\mathbf{V}}_t = \operatorname{diag}(\widehat{\mathbf{v}}_t) = \operatorname{diag}(\max(\widehat{\mathbf{v}}_{t-1}, \mathbf{v}_t, \epsilon)).$ **Bounding**  $I_1$ : 

$$\begin{split} I_{1} &= \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \eta \frac{\tilde{\Delta}_{t}}{\sqrt{\tilde{\mathbf{v}}_{t}}} \right\rangle\right] \\ &\leq \eta \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \frac{\sqrt{2} \cdot \tilde{\Delta}_{t}}{\sqrt{\mathbf{v}_{t} + \epsilon}} \right\rangle\right] \\ &= \sqrt{2}\eta \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \frac{\tilde{\Delta}_{t}}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \right\rangle\right] + \sqrt{2}\eta \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \frac{\tilde{\Delta}_{t}}{\sqrt{\mathbf{v}_{t} + \epsilon}} - \frac{\tilde{\Delta}_{t}}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \right\rangle\right], \end{split}$$
(C.3)

where the first inequality follows by the fact that  $\hat{\mathbf{v}}_t \geq \frac{\mathbf{v}_t + \epsilon}{2}$ . For the second term in C.2, then

 $\sqrt{2}\eta \mathbb{E}\left[\left\langle \nabla f(\theta_t), \frac{\tilde{\Delta}_t}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{\tilde{\Delta}_t}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\rangle\right]$ 

 $\leq \frac{\eta \sqrt{2(1-\beta_2)}G}{\epsilon} \mathbb{E}[\|\tilde{\Delta}_t\|^2],$ 

where the second inequality follows from Lemma F.1 and F.5, and we will further apply the bound for  $\mathbb{E}[\|\tilde{\Delta}_t\|^2]$  following

 $\leq \sqrt{2}\eta \mathbb{E} \|\nabla f(\theta_t)\| \mathbb{E} \left[ \left\| \frac{1}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{1}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\| \cdot \|\tilde{\Delta}_t\| \right]$ 

(C.4)

$$\begin{aligned}
&\sqrt{2}\eta \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \frac{\tilde{\Delta}_{t}}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}\right\rangle\right] \\
&= \sqrt{2}\eta \mathbb{E}\left[\left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}, \tilde{\Delta}_{t} + \eta K \nabla f(\theta_{t}) - \eta_{l} K \nabla f(\theta_{t})\right\rangle\right] \\
&= \sqrt{2}\eta \mathbb{E}\left[\left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}, \tilde{\Delta}_{t} + \eta K \nabla f(\theta_{t}) - \eta_{l} K \nabla f(\theta_{t})\right\rangle\right] \\
&= -\sqrt{2}\eta \eta_{t} K \mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}\right\|^{2}\right] + \sqrt{2}\eta \mathbb{E}\left[\left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}, \tilde{\Delta}_{t} + \eta_{l} K \nabla f(\theta_{t})\right\rangle\right] \\
&= -\sqrt{2}\eta \eta_{h} K \mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}\right\|^{2}\right] + \sqrt{2}\eta \left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}, \mathbb{E}\left[-\frac{1}{m}\sum_{i=1}^{m}\sum_{k=0}^{K-1} \eta \mathbf{g}_{t,k}^{i} + \eta_{l} K \nabla f(\theta_{t}) - \frac{1}{M_{t}}\sum_{i\in M_{t}} q_{t}^{i}\right]\right\rangle \\
&= -\sqrt{2}\eta \eta_{n} K \mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}\right\|^{2}\right] + \sqrt{2}\eta \left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}, \mathbb{E}\left[-\frac{\eta_{l}}{m}\sum_{i=1}^{m}\sum_{k=0}^{K-1} \mathbf{g}_{t,k}^{i} + \frac{\eta_{l} K}{m}\sum_{i=1}^{m} \nabla F_{i}(\theta_{t}) - \frac{1}{M_{t}}\sum_{i\in M_{t}} q_{t}^{i}\right]\right\rangle, \\
&= -\sqrt{2}\eta \eta_{n} K \mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}\right\|^{2}\right] + \sqrt{2}\eta \left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}, \mathbb{E}\left[-\frac{\eta_{l}}{m}\sum_{i=1}^{m}\sum_{k=0}^{K-1} \mathbf{g}_{t,k}^{i} + \frac{\eta_{l} K}{m}\sum_{i=1}^{m} \nabla F_{i}(\theta_{t}) - \frac{1}{M_{t}}\sum_{i\in M_{t}} q_{t}^{i}\right]\right\rangle, \\
&= -\sqrt{2}\eta \eta_{n} K \mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}\right\|^{2}\right] + \sqrt{2}\eta \left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}, \mathbb{E}\left[-\frac{\eta_{l}}{m}\sum_{i=1}^{m}\sum_{k=0}^{K-1} \mathbf{g}_{t,k}^{i} + \frac{\eta_{l} K}{m}\sum_{i=1}^{m} \nabla F_{i}(\theta_{t}) - \frac{1}{M_{t}}\sum_{i\in M_{t}} q_{t}^{i}\right]\right\rangle, \\ &= -\sqrt{2}\eta \eta_{n} K \mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}\right\|^{2}\right] + \sqrt{2}\eta \left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}, \mathbb{E}\left[-\frac{\eta_{l}}{m}\sum_{i=1}^{m}\sum_{k=0}^{K-1} \mathbf{g}_{t,k}^{i} + \frac{\eta_{l} K}{m}\sum_{i=1}^{m} \nabla F_{i}(\theta_{t}) - \frac{\eta_{l}}{M_{t}}\sum_{i\in M_{t}} q_{t}^{i}\right\right]\right\rangle,$$

where the third equality follows the local update rule. For the last term in C.5, we get

$$\begin{split} & \sqrt{2}\eta \left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}, \mathbb{E}\left[-\frac{\eta}{m}\sum_{i=1}^{m}\sum_{k=0}^{K-1}\mathbf{g}_{t,k}^{i} + \frac{\eta K}{m}\sum_{i=1}^{m}\nabla F_{i}(\theta_{t})\right]\right\rangle + \sqrt{2}\eta \mathbb{E}\left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}, -\frac{1}{M_{c}}\sum_{i\in M_{c}}q_{t}^{i}\right\rangle \\ & = \sqrt{2}\eta \left\langle \frac{\sqrt{\eta K}}{\sqrt{\beta_{2}\mathbf{v}_{t-1}+\epsilon}} \nabla f(\theta_{t}), -\frac{\sqrt{\eta K}}{Km}\frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}} \mathbb{E}\left[\sum_{i=1}^{m}\sum_{k=0}^{K-1}(\nabla F_{i}(\theta_{t,k}^{i}) - \nabla F_{i}(\theta_{t}))\right]\right\rangle \\ & + \sqrt{2}\eta \mathbb{E}\left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}, -\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{t}^{i}\right\rangle \\ & = \frac{\sqrt{2}\eta\eta K}{2}\mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\|^{2}\right] + \frac{\sqrt{2}\eta\eta t}{2Km^{2}}\mathbb{E}\left[\left\|\frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}} \sum_{i=1}^{m}\sum_{k=0}^{K-1}(\nabla F_{i}(\theta_{t,k}^{i}) - \nabla F_{i}(\theta_{t}))\right\|^{2}\right] \\ & = \frac{\sqrt{2}\eta\eta t}{2}\mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\|^{2}\right] + \frac{\sqrt{2}\eta\eta t}{2Km^{2}}\mathbb{E}\left[\left\|\frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}} \sum_{i=1}^{m}\sum_{k=0}^{K-1}(\nabla F_{i}(\theta_{t,k}^{i}) - \nabla F_{i}(\theta_{t}))\right\|^{2}\right] \\ & - \frac{\sqrt{2}\eta\eta t}{2Km^{2}}\mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\|^{2}\right] + \frac{\sqrt{2}\eta\eta t}{2m}\sum_{i=1}^{m}\sum_{k=0}^{K-1}\mathbb{E}\left[\left\|\frac{\nabla F_{i}(\theta_{t,k}^{i}) - \nabla F_{i}(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\|^{2}\right] \\ & - \frac{\sqrt{2}\eta\eta t}{2Km^{2}}\mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\|^{2}\right] + \frac{\sqrt{2}\eta\eta t}{2m}\sum_{i=1}^{m}\sum_{k=0}^{K-1}\mathbb{E}\left[\left\|\frac{\nabla F_{i}(\theta_{t,k}^{i}) - \nabla F_{i}(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\|^{2}\right] \\ & - \frac{\sqrt{2}\eta\eta t}{2Km^{2}}\mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\|^{2}\right] + \frac{\sqrt{2}\eta\eta t}{2m}\sum_{i=1}^{m}\sum_{k=0}^{K-1}\mathbb{E}\left[\left\|\frac{\nabla F_{i}(\theta_{t,k}^{i}) - \nabla F_{i}(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\|^{2}\right] \\ & - \frac{\sqrt{2}\eta\eta t}{\sqrt{2}}\mathbb{E}\left[\left\|\frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\|^{2}\right] + \frac{\sqrt{2}\eta\eta t}{2m}\sum_{i=1}^{m}\sum_{k=0}^{K-1}\mathbb{E}\left[\left\|\frac{\nabla F_{i}(\theta_{t,k}^{i}) - \nabla F_{i}(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\right]^{2}\right] \\ & - \frac{\sqrt{2}\eta\eta t}{2}\mathbb{E}\left[\left\|\frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\|^{2}\right] + \frac{\sqrt{2}\eta\eta t}{2m}\sum_{i=1}^{K-1}\mathbb{E}\left[\left\|\frac{\nabla F_{i}(\theta_{t,k}^{i}) - \nabla F_{i}(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\right\|^{2}\right] \\ & - \frac{\sqrt{2}\eta\eta t}{2}\mathbb{E}\left[\left\|\frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right\|^{2}\right] + \frac{\sqrt{2}\eta\eta t}{2m}\sum_{i=1}^{K-1}\mathbb{E}\left[\left(\frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}\right)\right] \\ & - \frac$$

where the second equation follows from  $\langle x, y \rangle = \frac{1}{2} [\|x\|^2 + \|y\|^2 - \|x - y\|^2]$ , the first inequality holds by applying Cauchy-Schwarz inequality, the second inequality follows from  $\langle x, y \rangle \leq \frac{1}{2} ||x||^2 + ||y||^2$ , the third inequality follows from Assumption 3.1.

 $\sqrt{2} \cdot \eta \left\langle \frac{\nabla f(\theta_t)}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}}, \mathbb{E} \left[ -\frac{\eta_l}{m} \sum_{i=1}^m \sum_{k=0}^{K-1} \mathbf{g}_{t,k}^i + \frac{\eta_l K}{m} \sum_{i=1}^m \nabla F_i(\theta_t) \right] \right\rangle$  $\leq \frac{3\sqrt{2}\eta\eta_l K}{4} \mathbb{E}\left[ \left\| \frac{\nabla f(\theta_t)}{\sqrt[4]{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\|^2 \right] + \frac{5\eta\eta_l^3 K^2 L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_g^2)$  $-\frac{\sqrt{2}\eta\eta_l}{2Km^2}\mathbb{E}\left[\left\|\frac{1}{\sqrt[4]{\beta_2\mathbf{v}_{t-1}+\epsilon}}(\sum_{i=1}^m\sum_{k=0}^{K-1}\nabla F_i(\boldsymbol{\theta}_{t,k}^i)\right\|^2\right]+\sqrt{2}\eta\mathbb{E}\left\langle\frac{\nabla f(\boldsymbol{\theta}_t)}{\sqrt{\beta_2\mathbf{v}_{t-1}+\epsilon}},-\frac{1}{M_t}\sum_{i\in M_t}q_t^i\right\rangle.$ (C.7)

Hence by applying Lemma F.14 with the local learning rate condition:  $\eta_l \leq \frac{1}{8KL}$ , then

Then merging pieces together,

$$\begin{split} I_{1} &\leq -\frac{\sqrt{2}\eta\eta K}{4} \mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}\right\|^{2}\right] + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}}{\sqrt{2\epsilon}}(\sigma_{l}^{2} + 6K\sigma_{g}^{2}) \\ &- \frac{\sqrt{2}\eta\eta_{l}}{2Km^{2}} \mathbb{E}\left[\left\|\frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}\right] + \frac{\eta\sqrt{2(1-\beta_{2})}G}{\epsilon} \mathbb{E}[\|\tilde{\Delta}_{t}\|^{2}] \\ &\leq -\frac{\eta\eta_{l}K - 2\sqrt{2}\eta}{4} \mathbb{E}\left[\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}\right\|^{2}\right] + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}}{\sqrt{2\epsilon}}(\sigma_{l}^{2} + 6K\sigma_{g}^{2}) \\ &- \frac{\sqrt{2}\eta\eta_{l}}{2Km^{2}} \mathbb{E}\left[\left\|\frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}\right] + \sqrt{2}\eta \mathbb{E}\left\langle\frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}}, -\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{t}^{i}\right\rangle. \end{split}$$
(C.8)

**Bounding**  $I_2$ : The bound for  $I_2$  mainly follows by the update rule and definition of virtual sequence  $\mathbf{z}_t$ ,

where the last inequality holds by applying Lemma F.5 and the fact of  $\hat{\mathbf{v}}_{t-1} \geq \epsilon$ . 

**Bounding**  $I_3$ : It can be bounded as follows:

$$\leq \eta^{2} L \mathbb{E} \left[ \left\| \widehat{\mathbf{V}}_{t}^{-1/2} \widetilde{\Delta}_{t} \right\|^{2} \right] + \eta^{2} L \mathbb{E} \left[ \left\| \frac{\beta_{1}}{1 - \beta_{1}} \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) \mathbf{m}_{t-1} \right\|^{2} \right]$$

$$\leq \eta^{2} L \mathbb{E} \left[ \left\| \widehat{\mathbf{V}}_{t}^{-1/2} \widetilde{\Delta}_{t} \right\|^{2} \right] + \eta^{2} L \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \eta^{2} K^{2} G^{2} \mathbb{E} \left[ \left\| \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right\|^{2} \right], \quad (C.10)$$

where the first inequality follows by Cauchy-Schwarz inequality, and the second one follows by Lemma F.5.

**1471 Bounding**  $I_4$ :

$$I_{4} = \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}) - \nabla f(\mathbf{x}_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\rangle \right]$$
$$\leq \mathbb{E}\left[ \|\nabla f(\mathbf{z}_{t}) - \nabla f(\mathbf{x}_{t})\| \|\eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t}\| \right]$$

 $I_3 = \frac{\eta^2 L}{2} \mathbb{E} \left[ \left\| \widehat{\mathbf{V}}_t^{-1/2} \widetilde{\Delta}_t + \frac{\beta_1}{1 - 2} \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_t^{-1/2} \right) \mathbf{m}_{t-1} \right\|^2 \right]$ 

1477  
1478 
$$\leq L\mathbb{E}\left[\|\mathbf{z}_t - \mathbf{x}_t\| \|\eta \widehat{\mathbf{V}}_t^{-1/2} \Delta_t\|\right]$$
1479

$$\leq \frac{\eta^2 L}{2} \mathbb{E} \bigg[ \left\| \widehat{\mathbf{V}}_t^{-1/2} \Delta_t \right\|^2 \bigg] + \frac{\eta^2 L}{2} \mathbb{E} \bigg[ \left\| \frac{\beta_1}{1 - \beta_1} \widehat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}_{t-1} \right\|^2 \bigg],$$

where the first inequality holds by the fact of  $\langle \mathbf{a}, \mathbf{b} \rangle \leq \|\mathbf{a}\| \|\mathbf{b}\|$ , the second one follows from Assumption 3.1 and the third one holds by the definition of virtual sequence  $\mathbf{z}_t$  and the fact of  $\|\mathbf{a}\| \|\mathbf{b}\| \leq \frac{1}{2} \|\mathbf{a}\|^2 + \frac{1}{2} \|\mathbf{b}\|^2$ . Then summing  $I_4$  over  $t = 1, \dots, T$ , then

 $\sum_{t=1}^{T} I_4 \leq \frac{\eta^2 L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\tilde{\Delta}_t\|^2] + \frac{\eta^2 L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E}\left[\left\|\frac{\beta_1}{1-\beta_1}\mathbf{m}_t\right\|^2\right]$ 

 $\leq \frac{\eta^2 L}{2\epsilon} \sum_{t=1}^T \mathbb{E}[\|\tilde{\Delta}_t\|^2] + \frac{\eta^2 L}{2\epsilon} \frac{\beta_1^2}{(1-\beta_1)^2} \sum_{t=1}^T \mathbb{E}[\|\mathbf{m}_t\|^2].$ 

(C.11)

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathbf{m}_{t}\|^{2}] \leq \frac{TK\eta_{l}^{2}}{m}\sigma_{l}^{2} + \frac{2\eta_{l}^{2}}{m^{2}}\sum_{t=1}^{T} \mathbb{E}\Big[\Big\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\Big] + \frac{2}{m^{2}}\sum_{t=1}^{T} \mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{t}^{i}\right\|^{2}$$

The summation of  $I_4$  term is bounded by

By Lemma F.10,

$$\sum_{t=1}^{T} I_4 \leq \frac{\eta^2 L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\tilde{\Delta}_t\|^2] + \frac{2\beta_1^2}{(1-\beta_1)^2} \frac{\eta^2 L}{2\epsilon} \frac{\eta_l^2}{m^2} \sum_{t=1}^{T} \mathbb{E} \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i) \right\|^2 + \frac{\beta_1^2}{(1-\beta_1)^2} \frac{\eta^2 L}{2\epsilon} \frac{TK\eta_l^2}{m} \sigma_l^2 + \frac{2\beta_1^2}{(1-\beta_1)^2} \frac{\eta^2 L}{2\epsilon} \frac{1}{m^2} \sum_{t=1}^{T} \mathbb{E} \left\| \frac{1}{M_t} \sum_{i \in M_t} q_t^i \right\|^2.$$
(C.12)

Merging pieces together: Substituting (C.8), (C.9) and (C.10) into (C.2), summing over from t = 1 to T and then adding (C.11), then

$$\begin{split} \mathbb{E}[f(\mathbf{z}_{T+1})] - f(\mathbf{z}_{1}) &= \sum_{t=1}^{T} [I_{1} + I_{2} + I_{3} + I_{4}] \\ &\leq -\frac{\eta\eta_{l}K - 2\sqrt{2}\eta}{4} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{\nabla f(\theta_{l})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \right\|^{2} \right] + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}T}{\sqrt{2\epsilon}} (\sigma_{l}^{2} + 6K\sigma_{g}^{2}) + \frac{\sqrt{2(1 - \beta_{2})}\eta G}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\tilde{\Delta}_{l}\|^{2}] \\ &- \frac{\eta\eta_{l}}{2Km^{2}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right) \right\|^{2} \right] \\ &+ \frac{\beta_{1}}{1 - \beta_{1}} \eta\eta_{l}KG^{2} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \widehat{\mathbf{v}}_{t-1}^{-1/2} - \widehat{\mathbf{v}}_{t}^{-1/2} \right\|_{1} \right] + \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \frac{\eta^{2}\eta_{l}^{2}K^{2}G^{2}}{\sqrt{\epsilon}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \widehat{\mathbf{v}}_{t-1}^{-1/2} - \widehat{\mathbf{v}}_{t}^{-1/2} \right\|_{1} \right] \\ &+ \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \eta^{2}\eta_{n}^{2}K^{2}LG^{2} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \widehat{\mathbf{v}}_{t-1}^{-1/2} - \widehat{\mathbf{v}}_{t}^{-1/2} \right\|^{2} \right] + \eta^{2}L \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \widehat{\mathbf{v}}_{t-1}^{-1/2} \widetilde{\Delta}_{t} \right\|^{2} \right] \\ &+ \frac{\eta^{2}L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\widetilde{\Delta}_{t}\|^{2}] + \frac{\eta^{2}L}{2\epsilon} \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \sum_{t=1}^{T} \mathbb{E}[\|\mathbf{m}_{t}\|^{2}]. \end{split}$$
(C.13)

By applying Lemma F.6 is into all terms containing the second moment estimate of model difference  $\tilde{\Delta}_t$  in (C.13),and the fact that  $(\sqrt{\beta_2 K^2 G^2 + \epsilon})^{-1} \|\theta\| \le (\sqrt{\beta_2 \eta_l^2 K^2 G^2 + \epsilon})^{-1} \|\theta\| \le \|\frac{\theta}{\sqrt{\beta_2 \mathbf{v}_t + \epsilon}}\| \le \epsilon^{-1/2} \|\theta\|$ , then

$$\begin{split} & \mathbb{E}[f(\mathbf{z}_{T+1})] - f(\mathbf{z}_{1}) \\ & \mathbb{E}[f(\mathbf{z}_{T+1})] - f(\mathbf{z}_{1}) \\ & \leq -\frac{\eta\eta_{t}K - 2\sqrt{2}\eta}{4} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \right\|^{2} \right] + \frac{5\eta\eta_{t}^{3}K^{2}L^{2}T}{\sqrt{2\epsilon}} (\sigma_{t}^{2} + 6K\sigma_{g}^{2}) + \frac{\beta_{1}}{1 - \beta_{1}} \frac{\eta\eta KG^{2}d}{\sqrt{\epsilon}} \\ & + \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \frac{2\eta^{2}\eta_{t}^{2}K^{2}LG^{2}d}{\epsilon} - \frac{2\eta\eta_{t}}{2Km^{2}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \right] \\ & + \left( \eta^{2}L + \frac{\eta^{2}L}{2} + \sqrt{2(1 - \beta_{2})}\eta G \right) \left[ \frac{KT\eta_{t}^{2}}{m\epsilon} \sigma_{t}^{2} + \frac{2\eta_{t}^{2}}{m^{2}\epsilon} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \right] \\ & + \frac{2\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \frac{\eta^{2}L}{2\epsilon} \frac{\eta_{t}^{2}}{m^{2}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{i,k}^{i}) \right\|^{2} \right] + \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \frac{\eta^{2}L}{2\epsilon} \frac{TK\eta_{t}^{2}}{m} \sigma_{t}^{2} \\ & + \frac{2\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \frac{\eta^{2}L}{2\epsilon} \frac{\eta_{t}^{2}}{m^{2}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{i,k}^{i}) \right\|^{2} \right] \\ & + \left( \eta^{2}L + \frac{\eta^{2}L}{2} + \sqrt{2(1 - \beta_{2})} + \eta G \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \frac{\eta^{2}L}{2} \right) \frac{2}{m^{2}\epsilon} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{1}{M_{t}} \sum_{i\in M_{t}} q_{t}^{i} \right\|^{2} \right] \\ & + \sqrt{2}\eta \sum_{t=1}^{T} \mathbb{E}\left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} - \frac{1}{M_{t}} \sum_{i\in M_{t}} q_{t}^{i} \right\rangle \\ & \leq -\frac{\eta\eta K - 2\sqrt{2}\eta}{\sqrt{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[ \| \nabla f(\theta_{t}) \|^{2} ] + \frac{5\eta\eta_{t}^{3}K^{2}L^{2}T}{\sqrt{2\epsilon}} (\sigma_{t}^{2} + 6K\sigma_{q}^{2}) \\ & + \frac{\beta_{1}}{1 - \beta_{1}} \frac{\eta K G^{2}d}{\sqrt{\epsilon}} + \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \frac{2\eta^{2}\eta_{t}^{2}K^{2}LG^{2}d}{\epsilon} \\ \end{cases}$$

$$\begin{aligned} & + \left( \eta^{2}L + \frac{\eta^{2}L}{2} + \sqrt{2(1-\beta_{2})}\eta G + \frac{\beta_{1}^{2}}{(1-\beta_{1})^{2}} \frac{\eta^{2}L}{2} \right) \frac{KT\eta_{l}^{2}}{m\epsilon} \sigma_{t}^{2} - \sum_{t=1}^{T} \mathbb{E} \left[ \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \right] \\ & \cdot \left[ \frac{\eta\eta}{2\sqrt{\beta_{2}K^{2}G^{2} + \epsilon Km^{2}}} - \left( \eta^{2}L + \frac{\eta^{2}L}{2} + \eta\sqrt{2(1-\beta_{2})}G + \frac{\beta_{1}^{2}}{(1-\beta_{1})^{2}} \frac{\eta^{2}L}{2} \right) \frac{2\eta_{t}^{2}}{m^{2}\epsilon} \right] \\ & + \left( \eta^{2}L + \frac{\eta^{2}L}{2} + \sqrt{2(1-\beta_{2})}\eta G + \frac{\beta_{1}^{2}}{(1-\beta_{1})^{2}} \frac{\eta^{2}L}{2} \right) \frac{T2K^{2}C^{2}\eta_{t}^{2}C^{2}}{\alpha^{2}m^{2}\epsilon} + \sqrt{2}\eta T\mathbb{E} \left\| \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}}v_{t-1} + \epsilon} \right\| \cdot \mathbb{E} \right\| - \frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i} \right\| \\ & = \frac{\eta\eta_{t}K - 2\sqrt{2}\eta}{4\sqrt{\beta_{2}}\eta_{t}^{2}K^{2}G^{2}} + \epsilon \sum_{t=1}^{T} \mathbb{E} [ \|\nabla f(\theta_{t})\|^{2} ] + \frac{5\eta\eta_{1}^{3}S^{2}L^{2}T}{\sqrt{2\epsilon}} (\sigma_{t}^{2} + 6K\sigma_{g}^{2}) \\ & + \frac{\beta_{1}}{1-\beta_{1}} \frac{\eta\eta_{t}KG^{2}}{\sqrt{\epsilon}} + \frac{\beta_{1}^{2}}{(1-\beta_{1})^{2}} \frac{2\eta^{2}\eta_{t}^{2}K^{2}LG^{2}}{\epsilon} + \frac{\sqrt{2}\eta\eta_{t}K^{2}G^{2}}{\sqrt{2}\eta_{t}^{2}K^{2}G^{2} + \epsilon} \\ & + \left( \eta^{2}L + \frac{\eta^{2}L}{2} + \eta\sqrt{2(1-\beta_{2})}G + \frac{\beta_{1}^{2}}{(1-\beta_{1})^{2}} \frac{\eta^{2}L}{2} \right) \frac{2TK^{2}C^{2}\eta_{t}^{2}G^{2}}{\alpha^{2}m^{2}\epsilon} \\ & + \left( \eta^{2}L + \frac{\eta^{2}L}{2} + \sqrt{2(1-\beta_{2})}G + \frac{\beta_{1}^{2}}{(1-\beta_{1})^{2}} \frac{\eta^{2}L}{2} \right) \frac{2TK^{2}C^{2}\eta_{t}^{2}G^{2}}{\alpha^{2}m^{2}\epsilon} \\ & + \left( \eta^{2}L + \frac{\eta^{2}L}{2} + \sqrt{2(1-\beta_{2})}G + \frac{\beta_{1}^{2}}{(1-\beta_{1})^{2}} \frac{\eta^{2}L}{2} \right) \frac{2TK^{2}C^{2}\eta_{t}^{2}G^{2}}{\alpha^{2}m^{2}\epsilon} \\ & \text{The last inequality holds due to additional constraint of local learning rate  $\eta_{l}$  with the inequality  $\frac{\eta\eta_{k}}{2\sqrt{\beta_{g}K^{2}G^{2} + \epsilon \cdot T}} \frac{\eta\eta_{k}K}{\sqrt{\beta_{g}K^{2}G^{2} + \epsilon(\beta^{2}-\beta_{1}^{2}-\beta_{1}^{2}-\beta_{1}^{2})} \frac{\eta^{2}L}{m^{2}} \right) \frac{\eta\eta_{k}}{m^{2}\epsilon} \geq 0, \text{ thus obtain the constraint } \\ \frac{\eta\eta_{k}K}{\sqrt{\beta_{g}K^{2}G^{2} + \epsilon \cdot T}} \sum_{k}^{T} \mathbb{E} [\|\nabla f(\theta_{k})\|^{2}] \end{aligned}$$$

$$\frac{\eta_{lk}}{4\sqrt{\beta_2\eta_l^2 K^2 G^2 + \epsilon} \cdot T} \sum_{t=1}^{\infty} \mathbb{E}[\|\nabla f(\theta_t)\|^2] \\
\frac{4\sqrt{\beta_2\eta_l^2 K^2 G^2 + \epsilon} \cdot T}{T} \sum_{t=1}^{\infty} \mathbb{E}[\|\nabla f(\theta_t)\|^2] \\
\frac{596}{1597} \leq \frac{f(\mathbf{z}_0) - \mathbb{E}[f(\mathbf{z}_T)]}{T} + \frac{5\eta\eta_l^3 K^2 L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_g^2) + [\frac{\sqrt{2}\eta}{2\epsilon} + (3 + C_1^2)\eta^2 L + 2\sqrt{2(1 - \beta_2)}\eta G] (\frac{2K\eta_l^2}{m\epsilon} \sigma_l^2 + \frac{2K^2 C^2 \eta_l^2 G^2}{\alpha^2 m^2 \epsilon}) \\
\frac{1598}{1599} + \frac{C_1\eta\eta_l K G^2 d}{T\sqrt{\epsilon}} + \frac{2C_1^2\eta^2 \eta_l^2 K^2 L G^2 d}{T\epsilon} + \frac{\sqrt{2}\eta\eta_l K G C}{\alpha m \epsilon}.$$

Therefore

1601

1606 1607 1608

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 4\sqrt{\beta_2 \eta_l^2 K^2 G^2 + \epsilon} \cdot \left[\frac{f_0 - f_*}{(\eta \eta_l K - 2\sqrt{2}\eta)T} + \frac{\Xi}{T} + \Omega\right],$$

where

$$\Xi = \frac{C_1 G^2 d}{\sqrt{\epsilon}} + \frac{2C_1^2 \eta \eta_l K L G^2 d}{\epsilon}$$

, and

$$\begin{array}{ll} \text{1609} & \text{, and} \\ \text{1610} & \Omega = \frac{5\eta_l^2 K^2 L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_g^2) + (3 + C_1^2)\eta^2 L + 2\sqrt{2(1 - \beta_2)}\eta G] (\frac{2\eta_l}{m\eta\epsilon}\sigma_l^2 + \frac{2KC^2\eta_l G^2}{\alpha^2\eta m^2\epsilon}) + \frac{\sqrt{2}GC}{\alpha m\epsilon} \\ \text{1612} & \text{, that } C_1 = \frac{\beta_1}{1 - \beta_1}. \end{array}$$

1614 The proof of Theorem 4.2 is similar to the above proof procedure and the detailed proof will not be given here.1616

# 1617 C.2 PROOF OF COROLLARY 4.1

1618  
1619 If pick 
$$\eta = \Theta(\frac{1}{\sqrt{TK}}, \eta = \Theta(\sqrt{Km}) \text{ and } T = \mathcal{O}(Km) \text{ ,then } \min_{t \in [T]} \mathbb{E}[\|\nabla f(\theta_t)\|^2] = \mathcal{O}(\frac{1}{T}).$$

**Assumption C.1.** (Compression Dissimilarity). For the biased compressor, there exists a constant  $\xi$ such that, for each iteration  $t \ge 0$ , that 

$$\left\| \mathcal{C}\left(\frac{1}{m}\sum_{i=1}^{m} [\Delta_{t}^{i} + \mathbf{e}_{t}^{i}]\right) - \frac{1}{m}\sum_{i=1}^{m} \mathcal{C}(\Delta_{t}^{i} + \mathbf{e}_{t}^{i}) \right\|$$
$$\leq \gamma \left\| \frac{1}{m}\sum_{i=1}^{m} \Delta_{t}^{i} \right\|$$
(C.14)

and

$$\left\|\frac{1}{M_t}\sum_{i\in M_t} \mathcal{C}(q_t^i)\right\| \le \lambda \left\|\frac{1}{m}\sum_{i=1}^m \Delta_t^i\right\|.$$
(C.15)

Here  $M_t$  denotes the set of all clients satisfying (3) or (4) at round t, and  $\mathcal{C}(q_t^i)$  is given in Algorithms 2 and 3. The assumption of bounded gradient is usually adopted in adaptive gradient methods Alistarh et al. (2018). 



Figure 21: Empirical justification for Assumption C.1. on CNN and MLP based on the Fashion-MINIST datset.

#### **PROOF OF THEOREM 4.3** C.3

Notations and equations: For partial participation, i.e.  $|S_t| = n, \forall t \in [T]$ . The global model difference is the average of local model difference from the subset  $S_t$ , i.e.,  $\tilde{\Delta}_t = \frac{1}{n} \sum_{i \in S_t} \tilde{\Delta}_i^t$ . Denote  $\bar{\Delta}_t = \frac{1}{m} \sum_{i=1}^m \tilde{\Delta}_t^i$ , and for convenience, we follow the previous notation of  $\widehat{\mathbf{V}}_t = \operatorname{diag}(\widehat{\mathbf{v}}_t + \epsilon)$ . Next we show that the global model difference  $\tilde{\Delta}_t$  is an unbiased estimator of  $\tilde{\Delta}_t$ ;

$$\mathbb{E}_{\mathcal{S}_t}[\tilde{\Delta}_t] = \frac{1}{n} \mathbb{E}_{\mathcal{S}_t}[\sum_{i=1}^n \tilde{\Delta}_t^{w_i}] = \mathbb{E}_{\mathcal{S}_t}[\tilde{\Delta}_t^{w_1}] = \frac{1}{m} \sum_{i=1}^m \tilde{\Delta}_t^i = \bar{\Delta}_t.$$

Define the virtual sequence  $\mathbf{z}_t$  same as previous: assume  $\theta_0 = \theta_1$ , for each  $t \ge 1$ , then 

B1 

$$\mathbf{z}_t = \theta_t + \frac{\beta_1}{1 - \beta_1} (\theta_t - \theta_{t-1}) = \frac{1}{1 - \beta_1} \theta_t - \frac{\beta_1}{1 - \beta_1} \theta_{t-1},$$

 $( \land 1/9 )$ 

B1

1671  
1672  
1673  

$$\mathbf{z}_{t+1} - \mathbf{z}_t = \eta \widehat{\mathbf{V}}_t^{-1/2} \widetilde{\Delta}_t - \eta \frac{\beta_1}{1 - \beta_1} \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_t^{-1/2} \right) \mathbf{m}_{t-1}$$
1673

1 /0 ~

By Assumption 3.1, we get

$$\begin{split} & \mathbb{E}[f(\mathbf{z}_{t+1})] - f(\mathbf{z}_{t}) \\ & \leq \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \widetilde{\Delta}_{t} \right\rangle\right] - \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \eta \frac{\beta_{1}}{1 - \beta_{1}} \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \mathbf{m}_{t-1} \right\rangle\right] \\ & + \frac{\eta^{2} L}{2} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t}^{-1/2} \widetilde{\Delta}_{t} - \frac{\beta_{1}}{1 - \beta_{1}} \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \mathbf{m}_{t-1}\right\|^{2}\right] \\ & = \underbrace{\mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \widetilde{\Delta}_{t} \right\rangle\right]}_{I_{1}'} - \eta \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\beta_{1}}{1 - \beta_{1}} \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \mathbf{m}_{t-1}\right\rangle\right]}_{I_{2}'} \\ & + \underbrace{\frac{\eta^{2} L}{2} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t}^{-1/2} \widetilde{\Delta}_{t} - \frac{\beta_{1}}{1 - \beta_{1}} \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \mathbf{m}_{t-1}\right\|^{2}\right]}_{I_{3}'} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \widetilde{\Delta}_{t}\right\rangle\right]}_{I_{4}'} \end{split}$$

Since  $\tilde{\Delta}_t$  is an unbiased estimator of  $\bar{\Delta}_t$ , the main difference of convergence analysis for partial participation cases is bounding  $\mathbb{E}[\|\tilde{\Delta}_t\|^2]$ .

1692 1693 The bound for  $I'_2$  is exactly the same as the bound for  $I_2$ . For the corresponding three terms,  $I'_1, I'_3$ 1694 and  $I'_4$  which include the second-order momentum estimate of  $\tilde{\Delta}_t$ . For  $I'_1$ , then

$$I_{1}^{\prime} = \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \eta \frac{\tilde{\Delta}_{t}}{\sqrt{\mathbf{\hat{v}}_{t}}} \right\rangle\right]$$

$$\leq \eta \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \frac{\sqrt{2} \cdot \tilde{\Delta}_{t}}{\sqrt{\mathbf{v}_{t} + \epsilon}} \right\rangle\right]$$

$$= \eta \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \frac{\sqrt{2} \cdot \tilde{\Delta}_{t}}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \right\rangle\right] + \eta \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \frac{\sqrt{2} \cdot \tilde{\Delta}_{t}}{\sqrt{\mathbf{v}_{t} + \epsilon}} - \frac{\sqrt{2} \cdot \tilde{\Delta}_{t}}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \right\rangle\right]. \quad (C.16)$$

The first term in (C.16) does not change in partial participation scheme. The second term is changed due to the variance of  $\tilde{\Delta}_t$  changes. For the second term of  $I'_1$ , 

$$\sqrt{2}\eta \mathbb{E}\left[\left\langle \nabla f(\theta_t), \frac{\tilde{\Delta}_t}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{\tilde{\Delta}_t}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}}\right\rangle\right] \le \frac{\sqrt{2(1 - \beta_2)}\eta G}{\epsilon} \mathbb{E}[\|\tilde{\Delta}_t\|^2].$$
(C.17)

For  $I'_3$ ,

$$\sum_{t=1}^{T} I_{3}^{\prime} \leq \frac{\eta^{2} L}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\tilde{\Delta}_{t}\|^{2}] + \eta^{2} L \frac{\beta_{1}^{2}}{(1-\beta_{1})^{2}} \eta_{l}^{2} G^{2} \sum_{t=1}^{T} \mathbb{E}\Big[\left\|\left(\hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2}\right)\right\|^{2}\Big], \quad (C.18)$$

and for  $I'_4$  similar to (C.11), We get

$$\sum_{t=1}^{T} I_{4}^{\prime} \leq \frac{\eta^{2} L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\tilde{\Delta}_{t}\|^{2}] + \frac{\eta^{2} L}{2\epsilon} \frac{\beta_{1}^{2}}{(1-\beta_{1})^{2}} \sum_{t=1}^{T} \mathbb{E}[\|\mathbf{m}_{t}\|^{2}].$$
(C.19)

From Lemma F.10, then

$$\sum_{t=1}^{1725} \mathbb{E}[\|\mathbf{m}_{t}\|^{2}] \leq \frac{KT\eta_{l}^{2}}{n}\sigma_{l}^{2} + \frac{2\eta_{l}^{2}}{n^{2}}\sum_{t=1}^{T}\mathbb{E}\left[\left\|\sum_{i\in\mathcal{S}_{t}}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}\right] + \frac{2}{n^{2}}\sum_{t=1}^{T}\mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in\mathcal{M}_{t}}q_{t}^{i}\right\|^{2}.$$

$$(C.20)$$

# Then substituting (C.20) into (C.19), then

$$\sum_{t=1}^{T} I_{4}^{\prime} \leq \frac{\eta^{2}L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\tilde{\Delta}_{t}\|^{2}] + \frac{2\beta_{1}^{2}}{(1-\beta_{1})^{2}} \frac{\eta^{2}\eta_{l}^{2}L}{2n^{2}\epsilon} \sum_{t=1}^{T} \mathbb{E}\left[\left\|\sum_{i\in\mathcal{S}_{t}}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}\right] \\ + \frac{\beta_{1}^{2}}{(1-\beta_{1})^{2}} \frac{\eta^{2}\eta_{l}^{2}KTL}{2n\epsilon} \sigma_{l}^{2} + \frac{1}{n^{2}} \mathbb{E}\left\|\frac{2}{M_{t}}\sum_{i\in\mathcal{M}_{t}}q_{t}^{i}\right\|^{2} \\ \leq \frac{\eta^{2}L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\tilde{\Delta}_{t}\|^{2}] + \frac{2\beta_{1}^{2}}{(1-\beta_{1})^{2}} \frac{\eta^{2}\eta_{l}^{2}L}{2n^{2}\epsilon} \sum_{t=1}^{T} \mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\mathbb{P}\{i\in\mathcal{S}_{t}\}\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}\right] \\ + \frac{\beta_{1}^{2}}{2\epsilon} \frac{\eta^{2}\eta_{l}^{2}KTL}{2\epsilon} \sigma_{t}^{2} + \frac{2}{2} \sum_{i=1}^{T} \mathbb{E}\left[\left\|\frac{1}{2}\sum_{i=1}^{K}\sum_{k=0}^{I}\mathbb{P}\{i\in\mathcal{S}_{t}\}\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}\right]$$

$$+ \frac{\beta_1}{(1-\beta_1)^2} \frac{\eta \eta_l K \Gamma L}{2n\epsilon} \sigma_l^2 + \frac{2}{n^2} \sum_{t=1}^{\infty} \mathbb{E} \left\| \frac{1}{M_t} \sum_{i \in M_t} q_t^i \right\| \quad , \tag{C.21}$$

where we will further apply the bound for  $\mathbb{E}[\|\tilde{\Delta}_t\|^2]$  following by Lemma F.8. The second term in (C.21) can be bounded from (F.11)). Therefore, summing up (C.17), (C.18) and (C.9), summing over from t = 1 to T, then adding (C.19),

$$\begin{aligned} & \mathbb{E}[f(\mathbf{z}_{T+1})] - f(\mathbf{z}_{1}) = \sum_{t=1}^{T} [I_{1}' + I_{2} + I_{3}' + I_{4}'] \\ & \mathbb{E}[f(\mathbf{z}_{T+1})] - f(\mathbf{z}_{1}) = \sum_{t=1}^{T} [I_{1}' + I_{2} + I_{3}' + I_{4}'] \\ & \mathbb{E}[f(\mathbf{z}_{T+1})] - f(\mathbf{z}_{1}) = \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2} \mathbf{v}_{t-1} + \epsilon}} \right\|^{2} \right] + \frac{5\eta \eta_{l}^{3} K^{2} L^{2} T}{\sqrt{2\epsilon}} (\sigma_{l}^{2} + 6K\sigma_{g}^{2}) + \frac{\sqrt{2(1 - \beta_{2})}\eta G}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\tilde{\Delta}_{t}\|^{2}] \\ & - \frac{\eta m}{2Km^{2}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{1}{\sqrt[4]{\beta_{2} \mathbf{v}_{t-1} + \epsilon}} \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \right] \\ & + \frac{\beta_{1}}{1 - \beta_{1}} \eta \pi K G^{2} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \hat{\nabla}_{t-1}^{-1/2} - \hat{\nabla}_{t}^{-1/2} \right\|_{1} \right] + \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \frac{\eta^{2} \eta_{l}^{2} K^{2} G^{2}}{\sqrt{\epsilon}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \hat{\nabla}_{t-1}^{-1/2} - \hat{\nabla}_{t}^{-1/2} \right\|_{1} \right] \\ & + \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \eta^{2} \eta_{l}^{2} K^{2} L G^{2} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \hat{\nabla}_{t-1}^{-1/2} - \hat{\nabla}_{t}^{-1/2} \right\|^{2} \right] \\ & + \frac{\eta^{2} L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \hat{\Delta}_{t} \right\|^{2} \right] + \frac{\eta^{2} L}{2\epsilon} \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \mathbf{m}_{t} \right\|^{2} \right] + \sqrt{2}\eta \mathbb{E}\left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2} \mathbf{v}_{t-1} + \epsilon}} - \frac{1}{M_{t}} \sum_{i \in M_{t}} q_{i}^{i} \right\rangle. \end{aligned} \right. \tag{C.22}$$

By applying Lemma F.6 into all terms containing the second moment estimate of model difference  $\tilde{\Delta}_t$ in (C.22), using the fact that  $(\sqrt{\beta_2 K^2 G^2 + \epsilon})^{-1} \|\theta\| \le (\sqrt{\beta_2 \eta_l^2 K^2 G^2 + \epsilon})^{-1} \|\theta\| \le \|\frac{\theta}{\sqrt{\beta_2 \gamma_2 + \epsilon}}\| \le \epsilon^{-1/2} \|\theta\|$ , and applying Lemma F.10, we get

$$\begin{split} \mathbb{E}[f(\mathbf{z}_{T+1})] &- f(\mathbf{z}_{1}) \\ &\leq -\frac{\eta\eta_{l}K}{4\sqrt{\beta_{2}\eta_{l}^{2}K^{2}G^{2} + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_{t})\|^{2}] + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}T}{\sqrt{2\epsilon}} (\sigma_{l}^{2} + 6K\sigma_{g}^{2}) \\ &+ \frac{\beta_{1}}{1 - \beta_{1}} \frac{\eta\eta_{l}KG^{2}d}{\sqrt{\epsilon}} + \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \frac{2\eta^{2}\eta_{l}^{2}K^{2}LG^{2}d}{\epsilon} \\ &+ \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1 - \beta_{1})^{2}} \eta^{2}L + \sqrt{2(1 - \beta_{2})}\eta G\right) \frac{KT\eta_{l}^{2}}{n\epsilon} \sigma_{l}^{2} - \sum_{t=1}^{T} \mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}\right] \\ &\cdot \left[\frac{\eta\eta_{l}}{2\sqrt{\beta_{2}K^{2}G^{2} + \epsilon}Km^{2}} - \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1 - \beta_{1})^{2}} \eta^{2}L + \sqrt{2(1 - \beta_{2})}\eta G\right) \frac{\eta_{l}^{2}(m - n)}{mn(m - 1)\epsilon}\right] \\ &+ \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1 - \beta_{1})^{2}} \eta^{2}L + \sqrt{2(1 - \beta_{2})}\eta G\right) \frac{\eta_{l}^{2}(m - n)}{mn(m - 1)\epsilon} \left[15mK^{3}L^{3}\eta_{l}^{2}(\sigma_{l}^{2} + 6K\sigma_{g}^{2})T\right] \\ &+ \left(90mK^{4}L^{2}\eta_{l}^{2} + 3mK^{2}\right)\sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_{t})\|^{2}] + 3mK^{2}T\sigma_{g}^{2}\right] \\ &+ \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1 - \beta_{1})^{2}} \eta^{2}L + \sqrt{2(1 - \beta_{2})}\eta G\right) \frac{2}{n^{2}\epsilon}\sum_{t=1}^{T} \mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{i}^{i}\right\|^{2} + \frac{\sqrt{2}\eta\eta_{l}KGCT}{\alpha n\epsilon}, \\ &+ \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1 - \beta_{1})^{2}} \eta^{2}L + \sqrt{2(1 - \beta_{2})}\eta G\right) \frac{2}{n^{2}\epsilon}\sum_{t=1}^{T} \mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{i}^{i}\right\|^{2} + \frac{\sqrt{2}\eta\eta_{l}KGCT}{\alpha n\epsilon}, \\ &+ \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1 - \beta_{1})^{2}} \eta^{2}L + \sqrt{2(1 - \beta_{2})}\eta G\right) \frac{2}{n^{2}\epsilon}\sum_{t=1}^{T} \mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{i}^{i}\right\|^{2} + \frac{\sqrt{2}\eta\eta_{l}KGCT}{\alpha n\epsilon}, \\ &+ \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1 - \beta_{1})^{2}} \eta^{2}L + \sqrt{2(1 - \beta_{2})}\eta G\right) \frac{2}{n^{2}\epsilon}\sum_{t=1}^{T} \mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{i}^{i}\right\|^{2} + \frac{\sqrt{2}\eta\eta_{l}KGCT}{\alpha n\epsilon}, \\ &+ \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1 - \beta_{1})^{2}} \eta^{2}L + \sqrt{2(1 - \beta_{2})}\eta G\right) \frac{2}{n^{2}\epsilon}\sum_{t=1}^{T} \mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{i}^{i}\right\|^{2} + \frac{\sqrt{2}\eta\eta_{l}KGCT}{\alpha n\epsilon}, \\ &+ \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1 - \beta_{1})^{2}} \eta^{2}L + \sqrt{2(1 - \beta_{2})}\eta G\right) \frac{2}{n^{2}\epsilon}\sum_{t=1}^{T} \mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{i}^{i}\right\|^{2} + \frac{M_{t}^{2}}{\alpha n\epsilon}\right)$$

then

$$\begin{aligned} & \mathbb{E}[f(\mathbf{z}_{T+1})] - f(\mathbf{z}_{1}) \\ & = -\frac{\eta\eta_{l}K}{4\sqrt{\beta_{2}\eta_{l}^{2}K^{2}G^{2} + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_{t})\|^{2}] + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}T}{\sqrt{2\epsilon}} (\sigma_{l}^{2} + 6K\sigma_{g}^{2}) + \frac{\beta_{1}}{1 - \beta_{1}} \frac{\eta\eta_{l}KG^{2}d}{\sqrt{\epsilon}} \\ & + \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \frac{2\eta^{2}\eta_{l}^{2}K^{2}LG^{2}d}{\epsilon} + \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1 - \beta_{1})^{2}}\eta^{2}L + \sqrt{2(1 - \beta_{2})}\eta G\right) \left(\frac{KT\eta_{l}^{2}}{n\epsilon}\sigma_{l}^{2} + \frac{2T\eta_{l}^{2}C^{2}K^{2}G^{2}}{\alpha^{2}n^{2}\epsilon}\right) \\ & + \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1 - \beta_{1})^{2}}\eta^{2}L + \sqrt{2(1 - \beta_{2})}\eta G\right) \frac{\eta_{l}^{2}(m - n)}{mn(m - 1)\epsilon} \left[15mK^{3}L^{3}\eta_{l}^{2}(\sigma_{l}^{2} + 6K\sigma_{g}^{2})T\right] \\ & + (90mK^{4}L^{2}\eta_{l}^{2} + 3mK^{2})\sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_{t})\|^{2}] + 3mK^{2}T\sigma_{g}^{2}\right] + \frac{\sqrt{2}\eta\eta_{l}KGCT}{\alpha n\epsilon}. \end{aligned}$$

By adopting additional constraint of local learning rate 
$$\eta_l$$
 with  
the inequality  $\left(\frac{3\eta^2 L}{2} + \frac{\beta_1}{2(1-\beta_1)^2}\eta^2 L + \sqrt{2(1-\beta_2)}\eta G\right)\frac{2\eta_1^2(n-1)}{mn(m-1)\epsilon} - \frac{\eta_l}{mn(m-1)\epsilon} - \frac{\eta_l}{2\sqrt{\beta_2 K^2 G^2 + \epsilon} Km^2} \leq 0$ , thus we obtain the constraint  $\eta_l \leq \frac{n(m-1)}{m(n-1)}\frac{\epsilon}{\sqrt{\beta_2 K^2 G^2 + \epsilon} K(3\eta L + C_1^2\eta L + 2\sqrt{2(1-\beta_2)}G)}$ , and we further need  $\eta_l$  satisfies  $\frac{\eta_l K}{4\sqrt{\beta_2 \eta_l^2 K^2 G^2 + \epsilon}} - \frac{\left(\frac{3\eta^2 L}{2} + \frac{\beta_1^2}{2(1-\beta_1)^2}\eta^2 L + \frac{\eta_l}{4\sqrt{\beta_2 \eta_l^2 K^2 G^2 + \epsilon}} + \sqrt{2(1-\beta_2)}\eta G\right)\frac{\eta_1^2(m-n)}{mn(m-1)\epsilon}(90mK^4 L^2 \eta_l^2 + \frac{3mK^2}{8\sqrt{\beta_2 \eta_l^2 K^2 G^2 + \epsilon}})$ . Hence we have the following condition on local learning rate

$$\begin{aligned} & 1832 \\ & 1833 \\ & 1834 \\ & 1834 \\ & 1835 \end{aligned} \qquad \eta \leq \frac{n(m-1)\epsilon}{48m(n-1)} \left[ K\sqrt{\beta_2 K^2 G^2 + \epsilon} \cdot \left( \frac{3\eta L}{2} + \frac{\beta_1^2}{2(1-\beta_1)^2} \eta L + \sqrt{2(1-\beta_2)} G \right) \right]^{-1}, \end{aligned}$$

then

$$\begin{split} & \frac{\eta\eta_{l}K}{8\sqrt{\beta_{2}\eta_{l}^{2}K^{2}G^{2}+\epsilon}\cdot T}\sum_{i=1}^{T}\mathbb{E}[\|\nabla f(\theta_{t})\|^{2}] \\ & \leq \frac{f(\mathbf{z}_{0})-\mathbb{E}[f(\mathbf{z}_{T})]}{T} + \frac{C_{1}\eta\eta_{l}KG^{2}d}{T\sqrt{\epsilon}} + \frac{2C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}LG^{2}d}{T\epsilon} \\ & + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}}{\sqrt{2\epsilon}}(\sigma_{l}^{2}+6K\sigma_{g}^{2}) + \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1-\beta_{1})^{2}}\eta^{2}L + \sqrt{2(1-\beta_{2})}\eta G\right)\left(\frac{K\eta_{l}^{2}}{n\epsilon}\sigma_{l}^{2} + \frac{2\eta_{l}^{2}C^{2}K^{2}G^{2}}{\alpha^{2}n^{2}\epsilon}\right) \\ & + \left(\frac{3\eta^{2}L}{2} + \frac{\beta_{1}^{2}}{2(1-\beta_{1})^{2}}\eta^{2}L + \sqrt{2(1-\beta_{2})}\eta G\right)\frac{\eta_{l}^{2}(m-n)}{mn(m-1)\epsilon}[15mK^{3}L^{2}\eta_{l}^{2}(\sigma_{l}^{2}+6K\sigma_{g}^{2}) + 3mK^{2}\sigma_{g}^{2}] + \frac{\sqrt{2}\eta\eta_{l}KGC}{\alpha n\epsilon} \end{split}$$

Therefore

 $\eta \eta_l K$ 

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 8\sqrt{\beta_2 \eta_l^2 K^2 G^2 + \epsilon} \bigg[ \frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega \bigg],$$

where  $\Xi = \frac{C_1 G^2 d}{\sqrt{\epsilon}} + \frac{2G_1^2 \eta \eta K L G^2 d}{\epsilon}, \Omega = \frac{5\eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K \sigma_q^2) + [(3 + C_l^2) \eta L + 2\sqrt{2(1 - \beta_2)}G](\frac{\eta_l}{n\eta\epsilon}\sigma_l^2 + \frac{2\eta_l K^2 C^2 G^2}{\alpha^2 \eta n^2 \epsilon}) + [(3 + C_1^2) \eta L + 2\sqrt{2(1 - \beta_2)}G]\frac{\eta_l(m-n)}{2n(m-1)\epsilon} [15K^2 L^2 \eta_l^2 (\sigma_l^2 + 6K\sigma_g^2) + 3K\sigma_g^2]\frac{1}{\eta\eta_l K} + \frac{\sqrt{2}GC}{\alpha n\epsilon} \text{ and } C_1 = \frac{\beta_1}{1 - \beta_1}.$ 

The proof of Theorem 4.4 is similar to the above proof procedure and the detailed proof will not be given here. 

C.4 **PROOF OF COROLLARY 4.2** 

If choose  $\eta_l = \Theta(\frac{1}{\sqrt{TK}})$  and  $\eta = \Theta(\sqrt{Kn})$ , we get  $\min_{t \in [T]} \mathbb{E}[\|\nabla f(\theta_t)\|^2] = \mathcal{O}(\frac{\sqrt{K}}{\sqrt{Tn}})$ .

#### **PROOF OF THEOREMS IN SECTIONA** D

#### D.1 **PROOF OF THEOREM A.1**

Notations and equations: From the update rule of Algorithm 3, then  $\mathbf{e}_1 = 0, \mathbf{e}_t =$  $\frac{1}{m}\sum_{i=1}^{m} \mathbf{e}_{t}^{i} \text{ and } \mathbf{m}_{t} = (1-\beta_{1})\sum_{i=1}^{t}\beta_{1}^{t-i}\widehat{\widehat{\Delta}}_{i}. \text{ Denote a global uncompressed difference } \Delta_{t} = \frac{1}{m}\sum_{i=1}^{m}\Delta_{t}^{i}. \text{ Denote a virtual momentum sequence: } \mathbf{m}_{t}' = \beta_{1}\mathbf{m}_{t-1}' + (1-\beta_{1})\Delta_{t}, \text{ hence we have } \mathbf{m}_{t}' = \beta_{1}\mathbf{m}_{t-1}' + (1-\beta_{1})\Delta_{t}.$  $\mathbf{m}_{t}^{m} = (1 - \beta_{1}) \sum_{i=1}^{t} \beta_{1}^{t-i} \Delta_{i}$ . By the aforementioned definition and notation, then 

$$\widehat{\Delta}_t - \Delta_t = \frac{1}{m} \sum_{i=1}^m (\widehat{\widehat{\Delta}}_t^i - \Delta_t^i) = \frac{1}{m} \sum_{i=1}^m (\widehat{\Delta}_t^i - \Delta_t^i) - \frac{1}{M_t} \sum_{i \in M_t} \mathcal{C}(q_t^i)$$

$$= \frac{1}{m} \sum_{i=1}^{m} (\mathbf{e}_{t}^{i} - \mathbf{e}_{t+1}^{i}) - \frac{1}{M_{t}} \sum_{i \in M_{t}} \mathcal{C}(q_{t}^{i}) = \mathbf{e}_{t} - \mathbf{e}_{t+1} - \frac{1}{M_{t}} \sum_{i \in M_{t}} \mathcal{C}(q_{t}^{i}).$$
(D.1)

Denote the weighted averaging error sequence  $\Gamma_t = (1 - \beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \mathbf{e}_r$ , with the input  $\mathbf{e}_1 = 0$ , we obtain the relation between  $\Gamma_t$  and  $\mathbf{m}_t$  as follows

$$\mathbf{m}_{t} - \mathbf{m}_{t}' = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\mathbf{e}_{\tau} - \mathbf{e}_{\tau+1}) - (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} \frac{1}{M_{t}} \sum_{i \in M_{t}} \mathcal{C}(q_{t}^{i}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\mathbf{e}_{\tau} - \mathbf{e}_{\tau+1}) - (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} \frac{1}{M_{t}} \sum_{i \in M_{t}} \mathcal{C}(q_{t}^{i}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\mathbf{e}_{\tau} - \mathbf{e}_{\tau+1}) - (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} \sum_{i \in M_{t}} \mathcal{C}(q_{t}^{i}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\mathbf{e}_{\tau} - \mathbf{e}_{\tau+1}) - (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} \sum_{i \in M_{t}} \mathcal{C}(q_{t}^{i}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\mathbf{e}_{\tau} - \mathbf{e}_{\tau+1}) - (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\mathbf{e}_{\tau} - \mathbf{e}_{\tau+1}) - (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_$$

$$=\Gamma_t - \Gamma_{t+1} - (1 - \beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \frac{1}{M_t} \sum_{i \in M_t} \mathcal{C}(q_t^i),$$
(D.2)

where the last step holds due to  $\Gamma_{t+1} = (1 - \beta_1) \sum_{\tau=1}^{t+1} \beta_1^{t-\tau} \mathbf{e}_{\tau+1} = (1 - \beta_1) \sum_{\tau=1}^{t} \beta_1^{t-\tau} \mathbf{e}_{\tau+1} + (1 - \beta_1) \sum_{\tau=1}^{t} \beta_1^{t-\tau} \mathbf{e}_{\tau+1} + (1 - \beta_1) \sum_{\tau=1}^{t} \beta_1^{t-\tau} \mathbf{e}_{\tau+1} = (1 - \beta_1) \sum_{\tau=1}^{t} \beta_1^{t-\tau} \mathbf{e}_{\tau+1} + (1 - \beta_1) \sum_{\tau=1}^{$  $\beta_1^t \mathbf{e}_1.$ 

Similar to previous works studied adaptive methods (Chen et al. (2018c); Wang et al. (2022a)), we introduce a Lyapunov sequence  $\mathbf{z}_t$ :assume  $\theta_0 = \theta_1$ , for each  $t \ge 1$ , 

$$\mathbf{z}_{t} = \theta_{t} + \frac{\beta_{1}}{1 - \beta_{1}}(\theta_{t} - \theta_{t-1}) = \frac{1}{1 - \beta_{1}}\theta_{t} - \frac{\beta_{1}}{1 - \beta_{1}}\theta_{t-1}$$

Therefore, by the update rule of  $\theta_t$ ,

1901  
1902 
$$\mathbf{y}_{t+1} = \theta_{t+1} + \eta \frac{\beta_1}{1-\beta_1} \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t - (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \frac{1}{M_t} \sum_{i \in M_t} \mathcal{C}(q_t^i)$$
]  
1903

$$=\theta_{t+1} + \eta \frac{\beta_1}{1-\beta_1} \widehat{\mathbf{V}}_t^{-1/2} [\mathbf{m}_t' + \Gamma_t - \Gamma_{t+1}] - (1-\beta_1) \sum_{\tau=1}^{\iota} \beta_1^{t-\tau} \frac{1}{M_t} \sum_{i \in M_t} \mathcal{C}(q_t^i)$$

$$=\theta_{t+1} + \eta \frac{\beta_1}{1-\beta_1} \widehat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t' + \eta \frac{\beta_1}{1-\beta_1} \widehat{\mathbf{V}}_t^{-1/2} \left[ \frac{\Gamma_{t+1} - (1-\beta_1)\mathbf{e}_{t+1}}{\beta_1} - \Gamma_{t+1} \right] - (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \frac{1}{M_t} \sum_{i \in M_t} \mathcal{C}(q_t^i) \sum_{\tau=1}^t \beta_1^{t-\tau} \frac{1}{M_t} \sum_{\tau=$$

$$=\theta_{t+1} + \eta \frac{\beta_1}{1-\beta_1} \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}'_t + \eta \hat{\mathbf{V}}_t^{-1/2} \Gamma_{t+1} - \eta \hat{\mathbf{V}}_t^{-1/2} \mathbf{e}_{t+1} - (1-\beta_1) \sum_{\tau=1}^{\iota} \beta_1^{t-\tau} \frac{1}{M_t} \sum_{i \in M_t} \mathcal{C}(q_t^i).$$
(D.3)

The third equation holds due to the fact that  $\Gamma_{t+1} = \beta_1 \Gamma_t + (1 - \beta_1) \mathbf{e}_{t+1}$ . We then introduce a new sequence based on the previous Lyapunov sequence  $y_t$  as follows 

1918  
1919 
$$\mathbf{z}_{t+1} = \mathbf{y}_{t+1} + (1-\beta_1) \sum_{\tau=1}^{t} \beta_1^{t-\tau} \frac{1}{M_t} \sum_{i \in M_t} C(q_t^i) + \eta \hat{\mathbf{V}}_t^{-1/2} \mathbf{e}_{t+1} = \theta_{t+1} + \eta \frac{\beta_1}{1-\beta_1} \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t' + \eta \hat{\mathbf{V}}_t^{-1/2} \Gamma_{t+1}.$$
(D.4)

The sequence difference  $\mathbf{z}_{t+1} - \mathbf{z}_t$  can be represented by 

$$\begin{aligned} \mathbf{z}_{t+1} &= \theta_{t+1} - \mathbf{z}_t = \theta_{t+1} - \theta_t + \eta \frac{\beta_1}{1 - \beta_1} \widehat{\mathbf{V}}_t^{-1/2} \mathbf{m}'_t - \eta \frac{\beta_1}{1 - \beta_1} \widehat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}'_{t-1} + \eta \widehat{\mathbf{V}}_t^{-1/2} \Gamma_{t+1} - \eta \widehat{\mathbf{V}}_{t-1}^{-1/2} \Gamma_t \\ \\ \mathbf{z}_{t+1} &= \eta \widehat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t + \eta \widehat{\mathbf{V}}_t^{-1/2} \Gamma_{t+1} + \eta \frac{\beta_1}{1 - \beta_t} \widehat{\mathbf{V}}_t^{-1/2} \mathbf{m}'_t - \eta \frac{\beta_1}{1 - \beta_t} \widehat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}'_{t-1} - \eta \widehat{\mathbf{V}}_{t-1}^{-1/2} \Gamma_t, \end{aligned}$$

where the second equation follows the update rule of  $\theta_{t+1}$ . Following (D.2), then combining likely terms and applying the definition of  $\mathbf{m}_t^{\prime}$ , then

$$\mathbf{z}_{t+1} - \mathbf{z}_t = \eta \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t' + \eta \hat{\mathbf{V}}_t^{-1/2} \Gamma_t + \eta \frac{\beta_1}{1 - \beta_1} \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t' - \eta \frac{\beta_1}{1 - \beta_1} \hat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}_{t-1}' - \eta \hat{\mathbf{V}}_{t-1}^{-1/2} \Gamma_t$$

$$= \eta \frac{1}{1-\beta_1} \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t' - \eta \frac{\beta_1}{1-\beta_1} \hat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}_{t-1}' + \eta \hat{\mathbf{V}}_t^{-1/2} \Gamma_t - \eta \hat{\mathbf{V}}_{t-1}^{-1/2} \Gamma,$$

$$=\eta \frac{1}{1-\beta_1} \widehat{\mathbf{V}}_t^{-1/2} [\beta_1 \mathbf{m}_{t-1}' + (1-\beta_1) \Delta_t] - \eta \frac{\beta_1}{1-\beta_1} \widehat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}_{t-1}' + \eta \widehat{\mathbf{V}}_t^{-1/2} \Gamma_t - \eta \widehat{\mathbf{V}}_{t-1}^{-1/2} \Gamma_t$$

1939  
1940
$$= \eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} - \eta \frac{\beta_{1}}{1 - \beta_{1}} \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) \mathbf{m}_{t-1}' - \eta \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) \Gamma_{t}. \quad (D.5)$$
1941

Therefore, we obtain a helpful Lyapunov sequence for our proof of FedCAMS. The proof of Fed-

CAMS in full participation settings has a similar outline with the proof of FedAMS. By Assumption 3.1, then

$$\begin{split} & \mathbb{E}[f(\mathbf{z}_{t+1})] - f(\mathbf{z}_{t}) \\ & \leq \mathbb{E}[\langle \nabla f(\mathbf{z}_{t}), \mathbf{z}_{t+1} - \mathbf{z}_{t} \rangle] + \frac{L}{2} \mathbb{E}[\|\mathbf{z}_{t+1} - \mathbf{z}_{t}\|^{2}] \\ & \leq \mathbb{E}\Big[\left\langle \nabla f(\mathbf{z}_{t}), \eta \widehat{\mathbf{V}_{t}^{-1/2}} \Delta_{t} \right\rangle \Big] \\ & - \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \eta \frac{\beta_{1}}{1 - \beta_{1}} \left(\widehat{\mathbf{V}_{t-1}^{-1/2}} - \widehat{\mathbf{V}_{t}^{-1/2}}\right) \mathbf{m}_{t-1}' + \eta \left(\widehat{\mathbf{V}_{t-1}^{-1/2}} - \widehat{\mathbf{V}_{t}^{-1/2}}\right) \mathbf{r}_{t} \right\rangle \Big] \\ & + \frac{\eta^{2}L}{2} \mathbb{E}\left[\left\|\widehat{\mathbf{V}_{t}^{-1/2}} \Delta_{t} - \frac{\beta_{1}}{1 - \beta_{1}} \left(\widehat{\mathbf{V}_{t-1}^{-1/2}} - \widehat{\mathbf{V}_{t}^{-1/2}}\right) \mathbf{m}_{t-1}' - \left(\widehat{\mathbf{V}_{t-1}^{-1/2}} - \widehat{\mathbf{V}_{t}^{-1/2}}\right) \Gamma_{t} \right\|^{2} \right] \\ & = \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}_{t}^{-1/2}} \Delta_{t} \right\rangle - \eta \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\beta_{1}}{1 - \beta_{1}} \left(\widehat{\mathbf{V}_{t-1}^{-1/2}} - \widehat{\mathbf{V}_{t}^{-1/2}}\right) \mathbf{m}_{t-1}' + \left(\widehat{\mathbf{V}_{t-1}^{-1/2}} - \widehat{\mathbf{V}_{t}^{-1/2}}\right) \Gamma_{t} \right\rangle \right] \\ & + \underbrace{\eta^{2}L}_{2} \mathbb{E}\left[\left\|\widehat{\mathbf{V}_{t}^{-1/2}} \Delta_{t} - \frac{\beta_{1}}{1 - \beta_{1}} \left(\widehat{\mathbf{V}_{t-1}^{-1/2}} - \widehat{\mathbf{V}_{t}^{-1/2}}\right) \mathbf{m}_{t-1}' - \left(\widehat{\mathbf{V}_{t-1}^{-1/2}} - \widehat{\mathbf{V}_{t}^{-1/2}}\right) \Gamma_{t} \right\|^{2}\right] \\ & - \frac{\tau_{3}}}{\tau_{3}} \\ & + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}_{t}^{-1/2}} \Delta_{t} \right\rangle\right]}_{T_{4}}. \end{split}$$
(D.6)

here we recall the notation  $\hat{\mathbf{V}}_t = \text{diag}(\widehat{\mathbf{v}}_t) = \text{diag}(\max(\widehat{\mathbf{v}}_{t-1}, \mathbf{v}_t, \epsilon)).$ **Bounding**  $T_1$ :

 $\mathbf{T}\mathbf{1} = \mathbb{E}\left[\left\langle \nabla f(\theta_t), \eta \frac{\Delta_t}{\sqrt{\hat{\mathbf{v}}_t}} \right\rangle\right]$ 

 $\leq \eta \mathbb{E}\left[\left\langle \nabla f(\theta_t), \frac{\sqrt{2} \cdot \Delta_t}{\sqrt{\mathbf{v}_t + \epsilon}} \right\rangle\right]$ 

where the first inequality follows by the fact that  $\hat{v}_t \geq \frac{v_t + \epsilon}{2}$ . For the second term in (D.7), then

 $= \sqrt{2}\eta \mathbb{E}\bigg[\left\langle \nabla f(\theta_t), \frac{\Delta_t}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\rangle\bigg] + \sqrt{2}\eta \mathbb{E}\bigg[\left\langle \nabla f(\theta_t), \frac{\Delta_t}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{\Delta_t}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\rangle\bigg],$ 

(D.7)

$$\sqrt{2} \cdot \eta \mathbb{E} \left[ \left\langle \nabla f(\theta_t), \frac{\Delta_t}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{\Delta_t}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\rangle \right]$$

$$\leq \sqrt{2} \cdot \eta \mathbb{E} \left\| \nabla f(\theta_t) \| \mathbb{E} \left[ \left\| \frac{1}{2} - \frac{1}{2} - \frac{1}{2} \right\| \frac{1}{2} \right] \right\|$$

$$\leq \sqrt{2} \cdot \eta \cdot \mathbb{E} \|\nabla f(\theta_t)\| \mathbb{E} \left[ \left\| \frac{1}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{1}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\| \cdot \|\Delta_t\| \right]$$
$$\leq \frac{\eta \sqrt{2(1 - \beta_2)}G}{\epsilon} \mathbb{E} [\|\Delta_t\|^2],$$

where the second inequality follows from Lemma F.1 and F.5, and we will further apply the bound for  $E[||\Delta_t||^2]$  by applying Lemma F.7. For the first term in (D.6), then

$$\begin{aligned} 1998\\ 1999\\ 2000\\ 2001\\ 2002\\ 2002\\ 2002\\ 2003\\ = \sqrt{2} \cdot \eta \mathbb{E}\left[\left\langle \nabla f(\theta_t), \frac{\Delta_t}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\rangle\right] \\ = \sqrt{2} \cdot \eta \mathbb{E}\left[\left\langle \frac{\nabla f(\theta_t)}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}}, \Delta_t + \eta K \nabla f(\theta_t) - \eta_t K \nabla f(\theta_t) \right\rangle\right] \\ 2004\\ 2005\\ = -\sqrt{2}\eta \eta_l K \mathbb{E}\left[\left\|\frac{\nabla f(\theta_t)}{\sqrt[4]{\beta_2 \mathbf{v}_{t-1} + \epsilon}}\right\|^2\right] + \sqrt{2}\eta \mathbb{E}\left[\left\langle \frac{\nabla f(\theta_t)}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}}, \Delta_t + \eta_l K \nabla f(\theta_t) \right\rangle\right] \\ 2008\\ = -\sqrt{2}\eta \eta_l K \mathbb{E}\left[\left\|\frac{\nabla f(\theta_t)}{\sqrt[4]{\beta_2 \mathbf{v}_{t-1} + \epsilon}}\right\|^2\right] + \sqrt{2}\eta \left\langle \frac{\nabla f(\theta_t)}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}}, \mathbb{E}\left[-\frac{\eta_l}{m}\sum_{i=1}^m\sum_{k=0}^{m} \mathbf{g}_{i,k}^i + \eta_l K \nabla f(\theta_t)\right]\right\rangle \\ 2012\\ = -\sqrt{2}\eta \eta_l K \mathbb{E}\left[\left\|\frac{\nabla f(\theta_t)}{\sqrt[4]{\beta_2 \mathbf{v}_{t-1} + \epsilon}}\right\|^2\right] + \sqrt{2}\eta \left\langle \frac{\nabla f(\theta_t)}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}}, \mathbb{E}\left[-\frac{\eta_l}{m}\sum_{i=1}^m\sum_{k=0}^{m-1} \mathbf{g}_{i,k}^i + \frac{\eta_l K}{m}\sum_{i=1}^m \nabla F_i(\theta_t)\right]\right\rangle \\ 2014\\ (D.8)
\end{aligned}$$

# For the last term in (D.8),

$$\begin{split} &\sqrt{2}\eta \left\langle \frac{\nabla f(\theta_{t})}{\sqrt{\beta_{2}\mathbf{v}_{t-1}+\epsilon}}, \mathbb{E}\bigg[ -\frac{\eta_{l}}{m}\sum_{i=1}^{m}\sum_{k=0}^{K-1}\mathbf{g}_{i,k}^{i} + \frac{\eta_{l}K}{m}\sum_{i=1}^{m}\nabla F_{i}(\theta_{t})\bigg] \right\rangle \\ &= \sqrt{2}\eta \left\langle \frac{\sqrt{\eta K}}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}} \nabla f(\theta_{t}), -\frac{\sqrt{\eta_{n}K}}{Km}\frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}} \mathbb{E}\bigg[\sum_{i=1}^{m}\sum_{k=0}^{K-1}(\nabla F_{i}(\theta_{i,k}^{i}) - \nabla F_{i}(\theta_{t}))\bigg] \right\rangle \\ &= \frac{\sqrt{2}\eta\eta_{l}K}{2} \left\| \frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}} \right\|^{2} + \frac{\sqrt{2}\eta\eta_{l}}{2Km^{2}} \mathbb{E} \left\| \frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}} \sum_{i=1}^{m}\sum_{k=0}^{K-1}(\nabla F_{i}(\theta_{i,k}^{i}) - \nabla F_{i}(\theta_{t})) \bigg\|^{2} \\ &- \frac{\sqrt{2}\eta\eta_{l}}{2Km^{2}} \mathbb{E} \left[ \left\| \frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}} \sum_{i=1}^{m}\sum_{k=0}^{K-1} \nabla F_{i}(\theta_{i,k}^{i}) \right\|^{2} \right] \\ &\leq \frac{\sqrt{2}\eta\eta_{l}K}{2} \left\| \frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}} \right\|^{2} + \frac{\sqrt{2}\eta\eta_{l}}{2m} \cdot \sum_{i=1}^{m}\sum_{k=0}^{K-1} \mathbb{E} \left[ \left\| \frac{\nabla F_{i}(\theta_{i,k}^{i}) - \nabla F_{i}(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}} \right\|^{2} \right] \\ &- \frac{\sqrt{2}\eta\eta_{l}}{2Km^{2}} \mathbb{E} \bigg[ \left\| \frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}+\epsilon}} \sum_{i=1}^{m}\sum_{k=0}^{K-1} \nabla F_{i}(\theta_{i,k}^{i}) \right\|^{2} \bigg], \end{split}$$

where the second equation follows from  $\langle \mathbf{x}, \mathbf{y} \rangle = \frac{1}{2} [\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2]$ , and the inequality holds by applying Cauchy-Schwarz inequality. Then by Assumption 3.1,

where the last inequality holds by applying Lemma F.14 and the constraint of local learning rate  $\eta_n \leq \frac{1}{8KL}$ . Then

$$T_{1} \leq -\frac{\sqrt{2} \cdot \eta \eta_{l} K}{4} \mathbb{E} \left[ \left\| \frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2} \mathbf{v}_{t-1} + \epsilon}} \right\|^{2} \right] + \frac{5 \eta \eta_{l}^{3} K^{2} L^{2}}{\sqrt{2\epsilon}} \left( \sigma_{l}^{2} + 6K \sigma_{g}^{2} \right) - \frac{\sqrt{2} \cdot \eta \eta_{l}}{2Km^{2}} \mathbb{E} \left[ \left\| \frac{1}{\sqrt[4]{\beta_{2} \mathbf{v}_{t-1} + \epsilon}} \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \right] + \frac{\eta \sqrt{2(1 - \beta_{2})} G}{\epsilon} \mathbb{E} \left[ \left\| \Delta_{t} \right\|^{2} \right] \leq -\frac{\eta \eta_{l} K}{4} \mathbb{E} \left[ \left\| \frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2} \mathbf{v}_{t-1} + \epsilon}} \right\|^{2} \right] + \frac{5 \eta \eta_{l}^{3} K^{2} L^{2}}{\sqrt{2\epsilon}} \left( \sigma_{l}^{2} + 6K \sigma_{g}^{2} \right) - \frac{\eta \eta_{l}}{2Km^{2}} \mathbb{E} \left[ \left\| \frac{1}{\sqrt[4]{\beta_{2} \mathbf{v}_{t-1} + \epsilon}} \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \right] + \frac{\eta \sqrt{2(1 - \beta_{2})} G}{\epsilon} \mathbb{E} \left[ \left\| \Delta_{t} \right\|^{2} \right].$$
(D.9)

**Bounding**  $T_2$ : The bound for  $T_2$  mainly follows by the update rule and definition of virtual sequence  $\mathbf{z}_t$ .

$$T_{2} = -\eta \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\beta_{1}}{1-\beta_{1}} \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \mathbf{m}_{t-1}' + \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \Gamma_{t}\right\rangle\right]$$

$$= \eta \mathbb{E}\left[\left\langle -\nabla f(\theta_{t}) + \nabla f(\theta_{t}) - \nabla f(\mathbf{z}_{t}), \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \left(\frac{\beta_{1}}{1-\beta_{1}}\mathbf{m}_{t-1}' + \Gamma_{t}\right)\right\rangle\right]$$

$$\leq \eta \mathbb{E}\left[\left\|\nabla f(\theta_{t})\right\| \left\| \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \left(\frac{\beta_{1}}{1-\beta_{1}}\mathbf{m}_{t-1}' + \Gamma_{t}\right)\right\|\right]$$

$$+ \eta^{2} L \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t-1}^{-1/2} \left(\frac{\beta_{1}}{1-\beta_{1}}\mathbf{m}_{t-1}' + \Gamma_{t}\right)\right\| \left\| \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \left(\frac{\beta_{1}}{1-\beta_{1}}\mathbf{m}_{t-1}' + \Gamma_{t}\right)\right\|\right]$$

$$\leq \eta C_{1} \eta_{l} K G^{2} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right\|_{1}\right] + \eta^{2} C_{1}^{2} L \eta_{l}^{2} K^{2} G^{2} \epsilon^{-1/2} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right\|_{1}\right], \quad (D.10)$$

where the last inequality holds by Lemma F.5, here  $C_1 = \frac{\beta_1}{1-\beta_1} + \sqrt{\frac{12q^2}{(1-q^2)^2} + \frac{(1-q^2)^2C^2}{\alpha^2m^2q^2}}$ . Bounding  $T_3$ :

$$T_{3} = \frac{\eta^{2}L}{2} \mathbb{E}\left[\left\|\hat{\mathbf{V}}_{t}^{-1/2}\Delta_{t} + \frac{\beta_{1}}{1-\beta_{1}}\left(\hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2}\right)\mathbf{m}_{t-1}' + \left(\hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2}\right)\Gamma_{t}\right\|^{2}\right]$$

$$\leq \eta^{2}L\mathbb{E}\left[\left\|\hat{\mathbf{V}}_{t}^{-1/2}\Delta_{t}\right\|^{2}\right] + \eta^{2}L\mathbb{E}\left[\left\|\frac{\beta_{1}}{1-\beta_{1}}\left(\hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2}\right)\mathbf{m}_{t-1}' + \left(\hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2}\right)\Gamma_{t}\right\|^{2}\right]$$

$$\leq \eta^{2}L\mathbb{E}\left[\left\|\hat{\mathbf{V}}_{t}^{-1/2}\Delta_{t}\right\|^{2}\right] + \eta^{2}LC_{1}^{2}\eta_{l}^{2}K^{2}G^{2}\mathbb{E}\left[\left\|\hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2}\right\|^{2}\right], \qquad (D.11)$$

where the first inequality follows by Cauchy-Schwarz inequality, and the second one follows by Lemma F.5, here  $C_1 = \frac{\beta_1}{1-\beta_1} + \sqrt{\frac{12q^2}{(1-q^2)^2} + \frac{(1-q^2)^2C^2}{\alpha^2m^2q^2}}$ .

**Bounding**  $T_4$ :

$$T_4 = \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_t) - \nabla f(\theta_t), \eta \widehat{\mathbf{V}}_t^{-1}\right.\right]$$

$$T_{4} = \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\rangle\right]$$
$$\leq \mathbb{E}\left[\left\|\nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t})\right\| \left\|\eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t}\right\|\right]$$

$$\leq L \mathbb{E} \left[ \| \mathbf{z}_t - heta_t \| \| \eta \widehat{\mathbf{V}}_t^{-1/2} \Delta_t \| 
ight.$$

$$\leq \frac{\eta^2 L}{2} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_t^{-1/2} \Delta_t\right\|^2\right] + \frac{\eta^2 L}{2} \mathbb{E}\left[\left\|\frac{\beta_1}{1-\beta_1}\widehat{\mathbf{V}}_{t-1}^{-1/2}\mathbf{m}_{t-1}' + \widehat{\mathbf{V}}_{t-1}^{-1/2}\mathbf{\Gamma}_t\right\|^2\right],$$

where the first inequality holds by the fact of  $\langle \mathbf{a}, \mathbf{b} \rangle \leq \|\mathbf{a}\| \|\mathbf{b}\|$ , the second one follows from Assumption 3.1 and the third one holds by the definition of virtual sequence  $z_t$  and the fact of  $\|\mathbf{a}\| \|\mathbf{b}\| \le \frac{1}{2} \|\mathbf{a}\|^2 + \frac{1}{2} \|\mathbf{b}\|^2$ . Then summing  $T_4$  over  $t = 1, \dots, T$ , 

$$\sum_{t=1}^{T} T_{4} \leq \frac{\eta^{2} L}{2} \sum_{t=1}^{T} \mathbb{E} \left[ \left\| \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\|^{2} \right] + \frac{\eta^{2} L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E} \left[ \left\| \frac{\beta_{1}}{1 - \beta_{1}} \mathbf{m}_{t-1}^{\prime} + \Gamma_{t} \right\|^{2} \right] \\ \leq \frac{\eta^{2} L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E} [\left\| \Delta_{t} \right\|^{2}] + \frac{\eta^{2} L}{\epsilon} \left[ \frac{\beta_{1}^{2}}{(1 - \beta_{1})^{2}} \sum_{t=1}^{T} \mathbb{E} \|\mathbf{m}_{t-1}^{\prime}\|^{2} + \sum_{t=1}^{T} \mathbb{E} \|\mathbf{\Gamma}_{t}\|^{2} \right].$$
(D.12)

By Lemma F.11,

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathbf{m}_{t-1}'\|^2] \le \frac{TK\eta_l^2}{m} \sigma_l^2 + \frac{\eta_l^2}{m^2} \sum_{t=1}^{T} \mathbb{E}\Big[\Big\|\sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i)\Big\|^2\Big],$$

and

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathbf{\Gamma}_t\|^2] \le \frac{4T(q+\gamma)^2}{(1-q^2)^2} \frac{K\eta_l^2}{m} \sigma_l^2 + \frac{\eta_l^2}{m^2} \frac{4(q+\gamma)^2}{(1-q^2)^2} \sum_{t=1}^{T} \mathbb{E}\left[\left\|\sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i)\right\|^2\right].$$

Therefore, the  $T_4$  term is bounded by

$$\sum_{t=1}^{T} T_4 \le \frac{\eta^2 L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\Delta_t\|^2] + \frac{C_2 \eta^2 L}{\epsilon} \frac{\eta_l^2}{m^2} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i) \right\|^2 \right] + \frac{C_2 \eta^2 L}{\epsilon} \frac{T K \eta_l^2}{m} \sigma_l^2, \tag{D.13}$$

where  $C_2 = \frac{4(q+\gamma+\frac{\lambda C}{\alpha m})^2}{(1-q^2)^2} + \frac{\beta_1^2}{(1-\beta_1)^2}.$ 

Merging pieces together: Substituting (D.9), (D.10) and (D.11) into (D.6), summing over from t = 1to T and then adding (D.13),

$$\mathbb{E}[f(\mathbf{z}_{T+1})] - f(\mathbf{z}_{1}) = \sum_{t=1}^{T} [T_{1} + T_{2} + T_{3} + T_{4}]$$

$$\leq -\frac{\eta\eta K}{4} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \right\|^{2} \right] + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}T}{\sqrt{2\epsilon}} (\sigma_{l}^{2} + 6K\sigma_{g}^{2}) + \frac{\sqrt{2(1 - \beta_{2})}\eta G}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\left\| \Delta_{t} \right\|^{2}]$$

$$- \frac{\eta m}{2Km^{2}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i})) \right\|^{2} \right] + C_{1}\eta\eta_{t}KG^{2} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \widehat{\mathbf{v}}_{t-1}^{-1/2} - \widehat{\mathbf{v}}_{t}^{-1/2} \right\|^{2} \right]$$

$$+ \frac{C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}G^{2}}{\sqrt{\epsilon}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \widehat{\mathbf{v}}_{t-1}^{-1/2} - \widehat{\mathbf{v}}_{t}^{-1/2} \right\|_{1}^{2} \right] + C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}LG^{2} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \widehat{\mathbf{v}}_{t-1}^{-1/2} - \widehat{\mathbf{v}}_{t}^{-1/2} \right\|^{2} \right]$$

$$+ \eta^{2}L \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \widehat{\mathbf{v}}_{t}^{-1/2}\Delta_{t} \right\|^{2} \right] + \frac{\eta^{2}L}{2} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \widehat{\mathbf{v}}_{t}^{-1/2}\Delta_{t} \right\|^{2} \right] + \frac{\eta^{2}L}{2} \sum_{t=1}^{T} \mathbb{E}[\left\| \mathbf{n}_{t}^{*} \right\|^{2}]. \quad (D.14)$$

Hence by organizing and applying Lemmas F.5, then

$$\begin{split} \mathbb{E}[f(\mathbf{z}_{T+1})] &- f(\mathbf{z}_{1}) \\ &\leq -\frac{\eta \eta_{l} K}{4} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2} \mathbf{v}_{t-1} + \epsilon}} \right\|^{2} \right] + \frac{5 \eta \eta_{l}^{3} K^{2} L^{2} T}{\sqrt{2\epsilon}} (\sigma_{l}^{2} + 6K \sigma_{g}^{2}) \\ &- \frac{\eta \eta_{l}}{2K m^{2}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{1}{\sqrt[4]{\beta_{2} \mathbf{v}_{t-1} + \epsilon}} \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i})) \right\|^{2} \right] + \frac{C_{1} \eta \eta_{l} K G^{2} d}{\sqrt{\epsilon}} + \frac{2C_{1}^{2} \eta^{2} \eta_{l}^{2} K^{2} L G^{2} d}{\epsilon} \\ &+ \left( \eta^{2} L + \frac{\eta^{2} L}{2} + \sqrt{2(1 - \beta_{2})} \eta G \right) \left[ \frac{K T \eta_{l}^{2}}{m \epsilon} \sigma_{l}^{2} + \frac{\eta_{l}^{2}}{m^{2} \epsilon} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \right] \right] \\ &+ \frac{\eta^{2} L}{\epsilon} \frac{\eta_{l}^{2} C_{2}}{m^{2}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \right] + \frac{\eta^{2} L}{\epsilon} \frac{T K \eta_{l}^{2} C_{2}}{m} \sigma_{l}^{2}, \end{split}$$

by applying Lemma F.6 into all terms containing the second moment estimate of model difference  $\Delta_t$  in (D.14), using the fact that  $\left(\sqrt{\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} K^2 G^2 + \epsilon}\right)^{-1} \|\theta\| \leq \left(\sqrt{\beta_2 \frac{(1+q^2)^3}{(1-q^2)} \eta_l^2 K^2 G^2 + \epsilon}\right)^{-1} \|\theta\| \leq \|\frac{\theta}{\sqrt{\beta_2 \mathbf{v} + \epsilon}}\| \leq \epsilon^{-1/2} \|\theta\|$ , and applying Lemma F.3 and F.13,

$$\mathbb{E}[f(\mathbf{z}_{T+1})] - f(\mathbf{z}_1) \\ \leq -\frac{\eta \eta_l K}{4\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 + \epsilon}} \sum_{t=1}^T \mathbb{E}[\|\nabla f(\theta_t)\|^2] + \frac{5\eta \eta_l^3 K^2 L^2 T}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_g^2)$$

$$\begin{aligned} & + \frac{C_{1}\eta\eta_{l}KG^{2}d}{\sqrt{\epsilon}} + \frac{2C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}LG^{2}d}{\epsilon} + \left(\frac{3\eta^{2}L}{2} + C_{2}\eta^{2}L + \sqrt{2(1-\beta_{2})}\eta G\right)\frac{KT\eta_{l}^{2}}{m\epsilon}\sigma_{l}^{2} + \frac{\eta_{l}\eta TC^{2}K^{2}G^{2}}{\alpha^{2}m^{2}\epsilon} \\ & = \sum_{t=1}^{T} \mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}\right]\left[\frac{\eta\eta_{l}}{2\sqrt{4\beta_{2}\frac{(1+q^{2})^{3}}{(1-q^{2})^{2}}\eta_{l}^{2}K^{2}G^{2} + \epsilon}Km^{2}} - \left(\frac{3\eta^{2}L}{2} + C_{2}\eta^{2}L + \sqrt{2(1-\beta_{2})}\eta G\right)\frac{\eta_{l}^{2}}{m^{2}\epsilon}\right] \\ & = \frac{\eta\eta_{l}K}{4C_{0}}\sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_{t})\|^{2}] + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}T}{\sqrt{2\epsilon}}(\sigma_{l}^{2} + 6K\sigma_{g}^{2}) \\ & = \frac{C_{1}\eta\eta_{l}KG^{2}d}{\sqrt{\epsilon}} + \frac{2C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}LG^{2}d}{\epsilon} + \left(\frac{3\eta^{2}L}{2} + C_{2}\eta^{2}L + \sqrt{2(1-\beta_{2})}\eta G\right)\frac{KT\eta_{l}^{2}}{m\epsilon}\sigma_{l}^{2}, \\ & \text{where the last inequality holds by } \eta_{l} \leq \frac{\epsilon}{1-2} \\ & = \frac{1}{2}\sum_{k=0}^{K}\sum_{l=1}^{K}\sum_{k=0}^{K}\sum_{l=1}^{K}\sum_{k=0}^{K}\sum_{l=1}^{K}\sum_{k=0}^{K}\sum_{l=1}^{K}\sum_{k=0}^{K}\sum_{l=1}^{K$$

 $\sqrt{4\beta_2(1+q^2)^3(1-q^2)^{-2}K^2G^2 + \epsilon \cdot K(3\eta L + 2C_2\eta L + 2\sqrt{2(1-\beta_2)}G)}$ Here

$$\frac{\eta\eta_{l}K}{4\sqrt{4\beta_{2}\frac{(1+q^{2})^{3}}{(1-q^{2})^{2}}\eta_{l}^{2}K^{2}G^{2}+\epsilon \cdot T}}\sum_{t=1}^{T}\mathbb{E}[\|\nabla f(\theta_{t})\|^{2}] \\
\leq \frac{f(\mathbf{z}_{0})-\mathbb{E}[f(\mathbf{z}_{T})]}{T}+\frac{5\eta\eta_{l}^{3}K^{2}L^{2}}{\sqrt{2\epsilon}}(\sigma_{l}^{2}+6K\sigma_{g}^{2})+\frac{C_{1}\eta\eta_{l}KG^{2}d}{T\sqrt{\epsilon}}+\frac{2C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}LG^{2}d}{T\epsilon} \\
+\left[3\eta^{2}L+2C_{2}\eta^{2}L+2\sqrt{2(1-\beta_{2})}\eta G\right]\frac{K\eta_{l}^{2}}{2m\epsilon}\sigma_{l}^{2},$$
(D.15)

where  $C_1 = \frac{\beta_1}{1-\beta_1} + \sqrt{\frac{12q^2}{(1-q^2)^2} + \frac{(1-q^2)^2C^2}{\alpha^2m^2q^2}}$  and  $C_2 = \frac{\beta_1^2}{(1-\beta_1)^2} + \frac{4(q+\gamma+\frac{\lambda C}{\alpha m})^2}{(1-q^2)^2}$ . (D.15) also implies,

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 4\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_t^2 K^2 G^2 + \epsilon} \Big[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\Big],$$

where  $\Xi = \frac{C_1 G^2 d}{\sqrt{\epsilon}} + \frac{2C_1^2 \eta \eta K L G^2 d}{\epsilon}, \Omega = \frac{5\eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_g^2) + [(3 + 2C_2)\eta L + 2\sqrt{2(1-\beta_2)}G] \frac{\eta_l}{2m\eta\epsilon} \sigma_l^2, C_1 = \frac{\beta_1}{1-\beta_1} + \sqrt{\frac{12q^2}{(1-q^2)^2} + \frac{(1-q^2)^2 C^2}{\alpha^2 m^2 q^2}} \text{ and } C_2 = \frac{\beta_1^2}{(1-\beta_1)^2} + \frac{4(q+\gamma+\frac{\lambda C}{\alpha m})^2}{(1-q^2)^2}.$ 

The proof of Theorem A.2 is similar to the above proof procedure and the detailed proof will not be given here.

D.2 PROOF OF COROLLARY A.1

> Let  $\eta_l = \Theta(\frac{1}{\sqrt{TK}}), \eta = \Theta(\sqrt{Km})$  and  $T = \mathcal{O}(Km)$ , the convergence rate under full participation scheme is  $\mathcal{O}(\frac{1}{T})$ .

D.3 PROOF OF THEOREM A.3 

#### **Proof of Theorem A.3:**

Notations and equations: From the update rule of Algorithm 3, then  $\mathbf{e}_1 = 0, \mathbf{e}_t =$  $\frac{1}{m}\sum_{i=1}^{m} \mathbf{e}_{t}^{i} \text{ and } \mathbf{m}_{t} = (1-\beta_{1})\sum_{i=1}^{t}\beta_{1}^{t-i}\widehat{\Delta}_{t}^{i}. \text{ Denote a global uncompressed difference } \Delta_{t} = \frac{1}{|S_{t}|}\sum_{i\in\mathcal{S}_{t}}\Delta_{t}^{i}. \text{ Denote a virtual momentum sequence: } \mathbf{m}_{t}^{\prime} = \beta_{1}\mathbf{m}_{t-1}^{\prime} + (1-\beta_{1})\Delta_{t}, \text{ hence}$  $\mathfrak{m}'_t = (1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} \Delta_i$ . Define additional two virtual sequences  $\Delta'_t = \frac{1}{n} \sum_{i=1}^m \Delta_t^i$  and  $\hat{\Delta}'_t = \frac{1}{n} \sum_{i=1}^m \hat{\Delta}^i_t$ . Note that when the client *i* does not take part in the round of participation at step t, we have  $\Delta_t^i = \widehat{\Delta}_t^i = 0$ , therefore,  $\Delta_t' = \Delta_t$  and  $\widehat{\Delta}_t' = \widehat{\Delta}_t$ . 

By the aforementioned definition and notation, define a subset  $S_t = \{w_1^t, w_2^t, ..., w_n^t\}$ , we have

$$\widehat{\Delta}_t - \Delta_t = \frac{1}{|\mathcal{S}_t|} \sum_{i \in \mathcal{S}_t} (\widehat{\Delta}_t^i - \Delta_t^i) = \frac{1}{n} \sum_{i=1}^m (\widehat{\Delta}_t^i - \Delta_t^i) = \frac{1}{n} \sum_{i=1}^m (\mathbf{e}_t^i - \mathbf{e}_{t+1}^i) = \mathbf{e}_t' - \mathbf{e}_{t+1}',$$

where the compression errors have the same structure,  $\mathbf{e}'_t = \frac{1}{n} \sum_{i=1}^{m} \mathbf{e}^i_t$ . Similar to the previous analysis, define the following sequence:

$$\Gamma_{t+1} := (1 - \beta_1) \sum_{\tau=1}^{t+1} \beta_1^{t+1-\tau} \mathbf{e}_{\tau}'.$$

and keep using the Lyapunov function  $\mathbf{z}_t$  from (D.4). For the expectation of model difference  $\Delta_t$ , then

$$\mathbb{E}_{\mathcal{S}_t}[\Delta_t] = \frac{1}{n} \mathbb{E}_{\mathcal{S}_t}\left[\sum_{i=1}^n \Delta_t^{w_i}\right] = \mathbb{E}_{\mathcal{S}_t}[\Delta_t^{w_1}] = \frac{1}{m} \sum_{i=1}^m \Delta_t^i = \bar{\Delta}_t.$$

The proof of FedCAMS in partial participation settings has a similar outline combing the proof of partial participation in FedAMS and full participation in FedCAMS. By Assumption 3.1, then

$$\mathbb{E}[f(\mathbf{z}_{t+1})] - f(\mathbf{z}_{t}) \\
\leq \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\rangle \right] \\
-\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \eta \frac{\beta_{1}}{1 - \beta_{1}} \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) m_{t-1}' + \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) \Gamma_{t} \right\rangle \right] \\
-\frac{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \eta \frac{\beta_{1}}{1 - \beta_{1}} \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) m_{t-1}' - \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) \Gamma_{t} \right\|^{2}\right] \\
+ \underbrace{\frac{\eta^{2} L}{2} \mathbb{E}\left[\left\| \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} - \frac{\beta_{1}}{1 - \beta_{1}} \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) m_{t-1}' - \left( \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right) \Gamma_{t} \right\|^{2}\right]}_{T_{4}'}$$

Note that the bound for  $T'_2$  is exactly the same as the bound for  $T_2$ . For the three corresponding terms,  $T'_1, T'_3$  and  $T'_4$  which include the second-order momentum estimate of  $\Delta_t$ . For  $T'_1$ , similar to the full participation settings,

$$T_{1}^{\prime} \leq \sqrt{2}\mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \eta \frac{\Delta_{t}}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \right\rangle\right] + \sqrt{2}\eta\mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \frac{\Delta_{t}}{\sqrt{\mathbf{v}_{t} + \epsilon}} - \frac{\Delta_{t}}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \right\rangle\right].$$
(D.16)

The first term in (D.16) does not change in partial participation scheme. The second term is changed due to the variance of  $\Delta_t$  changes. For the second term of  $T'_1$ , then

$$\sqrt{2}\eta \mathbb{E}\left[\left\langle \nabla f(\theta_t), \frac{\Delta_t}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{\Delta_t}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\rangle\right] \le \frac{\sqrt{2(1 - \beta_2)}\eta G}{\epsilon} \mathbb{E}[\|\Delta_t\|^2].$$

2318 For  $T'_3$ , similar to the proof of  $T_3$ , we have

$$\sum_{t=1}^{T} T_3' \leq \frac{\eta^2 L}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\Delta_t\|^2] + \eta^2 L C_1^2 \eta_l^2 K^2 G^2 \sum_{t=1}^{T} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_t^{-1/2}\right\|^2\right],$$

where  $C_1 = \frac{\beta_1}{1-\beta_1} + \frac{m}{n} \sqrt{\frac{12q^2}{(1-q^2)^2} + \frac{(1-q^2)^2 C^2}{\alpha^2 n^2 q^2}}$  in partial participation,  $T_4' = \eta \mathbb{E}\left[\left\langle f(\mathbf{z}_t) - f(\theta_t), \widehat{\mathbf{V}}_t^{-1/2} \Delta_t \right\rangle\right]$  $\leq \eta \mathbb{E} \Big[ \|f(\mathbf{z}_t) - f(\theta_t)\| \Big\| \widehat{\mathbf{V}}_t^{-1/2} \Delta_t \Big\| \Big]$  $\leq \eta^{2} L \mathbb{E} \left[ \left\| \frac{\beta_{1}}{1-\beta_{1}} \widehat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}_{t-1}' + \widehat{\mathbf{V}}_{t-1}^{-1/2} \Gamma_{t} \right\| \left\| \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\| \right]$  $\leq \frac{C_1 \eta^2 \eta_l^2 K^2 L G^2}{C_1 \eta^2 \eta_l^2 K^2 L G^2}.$ 

The summation from  $T'_1$  to  $T'_4$  over total iteration T is: 

 $\mathbb{E}[f(\mathbf{z}_{T+1})] - f(\mathbf{z}_1) = \sum^T [T_1' + T_2' + T_3' + T_4']$ 

 $+ \frac{\eta^2 L}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\Delta_t\|^2] + \frac{C_1 T \eta^2 \eta_l^2 K^2 L G^2}{\epsilon}$ 

The proof outline is similar with previous proof. We take the use of Lemma F.3.F.9.F.13 for corresponding terms. By additional constraints of local learning rate  $\eta_n$  with the inequality  $\left[\eta^2 L + \sqrt{2(1-\beta_2)}\eta G\right] \frac{\eta_l^2(n-1)}{mn(m-1)\epsilon} - \frac{\eta\eta_l}{2Km^2} \left[\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2}}\eta_l^2 K^2 G^2 + \epsilon\right]^{-1} \leq 0$ , we obtain the constraint  $\eta_l \leq \frac{n(m-1)}{m(n-1)} \frac{\epsilon}{2K\sqrt{4\beta_2(1+q^2)^3(1-q^2)^{-2}K^2 G^2 + \epsilon}[\eta L + \sqrt{2(1-\beta_2)G}]}$ , and we further need  $\eta_l$  satis-

 $\leq -\frac{\eta\eta_l K}{4} \sum_{i=1}^T \mathbb{E} \left[ \left\| \frac{\nabla f(\theta_t)}{\sqrt[4]{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\|^2 \right] + \frac{5\eta\eta_l^3 K^2 L^2 T}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_g^2) + \frac{\sqrt{2(1-\beta_2)}\eta G}{\epsilon} \sum_{t=1}^T \mathbb{E}[\|\Delta_t\|^2]$ 

 $-\frac{\eta\eta_l}{2Km^2}\sum_{i=1}^T \mathbb{E}\left[\left\|\frac{1}{\sqrt[4]{\beta_2 \mathbf{v}_{t-1} + \epsilon}}\sum_{i=1}^m\sum_{l=0}^{K-1}\nabla F_i(\theta_l)\right\|^2\right] + C_1\eta\eta_l KG^2\sum_{i=1}^T \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_t^{-1/2}\right\|_1\right]$ 

 $+ C_{1}^{2} \eta^{2} \eta_{l}^{2} K^{2} L G^{2} \epsilon^{-1/2} \sum_{i=1}^{T} \mathbb{E} \left[ \left\| \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right\|_{1} \right] + C_{1}^{2} \eta^{2} \eta_{l}^{2} K^{2} L G^{2} \sum_{i=1}^{T} \mathbb{E} \left[ \left\| \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right\|^{2} \right]$ 

 $\leq -\frac{\eta\eta_{l}K}{4\sqrt{4\beta_{2}\frac{(1+q^{2})^{3}}{(1-q^{2})^{2}}\eta_{l}^{2}K^{2}G^{2}+\epsilon}}\sum_{t=1}^{T}\mathbb{E}[\|\nabla f(\theta_{t})\|^{2}] + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}T}{\sqrt{2\epsilon}}(\sigma_{l}^{2}+6K\sigma_{g}^{2}) + \frac{C_{1}\eta\eta_{h}KG^{2}d}{T\sqrt{\epsilon}}$ 

 $+\frac{2C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}LG^{2}d}{T\epsilon}-\frac{\eta\eta_{l}}{2\sqrt{4\beta_{2}\frac{(1+q^{2})^{3}}{(1-2)^{2}}\eta_{l}^{2}K^{2}G^{2}+\epsilon Km^{2}}}\sum_{t=1}^{T}\mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t})\right)\right\|^{2}\right]$ 

 $+ \left(\frac{\eta^2 \eta_l^2 L K T}{n \epsilon} + \frac{\sqrt{2(1-\beta_2)} \eta \eta_l^2 K T G}{n \epsilon}\right) \sigma_l^2 + \frac{C_1 T \eta^2 \eta_l^2 K^2 L G^2}{\epsilon}$ 

+  $(90mK^4L^2\eta_l^2 + 3mK^2)\sum_{i=1}^{T}\mathbb{E}[\|\nabla f(\theta_t)\|^2] + 3mK^2T\sigma_g^2]$ 

 $+\left(\frac{\eta^2\eta_l^2L}{\epsilon}+\frac{\sqrt{2(1-\beta_2)}\eta\eta_l^2G}{\epsilon}\right)\frac{m-n}{mn(m-1)}\left[15mK^3L^3\eta_l^2(\sigma_l^2+6K\sigma_g^2)T\right]$ 

 $+\left(\eta^2 \eta_l^2 L + \sqrt{2(1-\beta_2)}\eta \eta_l^2 G\right) \frac{n-1}{mn(m-1)} \sum_{l=1}^T \mathbb{E}\left[\left\|\sum_{l=1}^m \sum_{k=1}^{K-1} \nabla F_i(\theta_t)\right\|^2\right].$ 

$$\begin{aligned} & \text{fies} \ \frac{\eta\eta_{l}K}{4\sqrt{4\beta_{2}(1+q^{2})^{3}(1-q^{2})^{-2}\eta_{l}^{2}K^{2}G^{2}+\epsilon}} - (\eta^{2}L + \sqrt{2(1-\beta_{2})}\eta G) \frac{\eta_{l}^{2}(m-n)}{mn(m-1)\epsilon} (90mK^{4}L^{2}\eta_{l}^{2} + 3mK^{2}) \geq \\ & \frac{\eta\eta_{l}K}{8\sqrt{4\beta_{2}(1+q^{2})^{3}(1-q^{2})^{-2}\eta_{l}^{2}K^{2}G^{2}+\epsilon}} \text{. Hence for the convergence rate, then} \end{aligned}$$

Therefore

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 8\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 + \epsilon} \left[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\right]$$

The proof of Theorem A.4 is similar to the above proof procedure and the detailed proof will not be given here. 

D.4 PROOF OF COROLLARY A.2 

If choose  $\eta_l = \Theta(\frac{1}{\sqrt{TK}})$  and  $\eta = \Theta(\sqrt{Kn})$ , we get  $\min_{t \in [T]} \mathbb{E}[\|\nabla f(\theta_t)\|^2] = \mathcal{O}(\frac{\sqrt{K}}{\sqrt{Tn}})$ . 

## E **PROOF OF THEOREMS IN SECTION 4.3 AND PARTIAL PARTICIPATION** SETTING FOR FEDBNLACA, FEDBACA

E.1 PROOF OF THEOREM 4.5 

Notations and equations: From the update rule of Algorithm 2, we get  $\mathbf{e}_1 = 0, \mathbf{e}_t = \frac{1}{m} \sum_{i=1}^{m} \mathbf{e}_t^i$ and  $\mathbf{m}_t = (1 - \beta_1) \sum_{i=1}^t \tilde{\beta}_1^{t-i} \widehat{\widehat{\Delta}}_i$ . Denote a global uncompressed difference  $\Delta_t = \frac{1}{m} \sum_{i=1}^m \Delta_t^i$ . Denote a virtual momentum sequence:  $\mathbf{m}'_t = \beta_1 \mathbf{m}'_{t-1} + (1 - \beta_1) \Delta_t$ , hence we have  $\mathbf{m}'_t = \mathbf{m}_t^{t-1} \mathbf{m}_t^{t-1} \mathbf{m}_t^{t-1}$ .  $(1-\beta_1)\sum_{i=1}^t \beta_1^{t-i}\Delta_i$ . By the aforementioned definition and notation, then

$$\widehat{\Delta}_t - \Delta_t = \frac{1}{m} \sum_{i=1}^m (\widehat{\widehat{\Delta}}_t^i - \Delta_t^i) = \frac{1}{m} \sum_{i=1}^m (\widehat{\Delta}_t^i - \Delta_t^i) - \frac{1}{M_t} \sum_{i \in M_t} \mathcal{C}(q_t^i)$$

$$= \frac{1}{m} \sum_{i=1}^{m} (\mathbf{e}_{t}^{i} - \mathbf{e}_{t+1}^{i}) - \frac{1}{M_{t}} \sum_{i \in M_{t}} \mathcal{C}(q_{t}^{i}) = \mathbf{e}_{t} - \mathbf{e}_{t+1} - \frac{1}{M_{t}} \sum_{i \in M_{t}} \mathcal{C}(q_{t}^{i}).$$
(E.1)

Denote the weighted averaging error sequence  $\Gamma_t = (1 - \beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \mathbf{e}_r$ , with the input  $\mathbf{e}_1 = 0$ , we obtain the relation between  $\Gamma_t$  and  $m_t$  as follows

2431  
2432  
2433  
2434
$$\mathbf{m}_{t} - \mathbf{m}_{t}' = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\mathbf{e}_{\tau} - \mathbf{e}_{\tau+1}) - (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} \frac{1}{M_{t}} \sum_{i \in M_{t}} C(q_{t}^{i})$$
2434
$$t = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\widehat{\Delta}_{\tau} - \Delta_{\tau}) = (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} (\mathbf{e}_{\tau} - \mathbf{e}_{\tau+1}) - (1 - \beta_{1}) \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} \frac{1}{M_{t}} \sum_{i \in M_{t}} C(q_{t}^{i})$$

$$= \Gamma_t - \Gamma_{t+1} - (1 - \beta_1) \sum_{\tau=1}^{\circ} \beta_1^{t-\tau} \frac{1}{M_t} \sum_{i \in M_t} \mathcal{C}(q_t^i),$$
(E.2)

where the last step holds due to  $\Gamma_{t+1} = (1 - \beta_1) \sum_{\tau=1}^{t+1} \beta_1^{t-\tau} \mathbf{e}_{\tau+1} = (1 - \beta_1) \sum_{\tau=1}^{t} \beta_1^{t-\tau} \mathbf{e}_{\tau+1} + \beta_1^t \mathbf{e}_1$ . Similar to previous works studied adaptive methods, we introduce a Lyapunov sequence  $z_t$ :assume  $\theta_0 = \theta_1$ , for each  $t \ge 1$ ,

2441  
2442 
$$\mathbf{z}_t = \theta_t + \frac{\beta_1}{1-\beta_1}(\theta_t - \theta_{t-1}) = \frac{1}{1-\beta_1}\theta_t - \frac{\beta_1}{1-\beta_1}\theta_{t-1}$$

Therefore, by the update rule of  $\theta_t$ , 

$$\mathbf{y}_{t+1} = \theta_{t+1} + \eta \frac{\beta_1}{1 - \beta_1} \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t - \eta \frac{\beta_1}{1 - \beta_1} \hat{\mathbf{V}}_t^{-1/2} (\hat{\hat{\theta}}_t - \theta_t) + (1 - \beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \frac{1}{M_t} \sum_{i \in M_t} \mathcal{C}(q_t^i)$$

 $=\theta_{t+1} + \eta \frac{\beta_1}{1-\beta_1} \widehat{\mathbf{V}}_t^{-1/2} [\mathbf{m}_t' + \Gamma_t - \Gamma_{t+1}] - (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \frac{1}{M_t} \sum_{i \in M_t} \mathcal{C}(q_t^i) + \eta \frac{\beta_1}{1-\beta_1} \widehat{\mathbf{V}}_t^{-1/2} (\widehat{\widehat{\theta}}_t - \theta_t)$ 

$$=\theta_{t+1} + \eta \frac{\beta_1}{1-\beta_1} \widehat{\mathbf{V}}_t^{-1/2} \mathbf{m}'_t + \eta \frac{\beta_1}{1-\beta_1} \widehat{\mathbf{V}}_t^{-1/2} \left[ \frac{\Gamma_{t+1} - (1-\beta_1)\mathbf{e}_{t+1}}{\beta_1} - \Gamma_{t+1} \right]$$

$$-(1-\beta_1)\sum_{\tau=1}^{1}\beta_1^{t-\tau}\frac{1}{M_t}\sum_{i\in M_t}\mathcal{C}(q_t^i) - \eta\frac{\beta_1}{1-\beta_1}\hat{\mathbf{V}}_t^{-1/2}(\widehat{\hat{\theta}}_t - \theta_t)$$

$$=\theta_{t+1} + \eta \frac{\beta_1}{1-\beta_1} \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}'_t + \eta \hat{\mathbf{V}}_t^{-1/2} \mathbf{\Gamma}_{t+1} - \eta \hat{\mathbf{V}}_t^{-1/2} \mathbf{e}_{t+1} - (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \frac{1}{M_t} \sum_{i \in M_t} \mathcal{C}(q_t^i) + \eta \frac{\beta_1}{1-\beta_1} \hat{\mathbf{V}}_t^{-1/2} (\widehat{\theta}_t - \theta_t).$$
(E.3)

> The third equation holds due to the fact that  $\Gamma_{t+1} = \beta_1 \Gamma_t + (1 - \beta_1) \mathbf{e}_{t+1}$ . We then introduce a new sequence based on the previous Lyapunov sequence  $y_t$  as follows

$$\mathbf{z}_{t+1} = \mathbf{y}_{t+1} + (1 - \beta_1) \sum_{\tau=1}^{t} \beta_1^{t-\tau} \frac{1}{M_t} \sum_{i \in M_t} \mathcal{C}(q_t^i) + \eta \hat{\mathbf{V}}_t^{-1/2} \mathbf{e}_{t+1}$$

$$=\theta_{t+1} + \eta \frac{\beta_1}{1-\beta_1} \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t' + \eta \hat{\mathbf{V}}_t^{-1/2} \Gamma_{t+1} + \eta \frac{\beta_1}{1-\beta_1} \hat{\mathbf{V}}_t^{-1/2} (\hat{\theta}_t - \theta_t).$$
(E.4)

The sequence difference  $\mathbf{z}_{t+1} - \mathbf{z}_t$  can be represented by 

$$\mathbf{z}_{t+1} - \mathbf{z}_t = \theta_{t+1} - \theta_t + \eta \frac{\beta_1}{1 - \beta_1} \widehat{\mathbf{V}}_t^{-1/2} \mathbf{m}'_t - \eta \frac{\beta_1}{1 - \beta_1} \widehat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}'_{t-1}$$

$$+\eta \widehat{\mathbf{V}}_t^{-1/2} \Gamma_{t+1} - \eta \widehat{\mathbf{V}}_{t-1}^{-1/2} \Gamma_t + \eta \frac{1}{1-\beta_1} \widehat{\mathbf{V}}_t^{-1/2} (\widehat{\widehat{\theta}}_t - \theta_t)$$

$$= \eta \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t + \eta \hat{\mathbf{V}}_t^{-1/2} \Gamma_{t+1} + \eta \frac{\beta_1}{1-\beta_t} \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t'$$

2479 
$$\beta_1 = \frac{1}{2}$$

2480  
2481  
2482 
$$-\eta \frac{\beta_1}{1-\beta_t} \hat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}'_{t-1} - \eta \hat{\mathbf{V}}_{t-1}^{-1/2} \Gamma_t + \eta \frac{1}{1-\beta_1} \hat{\mathbf{V}}_t^{-1/2} (\hat{\hat{\theta}}_t - \theta_t),$$

where the second equation follows the update rule of  $\theta_{t+1}$ . Following (E.2), then combining likely terms and applying the definition of  $m'_t$ , we have

 $\mathbf{z}_{t+1} - \mathbf{z}_t = \eta \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t' + \eta \hat{\mathbf{V}}_t^{-1/2} \Gamma_t + \eta \frac{\beta_1}{1 - \beta_1} \hat{\mathbf{V}}_t^{-1/2} \mathbf{m}_t'$ 

$$\begin{split} &-\eta \frac{\beta_{1}}{1-\beta_{1}} \hat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}_{t-1}' - \eta \hat{\mathbf{V}}_{t-1}^{-1/2} \Gamma_{t} + \eta \frac{1}{1-\beta_{1}} \hat{\mathbf{V}}_{t}^{-1/2} (\widehat{\widehat{\theta}_{t}} - \theta_{t}) \\ &= \eta \frac{1}{1-\beta_{1}} \hat{\mathbf{V}}_{t}^{-1/2} \mathbf{m}_{t}' - \eta \frac{\beta_{1}}{1-\beta_{1}} \hat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}_{t-1}' + \eta \hat{\mathbf{V}}_{t}^{-1/2} \Gamma_{t} + \eta \hat{\mathbf{V}}_{t-1}^{-1/2} \Gamma + \eta \frac{1}{1-\beta_{1}} \hat{\mathbf{V}}_{t}^{-1/2} (\widehat{\widehat{\theta}_{t}} - \theta_{t}), \\ &= \eta \frac{1}{1-\beta_{1}} \hat{\mathbf{V}}_{t}^{-1/2} [\beta_{1} \mathbf{m}_{t-1}' + (1-\beta_{1}) \Delta_{t}] - \eta \frac{\beta_{1}}{1-\beta_{1}} \hat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}_{t-1}' \\ &+ \eta \hat{\mathbf{V}}_{t}^{-1/2} \Gamma_{t} - \eta \hat{\mathbf{V}}_{t-1}^{-1/2} \Gamma_{t} + \eta \frac{1}{1-\beta_{1}} \hat{\mathbf{V}}_{t}^{-1/2} (\widehat{\widehat{\theta}_{t}} - \theta_{t}) \\ &= \eta \hat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} - \eta \frac{\beta_{1}}{1-\beta_{1}} \left( \hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2} \right) \mathbf{m}_{t-1}' - \eta \left( \hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2} \right) \Gamma_{t} + \eta \frac{1}{1-\beta_{1}} \hat{\mathbf{V}}_{t}^{-1/2} (\widehat{\widehat{\theta}_{t}} - \theta_{t}) \end{split}$$

Therefore, we obtain a helpful Lyapunov sequence for our proof of FedCAMS. The proof of Fed-CAMS in full participation settings has a similar outline with the proof of FedAMS. By Assumption 3.1,

 $-\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \eta \frac{\beta_{1}}{1-\beta_{t}} \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \mathbf{m}_{t-1}' + \eta \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \mathbf{\Gamma}_{t}\right\rangle\right]$ 

 $+\frac{\eta^{2}L}{2}\mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t}^{-1/2}\Delta_{t}-\frac{\beta_{1}}{1-\beta_{1}}\left(\widehat{\mathbf{V}}_{t-1}^{-1/2}-\widehat{\mathbf{V}}_{t}^{-1/2}\right)\mathbf{m}_{t-1}'-\left(\widehat{\mathbf{V}}_{t-1}^{-1/2}-\widehat{\mathbf{V}}_{t}^{-1/2}\right)\Gamma_{t}\right\|^{2}\right]$ 

 $+\underbrace{\frac{\eta^{2}L}{2}\mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t}^{-1/2}\Delta_{t}-\frac{\beta_{1}}{1-\beta_{1}}\left(\widehat{\mathbf{V}}_{t-1}^{-1/2}-\widehat{\mathbf{V}}_{t}^{-1/2}\right)\mathbf{m}_{t-1}^{\prime}-\left(\widehat{\mathbf{V}}_{t-1}^{-1/2}-\widehat{\mathbf{V}}_{t}^{-1/2}\right)\Gamma_{t}\right\|^{2}\right]}_{\mathbf{V}_{t}}$ 

 $=\underbrace{\mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\rangle\right]}_{T_{t}} \underbrace{-\eta \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\beta_{1}}{1-\beta_{1}} \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \mathbf{m}_{t-1}' + \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \Gamma_{t} \right\rangle\right]}_{T_{t}}$ 

 $+\underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\rangle\right]}_{T_{4}} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t}) \right\rangle\right]}_{T} + \underbrace{\frac{\eta^{2} L^{2}}{(1 - \beta)^{2}} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\|^{2}\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t}) \right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t})\right\rangle\right]}_{T} + \underbrace{\mathbb{E}$ 

(E.5)

 $\mathrm{E}[f(\mathbf{z}_{t+1})] - f(\mathbf{z}_t)$ 

 $\leq \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_t), \eta \widehat{\mathbf{V}}_t^{-1/2} \Delta_t \right\rangle \right]$ 

 $\leq \mathbb{E}[\langle \nabla f(\mathbf{z}_t), \mathbf{z}_{t+1} - \mathbf{z}_t \rangle] + \frac{L}{2} \mathbb{E}[\|\mathbf{z}_{t+1} - \mathbf{z}_t\|^2]$ 

here recall the notation  $\hat{\mathbf{V}}_t = \text{diag}(\hat{\mathbf{v}}_t) = \text{diag}(\max(\hat{\mathbf{v}}_{t-1}, \mathbf{v}_t, \epsilon)).$ Bounding  $T_1$ :,

$$\begin{split} \mathbf{T} &1 = \mathbb{E}\left[\left\langle \nabla f(\theta_t), \eta \frac{\Delta_t}{\sqrt{\hat{\mathbf{v}}_t}} \right\rangle\right] \\ &\leq \eta \mathbb{E}\left[\left\langle \nabla f(\theta_t), \frac{\sqrt{2} \cdot \Delta_t}{\sqrt{\mathbf{v}_t + \epsilon}} \right\rangle\right] \\ &= \sqrt{2}\eta \mathbb{E}\left[\left\langle \nabla f(\theta_t), \frac{\Delta_t}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\rangle\right] + \sqrt{2}\eta \mathbb{E}\left[\left\langle \nabla f(\theta_t), \frac{\Delta_t}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{\Delta_t}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\rangle\right], \end{split}$$
(E.6)

where the first inequality follows by the fact that  $\hat{v}_t \geq \frac{v_t + \epsilon}{2}$ . For the second term in (E.6),

$$\begin{split} \sqrt{2} \cdot \eta \mathbb{E} \left[ \left\langle \nabla f(\theta_t), \frac{\Delta_t}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{\Delta_t}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\rangle \right] \\ &\leq \sqrt{2} \cdot \eta \cdot \mathbb{E} \| \nabla f(\theta_t) \| \mathbb{E} \left[ \left\| \frac{1}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{1}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\| \cdot \| \Delta_t \| \right] \\ &\leq \frac{\eta \sqrt{2(1 - \beta_2)} G}{\epsilon} \mathbb{E} [\| \Delta_t \|^2], \end{split}$$

where the second inequality follows from Lemma F.1 and F.5, and we will further apply the bound for  $E[||\Delta_t||^2]$  by applying Lemma F.7. For the first term in (E.6),

$$\begin{aligned}
2573 \\
2574 \\
2575 \\
2576 \\
2576 \\
2577 \\
2576 \\
2577 \\
2578 \\
2579 \\
2580 \\
2580 \\
2581 \\
2581 \\
2582 \\
2582 \\
2582 \\
2583 \\
2582 \\
2583 \\
2584 \\
2584 \\
2584 \\
2584 \\
2584 \\
2584 \\
2585 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2587 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2587 \\
2587 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2587 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586 \\
2586$$

For the last term in (E.7), we get

 $\sqrt{2\eta} \left\langle \frac{\nabla f(\theta_t)}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}}, \mathbb{E} \left[ -\frac{\eta_l}{m} \sum_{i=1}^m \sum_{k=1}^{K-1} \mathbf{g}_{t,k}^i + \frac{\eta_l K}{m} \sum_{i=1}^m \nabla F_i(\theta_t) \right] \right\rangle$ 

 $-\frac{\sqrt{2}\eta\eta_l}{2Km^2}\mathbb{E}\left[\left\|\frac{1}{\sqrt[4]{\beta_2\mathbf{v}_{t-1}+\epsilon}}\sum_{i=1}^m\sum_{k=0}^{K-1}\nabla F_i(\theta_{t,k}^i)\right\|^2\right]$ 

 $-\frac{\sqrt{2\eta\eta_l}}{2Km^2}\mathbb{E}\bigg[\bigg\|\frac{1}{\sqrt[4]{\beta_2}\mathbf{v}_{t-1}+\epsilon}\sum_{i=1}^m\sum_{k=1}^{K-1}\nabla F_i(\theta_{t,k}^i)\bigg\|^2\bigg],$ 

 $=\sqrt{2}\eta \left\langle \frac{\sqrt{\eta K}}{\sqrt[4]{\beta_2 \mathbf{y}_{t-1}} + \epsilon} \nabla f(\theta_t), -\frac{\sqrt{\eta_n K}}{Km} \frac{1}{\sqrt[4]{\beta_2 \mathbf{y}_{t-1}} + \epsilon} \mathbb{E} \left[ \sum_{i=1}^m \sum_{k=1}^{K-1} (\nabla F_i(\theta_{t,k}^i) - \nabla F_i(\theta_t)) \right] \right\rangle$ 

 $=\frac{\sqrt{2}\eta\eta_{l}K}{2}\left\|\frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}}+\epsilon}\right\|^{2}+\frac{\sqrt{2}\eta\eta_{l}}{2Km^{2}}\mathbb{E}\left\|\frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1}}+\epsilon}\sum_{i=1}^{m}\sum_{j=1}^{K-1}(\nabla F_{i}(\theta_{t,k}^{i})-\nabla F_{i}(\theta_{t}))\right\|^{2}$ 

where the second equation follows from  $\langle \mathbf{x}, \mathbf{y} \rangle = \frac{1}{2} [\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2]$ , and the inequality

 $\leq \frac{\sqrt{2}\eta\eta_l K}{2} \left\| \frac{\nabla f(\theta_t)}{\sqrt[4]{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\|^2 + \frac{\sqrt{2}\eta\eta_l L^2}{2m} \sum_{i=1}^m \sum_{k=0}^{K-1} \mathbb{E} \left[ \left\| \frac{\theta_{t,k}^i - \theta_t}{\sqrt[4]{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\|^2 \right]$ 

 $\sqrt{2}\eta \left\langle \frac{\nabla f(\theta_t)}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}}, \mathbb{E} \left[ -\frac{\eta_t}{m} \sum_{i=1}^m \sum_{k=1}^{K-1} \mathbf{g}_{t,k}^i + \frac{\eta_l K}{m} \sum_{i=1}^m \nabla F_i(\theta_t) \right] \right\rangle$ 

 $\leq \frac{\sqrt{2}\eta\eta_l K}{2} \left\| \frac{\nabla f(\theta_t)}{\sqrt[4]{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\|^2 + \frac{\sqrt{2}\eta\eta_l}{2m} \cdot \sum_{i=1}^m \sum_{k=0}^{K-1} \mathbb{E} \left[ \left\| \frac{\nabla F_i(\theta_{t,k}^i) - \nabla F_i(\theta_t)}{\sqrt[4]{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\|^2 \right]$ 

holds by applying Cauchy-Schwarz inequality. Then by Assumption 3.1, then

$$\leq \frac{3\sqrt{2}\eta\eta_{l}K}{4} \left\| \frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \right\|^{2} + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}}{\sqrt{2\epsilon}}(\sigma_{l}^{2} + 6K\sigma_{g}^{2})$$
$$- \frac{\sqrt{2}\eta\eta_{l}}{2Km^{2}}\mathbb{E}\Big[ \left\| \frac{1}{\sqrt[4]{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \Big],$$

 $-\frac{\sqrt{2}\eta\eta_l}{2Km^2}\mathbb{E}\left[\left\|\frac{1}{\frac{4}{\sqrt{2}\pi}}\sum_{k=1}^m\sum_{k=1}^{K-1}\nabla F_i(\theta_{t,k}^i)\right\|^2\right]$ 

where the last inequality holds by applying Lemma F.13 and the constraint of local learning rate  $\eta_n \leq \frac{1}{8KL}$ . Then

$$-\frac{\eta\eta_l}{2Km^2}\mathbb{E}\Big[\left\|\frac{1}{\sqrt[4]{\beta_2\mathbf{v}_{t-1}+\epsilon}}\sum_{i=1}^m\sum_{k=0}^{K-1}\nabla F_i(\theta_{t,k}^i)\right\|^2\Big]+\frac{\eta\sqrt{2(1-\beta_2)}G}{\epsilon}\mathbb{E}[\|\Delta_t\|^2].$$
 (E.8)

 $-\frac{\sqrt{2}\cdot\eta\eta_l}{2Km^2}\mathbb{E}\left[\left\|\frac{1}{\sqrt[4]{\beta_2\mathbf{v}_{t-1}+\epsilon}}\sum_{i=1}^m\sum_{k=0}^{K-1}\nabla F_i(\theta_{t,k}^i)\right\|^2\right]+\frac{\eta\sqrt{2(1-\beta_2)}G}{\epsilon}\mathbb{E}[\|\Delta_t\|^2]$ 

**Bounding**  $T_2$ : The bound for  $T_2$  mainly follows by the update rule and definition of virtual sequence  $\mathbf{z}_t$ .

 $T_1 \leq -\frac{\sqrt{2} \cdot \eta \eta_l K}{4} \mathbb{E} \left[ \left\| \frac{\nabla f(\theta_t)}{\sqrt[4]{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\|^2 \right] + \frac{5\eta \eta_l^3 K^2 L^2}{\sqrt{2\epsilon}} \left( \sigma_l^2 + 6K \sigma_g^2 \right)$ 

 $\leq -\frac{\eta\eta_l K}{4} \mathbb{E}\left[ \left\| \frac{\nabla f(\theta_t)}{\sqrt[4]{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\|^2 \right] + \frac{5\eta\eta_l^3 K^2 L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K\sigma_g^2)$ 

$$T_{2} = -\eta \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\beta_{1}}{1-\beta_{1}} \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \mathbf{m}_{t-1}' + \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \Gamma_{t}\right\rangle\right]$$

$$= \eta \mathbb{E}\left[\left\langle -\nabla f(\theta_{t}) + \nabla f(\theta_{t}) - \nabla f(\mathbf{z}_{t}), \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \left(\frac{\beta_{1}}{1-\beta_{1}}\mathbf{m}_{t-1}' + \Gamma_{t}\right)\right\rangle\right]$$

$$\leq \eta \mathbb{E}\left[\|\nabla f(\theta_{t})\| \left\| \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \left(\frac{\beta_{1}}{1-\beta_{1}}\mathbf{m}_{t-1}' + \Gamma_{t}\right)\right\|\right]$$

$$+ \eta^{2} L \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t-1}^{-1/2} \left(\frac{\beta_{1}}{1-\beta_{1}}\mathbf{m}_{t-1}' + \Gamma_{t}\right)\right\| \left\| \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \left(\frac{\beta_{1}}{1-\beta_{1}}\mathbf{m}_{t-1}' + \Gamma_{t}\right)\right\|\right]$$

$$\leq \eta C_{1} \eta_{l} K G^{2} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right\|_{1}\right] + \eta^{2} C_{1}^{2} L \eta_{l}^{2} K^{2} G^{2} \epsilon^{-1/2} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right\|_{1}\right], \qquad (E.9)$$

where the last inequality holds by Lemma C.4, here  $C_1 = \frac{\beta_1}{1-\beta_1} + \sqrt{\frac{12q^2}{(1-q^2)^2} + \frac{(1-q^2)^2C^2}{\alpha^2m^2q^2}}$ . **Bounding** $T_3$ : It can be bounded as follows:

$$T_{3} = \frac{\eta^{2}L}{2} \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} + \frac{\beta_{1}}{1 - \beta_{1}} \left( \hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2} \right) \mathbf{m}_{t-1}' + \left( \hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2} \right) \Gamma_{t} \right\|^{2} \right]$$

$$\leq \eta^{2} L \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\|^{2} \right] + \eta^{2} L \mathbb{E} \left[ \left\| \frac{\beta_{1}}{1 - \beta_{1}} \left( \hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2} \right) \mathbf{m}_{t-1}' + \left( \hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2} \right) \Gamma_{t} \right\|^{2} \right]$$

$$\leq \eta^{2} L \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\|^{2} \right] + \eta^{2} L C_{1}^{2} \eta_{t}^{2} K^{2} G^{2} \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2} \right\|^{2} \right], \qquad (E.10)$$

where the first inequality follows by Cauchy-Schwarz inequality, and the second one follows by Lemma C.4, here  $C_1 = \frac{\beta_1}{1-\beta_1} + \sqrt{\frac{12q^2}{(1-q^2)^2} + \frac{(1-q^2)^2C^2}{\alpha^2m^2q^2}}$ . Bounding  $T_4$ :

$$T_{4} = \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\rangle\right]$$

$$\leq \mathbb{E}\left[\left\|\nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t})\right\| \left\|\eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t}\right\|\right]$$

$$\leq L\mathbb{E}\left[\left\|\mathbf{z}_{t} - \theta_{t}\right\| \left\|\eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t}\right\|\right]$$

$$\leq \frac{\eta^{2}L}{2} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t}\right\|^{2}\right] + \frac{\eta^{2}L}{2} \mathbb{E}\left[\left\|\frac{\beta_{1}}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}_{t-1}' + \widehat{\mathbf{V}}_{t-1}^{-1/2} \Gamma_{t}\right\|^{2}\right]$$

$$\|\mathbf{a}\| \|\mathbf{b}\| \leq \frac{1}{2} \|\mathbf{a}\|^2 + \frac{1}{2} \|\mathbf{b}\|^2. \text{ Then summing } T_4 \text{ over } t = 1, \cdots, T,$$
$$\sum_{t=1}^T T_4 \leq \frac{\eta^2 L}{2} \sum_{t=1}^T \mathbb{E} \left[ \left\| \widehat{\mathbf{V}}_t^{-1/2} \Delta_t \right\|^2 \right] + \frac{\eta^2 L}{2\epsilon} \sum_{t=1}^T \mathbb{E} \left[ \left\| \frac{\beta_1}{1 - \beta_1} \mathbf{m}_{t-1}' + \Gamma_t \right\|^2 \right]$$

$$\leq \frac{\eta^{2}L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\Delta_{t}\|^{2}] + \frac{\eta^{2}L}{\epsilon} \bigg[ \frac{\beta_{1}^{2}}{(1-\beta_{1})^{2}} \sum_{t=1}^{T} \mathbb{E}\|\mathbf{m}_{t-1}'\|^{2} + \sum_{t=1}^{T} \mathbb{E}\|\mathbf{\Gamma}_{t}\|^{2} \bigg].$$

where the first inequality holds by the fact of  $\langle \mathbf{a}, \mathbf{b} \rangle \leq \|\mathbf{a}\| \|\mathbf{b}\|$ , the second one follows from

Assumption 3.1 and the third one holds by the definition of virtual sequence  $z_t$  and the fact of

By Lemma F.11, then

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathbf{m}_{t-1}'\|^2] \le \frac{TK\eta_l^2}{m} \sigma_l^2 + \frac{\eta_l^2}{m^2} \sum_{t=1}^{T} \mathbb{E}\Big[\Big\|\sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i)\Big\|^2\Big],$$

and

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathbf{\Gamma}_t\|^2] \le \frac{4T(q+\gamma)^2}{(1-q^2)^2} \frac{K\eta_l^2}{m} \sigma_l^2 + \frac{\eta_l^2}{m^2} \frac{4(q+\gamma)^2}{(1-q^2)^2} \sum_{t=1}^{T} \mathbb{E}\left[\left\|\sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i)\right\|^2\right].$$

Therefore, the  $T_4$  term is bounded by

$$\sum_{t=1}^{T} T_4 \le \frac{\eta^2 L}{2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\Delta_t\|^2] + \frac{C_2 \eta^2 L}{\epsilon} \frac{\eta_l^2}{m^2} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i) \right\|^2 \right] + \frac{C_2 \eta^2 L}{\epsilon} \frac{T K \eta_l^2}{m} \sigma_l^2,$$
(E.11)

where 
$$C_2 = \frac{4(q+\gamma+\frac{\lambda C}{\alpha m})^2}{(1-q^2)^2} + \frac{\beta_1^2}{(1-\beta_1)^2}.$$

**Bound**  $T_5$  , there

$$\begin{aligned} & T_{5} = \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\rangle\right] = \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\rangle\right] \\ & = \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\rangle\right] + \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), -\nabla f(\theta_{t}), \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\rangle\right] \\ & = \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\nabla f(\theta_{t})\right\| \cdot \mathbb{E}\left\|\frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\| + \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t})\right\| \cdot \frac{\eta}{m} \sum_{i}^{m} \mathbb{E}\left\|\frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\| + \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t})\right\| \cdot \frac{\eta}{m} \sum_{i}^{m} \mathbb{E}\left\|\frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\| \\ & \leq \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\nabla f(\theta_{t})\right\| \cdot \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\| + L \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|z_{t} - \theta_{t}\right\| \cdot \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\| \\ & = \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\nabla f(\theta_{t})\right\| \cdot \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\| \\ & + L \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\nabla f(\theta_{t})\right\| \cdot \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\| \\ & + L \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} - \eta \frac{\beta_{1}}{1-\beta_{1}} \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) m_{t-1}' - \eta \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \Gamma_{t} + \eta \frac{1}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\| \\ & + L \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\| \\ & \cdot \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\| \\ & \leq \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\nabla f(\theta_{t})\right\| \cdot \mathbb{E}\left\|\frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\| + L \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t}\right\| \cdot \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\| \\ & \leq \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left\|\nabla f(\theta_{t})\right\| \cdot \mathbb{E}\left\|\frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{$$

$$\begin{aligned} & 2755 \\ & +L\frac{1}{m}\sum_{i}^{m} \mathbb{E} \left\| \eta \frac{\beta_{1}}{1-\beta_{1}} \left( \hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2} \right) \mathbf{m}_{t-1}' + \eta \left( \hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2} \right) \Gamma_{t} \right\| \cdot \frac{1}{m} \sum_{i}^{m} \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) \right\| \\ & +L\frac{1}{m} \sum_{i}^{m} \mathbb{E} \left\| \eta \frac{1}{1-\beta_{1}} \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) \right\| \cdot \frac{1}{m} \sum_{i}^{m} \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) \right\| \\ & +L\frac{1}{m} \sum_{i}^{m} \mathbb{E} \left\| \eta \frac{1}{1-\beta_{1}} \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) \right\| \\ & +L\frac{1}{m} \sum_{i}^{m} \mathbb{E} \left\| \theta_{t} \right\| \cdot \frac{1}{m} \sum_{i}^{m} \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) \right\| \\ & +L\frac{1}{m} \sum_{i}^{m} \mathbb{E} \left\| \theta_{t} \right\| \cdot \frac{1}{m} \sum_{i}^{m} \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) \right\| \\ & (E.12) \end{aligned}$$

$$& = \left[ G + \frac{L\eta\eta_{l}KG}{\sqrt{\epsilon}} + L\eta\eta_{l}C_{1}KG\mathbb{E} \| \hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2} \| \right] \cdot \frac{\eta(\gamma + \frac{C}{\alpha m})H}{(1-\beta)\sqrt{\epsilon}} + \frac{2L\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)^{2}\epsilon} + \frac{2L\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)\sqrt{\epsilon}} \\ & \text{here } C_{1} = \frac{\beta_{1}}{1-\beta_{1}} + \sqrt{\frac{12q^{2}}{(1-q^{2})^{2}} + \frac{(1-q^{2})^{2}C^{2}}{\alpha^{2}m^{2}q^{2}}}} \\ & \mathbf{Bound of } T_{6} \text{, here} \\ & T_{6} = \frac{\eta^{2}L^{2}}{(1-\beta)^{2}} \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t} - \theta_{t}) \right\|^{2} \right] = \frac{\eta^{2}L^{2}}{(1-\beta)^{2}} \frac{1}{m} \sum_{i}^{m} \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) - \hat{\mathbf{V}}_{t}^{-1/2} \frac{1}{m} \sum_{i}^{m} \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) - \hat{\mathbf{V}}_{t}^{-1/2} \frac{1}{m} \sum_{i}^{m} \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) - \hat{\mathbf{V}}_{t}^{-1/2} \frac{1}{m} \sum_{i}^{m} \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) - \hat{\mathbf{V}}_{t}^{-1/2} \frac{1}{m} \sum_{i}^{m} \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) - \hat{\mathbf{V}}_{t}^{-1/2} \frac{1}{m} \sum_{i}^{m} \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) - \hat{\mathbf{V}}_{t}^{-1/2} \frac{1}{m} \sum_{i}^{m} \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_{t}) - \hat{\mathbf{V}}_{t}^{-1/2} \frac{1}{m} \sum_{i}^{m} \mathbb{E} \left[ \left\| \hat{\mathbf{V}}_{t}^{-1/2} (\hat{\theta}_{t}^{-} - \theta_$$

$$T_{6} = \frac{\eta^{2}L^{2}}{(1-\beta)^{2}} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}}-\theta_{t})\right\|^{2}\right] = \frac{\eta^{2}L^{2}}{(1-\beta)^{2}} \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}}-\theta_{t})\right\|^{2}\right]$$
$$= \frac{\eta^{2}L^{2}}{(1-\beta)^{2}} \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}-\theta_{t}) - \widehat{\mathbf{V}}_{t}^{-1/2}\frac{1}{M_{t}} \sum_{i \in M_{t}} Q_{t}^{i}\right\|^{2}\right]$$
$$\leq \frac{2\eta^{2}L^{2}}{(1-\beta)^{2}} \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}-\theta_{t})\right\|^{2}\right] + \frac{2\eta^{2}L^{2}}{(1-\beta)^{2}} \frac{1}{m} \sum_{i}^{m} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t}^{-1/2}\frac{1}{M_{t}} \sum_{i \in M_{t}} Q_{t}^{i}\right\|^{2}\right]$$
$$\leq \frac{2L^{2}\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)^{2}\epsilon} (\frac{C^{2}}{\alpha^{2}m^{2}} + 1)$$
(E.14)

Merging pieces together: Substituting (E.8)), (E.9), (E.10), (E.12),, (E.14) and , (E.11) into (E.5), summing over from t = 1 to T, then 

2808 Hence by organizing and applying Lemmas, then 2809

2810

$$\begin{split} \mathbb{E}[f(\mathbf{z}_{T+1})] &- f(\mathbf{z}_{1}) \\ &\leq -\frac{\eta \eta_{l} K}{4} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2} \mathbf{v}_{t-1} + \epsilon}} \right\|^{2} \right] + \frac{5 \eta \eta_{l}^{3} K^{2} L^{2} T}{\sqrt{2\epsilon}} (\sigma_{l}^{2} + 6K \sigma_{g}^{2}) \\ &- \frac{\eta \eta_{l}}{2K m^{2}} \sum_{t=1}^{T} \mathbb{E} \left[ \left\| \frac{1}{\sqrt[4]{\beta_{2} \mathbf{v}_{t-1} + \epsilon}} \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i})) \right\|^{2} \right] + \frac{C_{1} \eta \eta_{l} K G^{2} d}{\sqrt{\epsilon}} + \frac{2C_{1}^{2} \eta^{2} \eta_{l}^{2} K^{2} L G^{2} d}{\epsilon} \\ &+ \left( \eta^{2} L + \frac{\eta^{2} L}{2} + \sqrt{2(1 - \beta_{2})} \eta G \right) \left[ \frac{K T \eta_{l}^{2}}{m \epsilon} \sigma_{l}^{2} + \frac{\eta_{l}^{2}}{m^{2} \epsilon} \sum_{t=1}^{T} \mathbb{E} \left[ \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \right] \right] \\ &+ \frac{\eta^{2} L}{\epsilon} \frac{\eta_{l}^{2} C_{2}}{m^{2}} \sum_{t=1}^{T} \mathbb{E} \left[ \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \right] + \frac{\eta^{2} L}{\epsilon} \frac{T K \eta_{l}^{2} C_{2}}{m} \sigma_{l}^{2} + \sum_{t=1}^{T} T_{5} + \sum_{t=1}^{T} T_{6}, \end{split}$$

by applying Lemma F.7 into all terms containing the second moment estimate of model difference  $\Delta_t$  in (E.15), using the fact that  $\left(\sqrt{\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} K^2 G^2 + \epsilon}\right)^{-1} \|\theta\| \leq 1$  $\left(\sqrt{\beta_2 \frac{(1+q^2)^3}{(1-q^2)}} \eta_l^2 K^2 G^2 + \epsilon\right)^{-1} \|\theta\| \le \|\frac{\theta}{\sqrt{\beta_2 \mathbf{v} + \epsilon}}\| \le \epsilon^{-1/2} \|\theta\|, \text{ and applying Lemma F.3 and F.13,}$ then

$$\begin{split} \mathbb{E}[f(\mathbf{z}_{T+1})] &- f(\mathbf{z}_{1}) \\ &\leq -\frac{m_{l}K}{4\sqrt{4\beta_{2}\frac{(1+q^{2})^{3}}{(1-q^{2})^{2}}\eta_{l}^{2}K^{2}G^{2} + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_{t})\|^{2}] + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}T}{\sqrt{2\epsilon}} (\sigma_{t}^{2} + 6K\sigma_{g}^{2}) \\ &\qquad + \frac{C_{1}\eta\eta_{l}KG^{2}d}{\sqrt{\epsilon}} + \frac{2C_{1}^{2}\eta_{l}^{2}\eta_{l}^{2}K^{2}G^{2}}{\epsilon} + \left(\frac{3\eta^{2}L}{2} + C_{2}\eta^{2}L + \sqrt{2(1-\beta_{2})}\eta_{G}}\right) \frac{KT\eta_{l}^{2}}{m\epsilon} \sigma_{l}^{2} + \frac{\eta_{l}\eta_{T}C^{2}K^{2}G^{2}}{\alpha^{2}m^{2}\epsilon} \\ &\qquad -\sum_{t=1}^{T} \mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1} \nabla F_{i}(\theta_{i,k}^{i})\right\|^{2}\right] \left[\frac{\eta_{l}}{2\sqrt{4\beta_{2}\frac{(1+q^{2})^{3}}{(1-q^{2})^{2}}\eta_{l}^{2}K^{2}G^{2} + \epsilon}Km^{2}} - \left(\frac{3\eta^{2}L}{2} + C_{2}\eta^{2}L + \sqrt{2(1-\beta_{2})}\eta_{G}}\right)\frac{\eta_{l}^{2}}{m^{2}\epsilon}\right] \\ &\qquad + \left[G + \frac{L\eta\eta_{l}KG}{\sqrt{\epsilon}} + L\eta\eta_{l}C_{1}KG\mathbb{E}\|\widehat{\nabla}_{t-1}^{-1/2} - \widehat{\nabla}_{t}^{-1/2}\|\right] \cdot \frac{T\eta(\gamma + \frac{C}{\alpha m})H}{(1-\beta)\sqrt{\epsilon}} \\ &\qquad + \frac{3TL\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)^{2}\epsilon} + \frac{2TL\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)\sqrt{\epsilon}} \\ &\leq -\frac{\eta\eta_{l}K}{4C_{0}}\sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_{t})\|^{2}] + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}T}{\sqrt{2\epsilon}}(\sigma_{t}^{2} + 6K\sigma_{g}^{2}) \\ &\qquad + \left[G + \frac{L\eta\eta_{l}KG^{2}}{\sqrt{\epsilon}} + \frac{2C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}G^{2}d}{\epsilon} + \left(\frac{3\eta^{2}L}{2} + C_{2}\eta^{2}L + \sqrt{2(1-\beta_{2})}\eta_{G}\right)\frac{KT\eta_{l}^{2}}{m\epsilon}\sigma_{l}^{2} \\ &\qquad + \frac{G\eta_{l}KG^{2}}{4C_{0}} + \frac{2C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}LG^{2}d}{(1-\beta)\sqrt{\epsilon}} + \frac{2TL\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)\sqrt{\epsilon}} \\ &\qquad + \left[G + \frac{L\eta\eta_{l}KG^{2}}{4C_{0}} + \frac{2C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}LG^{2}d}{\epsilon} + \left(\frac{3\eta^{2}L}{2} + C_{2}\eta^{2}L + \sqrt{2(1-\beta_{2})}\eta_{G}\right)\frac{KT\eta_{l}^{2}}{m\epsilon}\sigma_{l}^{2} \\ &\qquad + \left[G + \frac{L\eta\eta_{l}KG}{\sqrt{\epsilon}} + \frac{L\eta\eta_{l}C_{1}KGd}{\epsilon}\right] \cdot \frac{\eta(\gamma + C}{\alpha}m)H}{(1-\beta)\sqrt{\epsilon}} + \frac{2TL\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)\sqrt{\epsilon}} + \frac{2TL\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)\sqrt{\epsilon}} \\ &\qquad + \frac{2TL^{2}\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)\sqrt{\epsilon}} + \frac{2TL\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)\sqrt{\epsilon}} \\ &\qquad + \frac{2TL^{2}\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)\sqrt{\epsilon}} + \frac{2}{\sqrt{4\beta^{2}(1-\beta^{2})^{2}}} \\ &\qquad + \frac{2TL^{2}\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)\sqrt{\epsilon}} \\ \\ &\qquad + \frac{2TL^{2}\eta^{2}(\gamma^$$

Hence

$$\begin{split} &\frac{\eta\eta_{l}K}{4\sqrt{4\beta_{2}\frac{(1+q^{2})^{3}}{(1-q^{2})^{2}}\eta_{l}^{2}K^{2}G^{2}+\epsilon\cdot T}}\sum_{t=1}^{T}\mathbb{E}[\|\nabla f(\theta_{t})\|^{2}]\\ &\leq \frac{f(\mathbf{z}_{0})-\mathbb{E}[f(\mathbf{z}_{T})]}{T}+\frac{5\eta\eta_{l}^{3}K^{2}L^{2}}{\sqrt{2\epsilon}}(\sigma_{l}^{2}+6K\sigma_{g}^{2})+\frac{C_{1}\eta\eta_{l}KG^{2}d}{T\sqrt{\epsilon}}+\frac{2C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}LG^{2}d}{T\epsilon}\\ &+\left[G+\frac{L\eta\eta_{l}KG}{\sqrt{\epsilon}}+\frac{L\eta\eta_{l}C_{1}KGd}{\epsilon}\right]\cdot\frac{\eta(\gamma+\frac{C}{\alpha m})H}{(1-\beta)\sqrt{\epsilon}}+\frac{L\eta^{2}(\gamma^{2}+\frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)^{2}\epsilon}+\frac{2L\eta^{2}(\gamma^{2}+\frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)\sqrt{\epsilon}}\\ &+\left[3\eta^{2}L+2C_{2}\eta^{2}L+2\sqrt{2(1-\beta_{2})}\eta G\right]\frac{K\eta_{l}^{2}}{2m\epsilon}\sigma_{l}^{2}+\frac{2L^{2}\eta^{2}(\gamma^{2}+\frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)^{2}\epsilon}(\frac{C^{2}}{\alpha^{2}m^{2}}+1),\\ &\text{where }C_{1}=\frac{\beta_{1}}{1-\beta_{1}}+\sqrt{\frac{12q^{2}}{(1-q^{2})^{2}}+\frac{(1-q^{2})^{2}C^{2}}{\alpha^{2}m^{2}q^{2}}}\text{ and }C_{2}=\frac{\beta_{1}^{2}}{(1-\beta_{1})^{2}}+\frac{4(q+\gamma+\frac{\lambda C}{\alpha m})^{2}}{(1-q^{2})^{2}}. \text{ then,} \end{split}$$

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 4\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 + \epsilon \Big[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\Big]}$$

 $\begin{array}{lll} \text{where } \Xi &= \frac{C_1 G^2 d}{\sqrt{\epsilon}} + \frac{2C_1^2 \eta \eta K L G^2 d}{\epsilon}, \Omega &= \left[G + \frac{L \eta \eta_l K G}{\sqrt{\epsilon}} + \frac{L \eta \eta_l C_1 K G d}{\epsilon}\right] \cdot \frac{\eta (\gamma + \frac{C}{\alpha m}) H}{(1 - \beta) \sqrt{\epsilon}} + \\ \frac{2L \eta^2 (\gamma^2 + \frac{C^2}{\alpha^2 m^2}) H^2}{(1 - \beta)^2 \epsilon} &+ \frac{2L \eta^2 (\gamma^2 + \frac{C^2}{\alpha^2 m^2}) H^2}{(1 - \beta) \sqrt{\epsilon}} \frac{5 \eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K \sigma_g^2) + \left[(3 + 2C_2) \eta L + 2\sqrt{2(1 - \beta_2)}G\right] \frac{\eta_l}{2m\eta \epsilon} \sigma_l^2, C_1 &= \frac{\beta_1}{1 - \beta_1} + \sqrt{\frac{12q^2}{(1 - q^2)^2} + \frac{(1 - q^2)^2 C^2}{\alpha^2 m^2 q^2}} + \frac{2T L^2 \eta^2 (\gamma^2 + \frac{C^2}{\alpha^2 m^2}) H^2}{(1 - \beta)^2 \epsilon} (\frac{C^2}{\alpha^2 m^2} + 1) \\ \text{and } C_2 &= \frac{\beta_1^2}{(1 - \beta_1)^2} + \frac{4(q + \gamma + \frac{\lambda C}{\alpha m})^2}{(1 - q^2)^2}. \end{array}$ 

The proof of Theorem 4.6 is similar to the above proof procedure and the detailed proof will not be given here.

#### E.2 PROOF OF COROLLARY 4.3

Let  $\eta_l = \Theta(\frac{1}{\sqrt{T_K}}), T = \mathcal{O}(Km)$  and  $\eta = \Theta(\sqrt{Km})$ , the convergence rate under full participation scheme is  $\mathcal{O}(\frac{1}{T})$ . 

#### E.3 ANALYSIS ON THE PARTIAL PARTICIPATION SETTING FOR FEDBNLACA

Similar to partial participation scheme in Section 3, we have the following convergence analysis. **Theorem E.1.** Under Assumption 3.1-3.4, if the local learning rate  $\eta_l$  satisfies the following condition:  $\eta_l \leq \min\left\{\frac{1}{8KL}, \frac{n(m-1)\epsilon}{48m(n-1)} [K\sqrt{4\beta_2(1+q^2)^3(1-q^2)^{-2}K^2G^2} + \epsilon(\eta L + \sqrt{2(1-\beta_2)}G)]^{-1}\right\},$ then the iterates of Algorithm 2 under partial participation scheme satisfy 

$$\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 8\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 + \epsilon} \left[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\right]$$

 $\begin{array}{ll} \text{, where } \Xi &=& \frac{C_1 G^3 d}{\sqrt{\epsilon}} + \frac{2C_1^2 \eta \eta_l K L G^2 d}{\epsilon}, \Omega &=& \left[G + \frac{L \eta \eta_l K G}{\sqrt{\epsilon}} + \frac{L \eta \eta_l C_1 K G d}{\epsilon} \|\right] \cdot \frac{\eta (\gamma + \frac{C}{\alpha n}) H}{(1 - \beta) \sqrt{\epsilon}} + \frac{4L \eta^2 (\gamma^2 + \frac{C^2}{\alpha^2 n^2}) H^2}{(1 - \beta)^2 \epsilon} + \frac{2L \eta^2 (\gamma^2 + \frac{C^2}{\alpha^2 n^2}) H^2}{(1 - \beta)^2 \epsilon} (\frac{C^2}{\alpha^2 n^2} + 1) + \frac{C_1 \eta \eta_l K L G^2}{\epsilon} + \frac{5\eta^2 K L^2}{\sqrt{2\epsilon}} (\sigma_l^2 + 6K \sigma_g^2) + \left[\eta L + \sqrt{2(1 - \beta_2)} G\right] \frac{\eta_l}{\eta_{n\epsilon}} \sigma_l^2 + \left[\eta L + \sqrt{2(1 - \beta_2)} G\right] \frac{\eta_l (m - n)}{n(m - 1)\epsilon} \left[15K^2 L^2 \eta_l^2 (\sigma_l^2 + 6K \sigma_g^2) + 6K \sigma_g^2\right] + \frac{1}{2} \left[\frac{1}{2} \left[\frac{\eta_l}{\eta_{n\epsilon}} \sigma_l^2 + \frac{1}{2} \left[\frac{\eta_l}{\eta_{n$  $3K\sigma_g^2$  and  $C_1 = \frac{\beta_1}{1-\beta_1} + \frac{m}{n}\sqrt{\frac{12q^2}{(1-q^2)^2} + \frac{(1-q^2)^2C^2}{\alpha^2n^2\sigma^2}}$ Theorem E.2.  $\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 8\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 + \epsilon} \left[\frac{f_0 - f_*}{\eta m KT} + \frac{\Xi}{T} + \Omega\right]$ 

2916  
2917 , where 
$$\Xi = \frac{C_1 G^3 d}{\sqrt{\epsilon}} + \frac{2C_1^2 \eta \eta_l K L G^2 d}{\epsilon}, \Omega = \left[G + \frac{L \eta \eta_l K G}{\sqrt{\epsilon}} + \frac{L \eta \eta_l C_1 K G d}{\epsilon}\right] \cdot \frac{\eta(\gamma+1)H}{(1-\beta)\sqrt{\epsilon}} + \frac{4L \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} + \frac{2L^2 \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} (2) + \frac{C_1 \eta \eta_l K L G^2}{\epsilon} + \frac{5\eta^2 K L^2}{(2)} (\sigma_l^2 + 6K \sigma_a^2) + [\eta L + \frac{2L^2 \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} + \frac{2L^2 \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} (2) + \frac{C_1 \eta \eta_l K L G^2}{\epsilon} + \frac{5\eta^2 K L^2}{(2)} (\sigma_l^2 + 6K \sigma_a^2) + [\eta L + \frac{2L^2 \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} + \frac{2L^2 \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} + \frac{2L^2 \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} (2) + \frac{C_1 \eta \eta_l K L G^2}{\epsilon} + \frac{5\eta^2 K L^2}{(2)} (\sigma_l^2 + 6K \sigma_a^2) + [\eta L + \frac{2L^2 \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} + \frac{2L^2 \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} + \frac{C_1 \eta \eta_l K L G^2}{\epsilon} + \frac{5\eta^2 K L^2}{(2)} (\sigma_l^2 + 6K \sigma_a^2) + [\eta L + \frac{2L^2 \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} + \frac{2L^2 \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} + \frac{C_1 \eta \eta_l K L G^2}{\epsilon} + \frac{5\eta^2 K L^2}{(2)} (\sigma_l^2 + 6K \sigma_a^2) + [\eta L + \frac{2L^2 \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} + \frac{2L^2 \eta^2 (\gamma^2+1)H^2}{(1-\beta)\sqrt{\epsilon}} + \frac{C_1 \eta \eta_l K L G^2}{\epsilon} + \frac{5\eta^2 K L^2}{(1-\beta)\sqrt{\epsilon}} + \frac{C_1 \eta \eta_l K L G^2}{\epsilon} + \frac$$

 $\begin{array}{c} (1-\beta)^{2}\epsilon & (1-\beta)\sqrt{\epsilon} & (1-\beta)\sqrt{\epsilon} & (1-\beta)^{2}\epsilon & (2) & \epsilon & (-\gamma)\sqrt{2\epsilon} & (0_{l}+0K\sigma_{g})+[\eta L+2\sigma_{g}^{2}] \\ \sqrt{2(1-\beta_{2})}G]\frac{\eta_{l}}{\eta_{l}\epsilon}\sigma_{l}^{2} + [\eta L + \sqrt{2(1-\beta_{2})}G]\frac{\eta_{l}(m-n)}{n(m-1)\epsilon}[15K^{2}L^{2}\eta_{l}^{2}(\sigma_{l}^{2}+6K\sigma_{g}^{2})+3K\sigma_{g}^{2}] \text{ and} \\ \\ \textbf{2920} & C_{1} = \frac{\beta_{1}}{1-\beta_{1}} + \frac{m}{n}\sqrt{\frac{12q^{2}}{(1-q^{2})^{2}} + \frac{(1-q^{2})^{2}}{q^{2}}}. \end{array}$ 

*Remark* E.1. When he parameters C = D,  $\frac{C}{\alpha n} = 1$ , the result of Theorem E.1 becomes the result of Theorem E.2. The upper bound for  $\min_{t \in [T]} \mathbb{E} ||\nabla f(\theta_t)||^2$  of partial participation is similar to full participation case but with a larger variance term  $\Omega$ . This is due to the fact that random sampling of participating workers introduces an additional variance during sampling.

**Proof of Theorem E.1:** Notations and equations: From the update rule of Algorithm 2, we have  $\mathbf{e}_1 = 0$ ,  $\mathbf{e}_t = \frac{1}{m} \sum_{i=1}^m \mathbf{e}_t^i$  and  $\mathbf{m}_t = (1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} \widehat{\Delta}_t^i$ . Denote a global uncompressed difference  $\Delta_t = \frac{1}{|S_t|} \sum_{i \in S_t} \Delta_t^i$ . Denote a virtual momentum sequence:  $\mathbf{m}'_t = \beta_1 \mathbf{m}'_{t-1} + (1 - \beta_1) \Delta_t$ , hence we have  $\mathbf{m}'_t = (1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} \Delta_i$ . Define additional two virtual sequences  $\Delta'_t = \frac{1}{n} \sum_{i=1}^m \Delta_t^i$  and  $\widehat{\Delta}'_t = \frac{1}{n} \sum_{i=1}^m \widehat{\Delta}_t^i$ . Note that when the client *i* does not take part in the round of participation at step *t*, we have  $\Delta_t^i = \widehat{\Delta}_t^i = 0$ , therefore,  $\Delta'_t = \Delta_t$  and  $\widehat{\Delta}'_t = \widehat{\Delta}_t$ .

By the aforementioned definition and notation, define a subset  $S_t = \{w_1^t, w_2^t, ..., w_n^t\}$ , then

$$\widehat{\Delta}_t - \Delta_t = \frac{1}{|\mathcal{S}_t|} \sum_{i \in \mathcal{S}_t} (\widehat{\Delta}_t^i - \Delta_t^i) = \frac{1}{n} \sum_{i=1}^m (\widehat{\Delta}_t^i - \Delta_t^i) = \frac{1}{n} \sum_{i=1}^m (\mathbf{e}_t^i - \mathbf{e}_{t+1}^i) = \mathbf{e}_t' - \mathbf{e}_{t+1}'$$

where the compression errors have the same structure,  $\mathbf{e}'_t = \frac{1}{n} \sum_{i=1}^{m} \mathbf{e}^i_t$ . Similar to the previous analysis, we define the following sequence:

$$\Gamma_{t+1} := (1 - \beta_1) \sum_{\tau=1}^{t+1} \beta_1^{t+1-\tau} \mathbf{e}_{\tau}',$$

and keep using the Lyapunov function  $\mathbf{z}_t$  from (E.4). For the expectation of model difference  $\Delta_t$ ,

$$\mathbb{E}_{\mathcal{S}_t}[\Delta_t] = \frac{1}{n} \mathbb{E}_{\mathcal{S}_t} \left[ \sum_{i=1}^n \Delta_t^{w_i} \right] = \mathbb{E}_{\mathcal{S}_t}[\Delta_t^{w_1}] = \frac{1}{m} \sum_{i=1}^m \Delta_t^i = \bar{\Delta}_t.$$

The proof of FedCAMS in partial participation settings has a similar outline combing the proof of partial participation in FedAMS and full participation in FedCAMS. By Assumption 3.1, then

$$\mathbb{E}[f(\mathbf{z}_{t+1})] - f(\mathbf{z}_{t}) \\
\leq \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\rangle\right] \\
-\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \eta \frac{\beta_{1}}{1 - \beta_{1}} \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) m_{t-1}' + \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \Gamma_{t} \right\rangle\right] \\
-\frac{1}{T_{2}'} \\
+ \frac{\eta^{2}L}{2} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} - \frac{\beta_{1}}{1 - \beta_{1}} \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) m_{t-1}' - \left(\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right) \Gamma_{t} \right\|^{2}\right] \\
+ \underbrace{\mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t}), \eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\rangle\right]}_{T_{t}'}$$

2968 Note that the bound for  $T'_2$  is exactly the same as the bound for  $T_2$ . For the three corresponding terms, 2969  $T'_1, T'_3$  and  $T'_4$  which include the second-order momentum estimate of  $\Delta_t$ . For  $T'_1$ , similar to the full participation settings, we have

$$T_{1}^{\prime} \leq \sqrt{2}\mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \eta \frac{\Delta_{t}}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \right\rangle\right] + \sqrt{2}\eta\mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \frac{\Delta_{t}}{\sqrt{\mathbf{v}_{t} + \epsilon}} - \frac{\Delta_{t}}{\sqrt{\beta_{2}\mathbf{v}_{t-1} + \epsilon}} \right\rangle\right].$$
(E.16)

The first term in (E.16) does not change in partial participation scheme. The second term is changed due to the variance of  $\Delta_t$  changes. For the second term of  $T'_1$ , then

$$\sqrt{2}\eta \mathbb{E}\left[\left\langle \nabla f(\theta_t), \frac{\Delta_t}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{\Delta_t}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}} \right\rangle\right] \le \frac{\sqrt{2(1 - \beta_2)}\eta G}{\epsilon} \mathbb{E}[\|\Delta_t\|^2].$$

For  $T'_3$ , similar to the proof of  $T_3$ , we get

$$\sum_{t=1}^{T} T_{3}^{\prime} \leq \frac{\eta^{2} L}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\Delta_{t}\|^{2}] + \eta^{2} L C_{1}^{2} \eta_{l}^{2} K^{2} G^{2} \sum_{t=1}^{T} \mathbb{E}\left[\left\|\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}\right\|^{2}\right],$$

 where  $C_1 = \frac{\beta_1}{1-\beta_1} + \frac{m}{n} \sqrt{\frac{12q^2}{(1-q^2)^2} + \frac{(1-q^2)^2 C^2}{\alpha^2 n^2 q^2}}$  in partial participation, then

$$T_{4}' = \eta \mathbb{E}\left[\left\langle f(\mathbf{z}_{t}) - f(\theta_{t}), \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} \right\rangle\right]$$
  

$$\leq \eta \mathbb{E}\left[\left\|f(\mathbf{z}_{t}) - f(\theta_{t})\right\| \left\|\widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t}\right\|\right]$$
  

$$\leq \eta^{2} L \mathbb{E}\left[\left\|\frac{\beta_{1}}{1 - \beta_{1}} \widehat{\mathbf{V}}_{t-1}^{-1/2} \mathbf{m}_{t-1}' + \widehat{\mathbf{V}}_{t-1}^{-1/2} \Gamma_{t}\right\| \left\|\widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t}\right\|\right]$$
  

$$\leq \frac{C_{1} \eta^{2} \eta_{t}^{2} K^{2} L G^{2}}{\epsilon}.$$

**Bound** of  $T'_5$  , there

$$\begin{aligned} T_{5}^{\prime} &= \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\rangle\right] = \frac{1}{n} \sum_{i}^{n} \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\rangle\right] \\ &= \frac{1}{n} \sum_{i}^{n} \mathbb{E}\left[\left\langle \nabla f(\theta_{t}), \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\rangle\right] + \frac{1}{n} \sum_{i}^{n} \mathbb{E}\left[\left\langle \nabla f(\mathbf{z}_{t}), -\nabla f(\theta_{t}), \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\rangle\right] \\ &\leq \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \nabla f(\theta_{t}) \right\| \cdot \mathbb{E}\left\| \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\| + \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t}) \right\| \cdot \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\| + \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \nabla f(\mathbf{z}_{t}) - \nabla f(\theta_{t}) \right\| \cdot \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\| \\ &\leq \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \nabla f(\theta_{t}) \right\| \cdot \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\| \\ &= \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \nabla f(\theta_{t}) \right\| \cdot \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\| \\ &= \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \nabla f(\theta_{t}) \right\| \cdot \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\| \\ &+ L \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \eta \widehat{\mathbf{V}}_{t}^{-1/2} \Delta_{t} - \eta \frac{\beta_{1}}{1-\beta_{1}} (\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}) \right\| \\ &+ L \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \eta \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\| \\ &= \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \nabla f(\theta_{t}) \right\| \cdot \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} - \widehat{\mathbf{V}}_{t-1}^{-1/2}\right\| \right\| \\ &+ L \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\| \\ &\leq \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\| \\ &\leq \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \nabla f(\theta_{t}) \right\| \cdot \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\| \\ &+ L \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \eta \widehat{\mathbf{V}}_{t}^{-1/2} \widehat{\theta}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t}) \right\| \\ &\leq \frac{1}{n} \sum_{i}^{n} \mathbb{E} \left\| \nabla f(\theta_{t}) \right\| \cdot \mathbb{E} \left\| \frac{\eta}{1-\beta_{1}} \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{$$

$$\begin{split} &+L\frac{1}{n}\sum_{i}^{n}\mathbb{E}\left\|\eta\frac{\beta_{1}}{1-\beta_{1}}\left(\widehat{\mathbf{V}}_{t-1}^{-1/2}-\widehat{\mathbf{V}}_{t}^{-1/2}\right)\mathbf{m}_{t-1}^{\prime}+\eta\left(\widehat{\mathbf{V}}_{t-1}^{-1/2}-\widehat{\mathbf{V}}_{t}^{-1/2}\right)\Gamma_{t}\right\|\cdot\frac{1}{n}\sum_{i}^{n}\mathbb{E}\left\|\frac{\eta}{1-\beta_{1}}\widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\|\\ &+L\frac{1}{n}\sum_{i}^{n}\mathbb{E}\left\|\eta\frac{1}{1-\beta_{1}}\widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\|\cdot\frac{1}{n}\sum_{i}^{n}\mathbb{E}\left\|\frac{\eta}{1-\beta_{1}}\widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\|\\ &+L\frac{1}{n}\sum_{i}^{n}\mathbb{E}\left\|\theta_{t}\right\|\cdot\frac{1}{n}\sum_{i}^{n}\mathbb{E}\left\|\frac{\eta}{1-\beta_{1}}\widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t}^{-}-\theta_{t})\right\|\\ &\leq\left[G+\frac{L\eta\eta_{l}KG}{\sqrt{\epsilon}}+L\eta\eta_{l}C_{1}KG\mathbb{E}\|\widehat{\mathbf{V}}_{t-1}^{-1/2}-\widehat{\mathbf{V}}_{t}^{-1/2}\|\right]\cdot\frac{\eta(\gamma+\frac{C}{\alpha m})H}{(1-\beta)\sqrt{\epsilon}}+\frac{2L\eta^{2}(\gamma^{2}+\frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)^{2}\epsilon}+\frac{2L\eta^{2}(\gamma^{2}+\frac{C^{2}}{\alpha^{2}m^{2}})H^{2}}{(1-\beta)\sqrt{\epsilon}}\end{split}$$

here  $C_1 = \frac{\beta_1}{1-\beta_1} + \frac{m}{n} \sqrt{\frac{12q^2}{(1-q^2)^2} + \frac{(1-q^2)^2C^2}{\alpha^2 n^2 q^2}}$ . Bound of  $T_6'$ , there

$$\begin{split} T_{6}^{\prime} &= \frac{\eta^{2}L^{2}}{(1-\beta)^{2}} \mathbb{E}\left[ \left\| \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t}) \right\|^{2} \right] = \frac{\eta^{2}L^{2}}{(1-\beta)^{2}} \frac{1}{n} \sum_{i}^{n} \mathbb{E}\left[ \left\| \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\widehat{\theta}_{t}} - \theta_{t}) \right\|^{2} \right] \\ &= \frac{\eta^{2}L^{2}}{(1-\beta)^{2}} \frac{1}{n} \sum_{i}^{n} \mathbb{E}\left[ \left\| \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t} - \theta_{t}) - \widehat{\mathbf{V}}_{t}^{-1/2} \frac{1}{M_{t}} \sum_{i \in M_{t}} Q_{t}^{i} \right\|^{2} \right] \\ &\leq \frac{2\eta^{2}L^{2}}{(1-\beta)^{2}} \frac{1}{n} \sum_{i}^{n} \mathbb{E}\left[ \left\| \widehat{\mathbf{V}}_{t}^{-1/2}(\widehat{\theta}_{t} - \theta_{t}) \right\|^{2} \right] + \frac{2\eta^{2}L^{2}}{(1-\beta)^{2}} \frac{1}{n} \sum_{i}^{n} \mathbb{E}\left[ \left\| \widehat{\mathbf{V}}_{t}^{-1/2} \frac{1}{M_{t}} \sum_{i \in M_{t}} Q_{t}^{i} \right\|^{2} \right] \\ &\leq \frac{2L\eta^{2}(\gamma^{2} + \frac{C^{2}}{\alpha^{2}n^{2}})H^{2}}{(1-\beta)^{2}\epsilon} \end{split}$$

Hence, the summation from  $T'_1$  to  $T'_6$  over total iteration T is:

$$\begin{split} & \sum_{k=1}^{3064} \mathbb{E}[f(\mathbf{z}_{T+1})] - f(\mathbf{z}_{1}) = \sum_{t=1}^{T} [T_{1}' + T_{2}' + T_{3}' + T_{4}' + + T_{5}' + T_{6}'] \\ & \sum_{t=1}^{3066} \mathbb{E}[f(\mathbf{z}_{T+1})] - f(\mathbf{z}_{1}) = \sum_{t=1}^{T} [T_{1}' + T_{2}' + T_{3}' + T_{4}' + + T_{5}' + T_{6}'] \\ & \leq -\frac{\eta\eta_{l}K}{4} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{\nabla f(\theta_{t})}{\sqrt[4]{\beta_{2}}\mathbf{v}_{t-1} + \epsilon} \right\|^{2} \right] + \frac{5\eta\eta_{l}^{3}K^{2}L^{2}T}{\sqrt{2\epsilon}} (\sigma_{l}^{2} + 6K\sigma_{g}^{2}) + \frac{\sqrt{2(1 - \beta_{2})}\eta G}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\Delta_{t}\|^{2}] \\ & -\frac{\eta\eta_{l}}{2Km^{2}} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \frac{1}{\sqrt[4]{\beta_{2}}\mathbf{v}_{t-1} + \epsilon} \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t}) \right) \right\|^{2} \right] + C_{1}\eta\eta_{l}KG^{2} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \widehat{\nabla}_{t-1}^{-1/2} - \widehat{\nabla}_{t}^{-1/2} \right\|_{1} \right] \\ & + C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}LG^{2}\epsilon^{-1/2} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \widehat{\nabla}_{t-1}^{-1/2} - \widehat{\nabla}_{t}^{-1/2} \right\|_{1} \right] + C_{1}^{2}\eta^{2}\eta_{l}^{2}K^{2}LG^{2} \sum_{t=1}^{T} \mathbb{E}\left[ \left\| \widehat{\nabla}_{t-1}^{-1/2} - \widehat{\nabla}_{t}^{-1/2} \right\|^{2} \right] \\ & + \frac{\eta^{2}L}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\Delta_{t}\|^{2}] + \frac{C_{1}T\eta^{2}\eta_{l}^{2}K^{2}LG^{2}}{\epsilon} \end{split}$$

$$\begin{split} &+ \left[G + \frac{L\eta\eta_{K}KG}{\sqrt{\epsilon}} + L\eta\eta_{C}_{1}KG\mathbb{E}\||\widehat{\nabla}_{t-1}^{-1/2} - \widehat{\nabla}_{t}^{-1/2}\|\right] \cdot \frac{\eta(\gamma + \frac{G}{\alpha_{N}})H}{(1-\beta)\sqrt{\epsilon}} + \frac{4L\eta^{2}(\gamma^{2} + \frac{G^{2}}{\alpha_{N}})H^{2}}{(1-\beta)^{2}\epsilon} + \frac{2L\eta^{2}(\gamma^{2} + \frac{G^{2}}{\alpha_{N}})H^{2}}{(1-\beta)\sqrt{\epsilon}} \\ &\leq -\frac{\eta\eta_{K}}{4\sqrt{4\beta_{2}\binom{(1-q)}{(1-q)^{2}\eta_{K}^{2}K^{2}G^{2}} + \epsilon} \sum_{i=1}^{T} \mathbb{E}[\|\nabla f(\theta_{i})\|^{2}] + \frac{5\eta\eta_{i}^{2}K^{2}L^{2}T}{\sqrt{2\epsilon}}(\sigma_{i}^{2} + 6K\sigma_{g}^{2}) + \frac{C_{1}\eta\eta_{K}KG^{2}d}{T\sqrt{\epsilon}} \\ &+ \frac{2C_{2}^{2}\eta^{2}\eta_{K}^{2}LG^{2}d}{T\epsilon} - \frac{\eta\eta_{K}}{2\sqrt{4\beta_{2}\binom{(1-q)}{(1-q)^{2}\eta_{K}^{2}}R^{2}G^{2} + \epsilon Km^{2}} \sum_{i=1}^{T} \mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1} \nabla F_{i}(\theta_{i})\right)\right\|^{2}\right] \\ &+ \left(\frac{\eta^{2}\eta_{K}^{2}LK}{n\epsilon} + \frac{\sqrt{2(1-\beta_{2})}\eta\eta_{K}^{2}G}{n\epsilon}\right) \frac{m-n}{m(m-1)} \left[15mK^{3}L^{3}\eta_{L}^{2}(\sigma_{i}^{2} + 6K\sigma_{g}^{2})T\right] \\ &+ \left(9\alpha K^{4}L^{2}\eta_{i}^{2} + 3mK^{2}\right)\sum_{i=1}^{T} \mathbb{E}\left[\left\||\nabla f(\theta_{i})|^{2}\right] + 3mK^{2}T\sigma_{g}^{2}\right] \\ &+ \left(\eta^{2}\eta_{i}^{2}L + \sqrt{2(1-\beta_{2})}\eta\eta_{K}^{2}G\right) \frac{m-1}{m(m-1)}\sum_{i=1}^{T} \mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1} \nabla F_{i}(\theta_{i})\right\right)\right\|^{2}\right] \\ &+ \left[G + \frac{L\eta\eta_{K}KG}{\sqrt{\epsilon}} + \frac{L\eta\eta_{C}G_{K}KGd}{\epsilon}\right] \frac{m-1}{m(m-1)}\sum_{i=1}^{T} \mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1} \nabla F_{i}(\theta_{i})\right\right)\right|^{2}\right] \\ &+ \left[G + \frac{L\eta\eta_{K}G}{\sqrt{\epsilon}} + \frac{L\eta\eta_{K}G_{i}KGd}{\epsilon}\right] \frac{1}{\eta(\gamma + \frac{dm}{m})} + \frac{\eta(\gamma + \frac{dm}{m})}{(1-\beta)\sqrt{\epsilon}} + \frac{4L\eta^{2}(\gamma^{2} + \frac{d^{2}}{m^{2}})H^{2}}{(1-\beta)^{2}q}} + \frac{2L\eta^{2}(\gamma^{2} + \frac{d^{2}}{m^{2}})H^{2}}{(1-\beta)\sqrt{\epsilon}}. \\ The proof outine is similar with previous proof. We take the use of Lemma F2,F9,F13 for constraints of tocal learning rate  $\eta_{n}$ , with the inequality  $[\eta^{2}L + \sqrt{2(1-\beta_{2})}\eta] \frac{1}{\eta^{2}(m-1)}} - \frac{\eta^{2}}{2}M_{i}\sqrt{4\beta_{2}(\frac{(1+q^{2})}{1-q^{2}}}\eta_{i}^{2}K^{2}G^{2} + \epsilon}} \right]^{-1} \leq 0.00 \text{ obtain the constraints} S.$ 

$$fies \frac{\eta(M-1)}{\sqrt{4\beta_{2}(\frac{(1+q^{2})}{1-q^{2}})(1-q^{2})^{-2}(R^{2}G^{2}i+\epsilon}) - \eta^{2}R^{2}(Q^{2}i+\epsilon)} + \frac{2}{2}R^{2}(Q^{2}i+\epsilon)} + \frac{\eta^{2}}{2}R^{2}(Q^{2}i+\epsilon)} + \frac{2}{2}R^{2}(Q^{2}i+\epsilon)} + \frac{1}{2}R^{2}[\|\nabla f(\theta_{i})\|^{2}] \\ \leq \frac{f(\alpha_{0}) = \frac{1}{2}[\eta^{2}(K^{2}G^{2}i+\epsilon)} + \frac{1}{2}R^{2}[\|\nabla f(\theta_{i})\|^{2}] \\ \leq \frac{f(\alpha_{0}) = \frac{1}{2}R^{2}}(Q^{2}i+\epsilon)} + \frac{1}{2}R^{2}[R^{2}G^{2}i+\epsilon} + \frac{$$$$

Therefore  $\min \mathbb{E}[\|\nabla f(\theta_t)\|^2] \le 8\sqrt{4\beta_2 \frac{(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 + \epsilon} \left[\frac{f_0 - f_*}{\eta \eta_l KT} + \frac{\Xi}{T} + \Omega\right]$ 

$$\begin{aligned} \| \begin{array}{l} & \text{where } \Xi = \frac{c_1 C^3 L_2}{\sqrt{\epsilon}} + \frac{2C_1^2 m_t KLG^2 d}{\sqrt{\epsilon}} \Omega = \left[ G + \frac{L m_t KG}{\sqrt{\epsilon}} + \frac{L m_t C_1 KG}{\sqrt{\epsilon}} \| \right] \cdot \frac{n(\tau) \frac{L}{2} \Omega_1^{1/2}}{(1-\beta)^2 \epsilon^2} + \frac{2L^2 \eta^2 (\tau^2 + \frac{C_1^2}{2})^{1/2}}{(1-\beta)^2 \epsilon^2} \frac{(1-\beta)^2 \epsilon^2}{(1-\beta)^2 \epsilon^2} + \frac{2L^2 \eta^2 (\tau^2 + \frac{C_1^2}{2})^{1/2}}{(1-\beta)^2 \epsilon^2} \frac{(1-\beta)^2 \epsilon^2}{(1-\beta)^2 \epsilon^2} + \frac{2L^2 \eta^2 (\tau^2 + \frac{C_1^2}{2})^{1/2}}{(1-\beta)^2 \epsilon^2} + \frac{2L^2 \eta^2 (\tau^2 + \frac{C_1^2}{2})^{1/2}}{(1-2\gamma)^2 \epsilon^2} + \frac{2L^2 \eta^2 (\tau^2 + \frac{C_1^2}{2})^{1/2}} + \frac{2L^2 \eta^2 (\tau^2 + \frac{C_1^2}{2})^{1/2}}{(1-2\gamma)^2 \epsilon^2} + \frac{L^2 \eta^2 (\tau^2 + \frac{C_1^2}{2})^{1/2}}{(1-2\gamma)^2 \epsilon^2} + \frac{$$

**185** Lemma F.2. For the element-wise difference,  $W_t = \frac{1}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{1}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}}, \|W_t\| \le \frac{\sqrt{1-\beta_2}}{\epsilon} \|\Delta_t\|.$ 

Proof. Note that :  $\|W_t\| = \left\|\frac{1}{\sqrt{\mathbf{v}_t + \epsilon}} - \frac{1}{\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}}\right\|$  $= \left\| \frac{(\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon} - \sqrt{\mathbf{v}_t + \epsilon})(\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon} + \sqrt{\mathbf{v}_t + \epsilon})}{\sqrt{\mathbf{v}_t + \epsilon}\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}(\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon} + \sqrt{\mathbf{v}_t + \epsilon})} \right\|$  $= \left\| \frac{\beta_2 \mathbf{v}_{t-1} - \mathbf{v}_t}{\sqrt{\mathbf{v}_t + \epsilon} \sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon} (\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon} + \sqrt{\mathbf{v}_t + \epsilon})} \right\|$  $= \left\| \frac{-(1-\beta_2)\Delta_t^2}{\sqrt{\mathbf{v}_t + \epsilon}\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}(\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon} + \sqrt{\mathbf{v}_t + \epsilon})} \right\|$  $\leq \left\| \frac{(1-\beta_2)\Delta_t^2}{\sqrt{\mathbf{v}_t + \epsilon}\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon}\sqrt{1-\beta_2}\Delta_t} \right\|$  $\leq \frac{\sqrt{1-\beta_2}}{\epsilon} \|\Delta_t\|,$ (F.2)

where the forth equation holds by the update rule of  $v_t$ , i.e.,  $v_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \Delta_t^2$ , and the first inequality holds due  $\sqrt{\mathbf{v}_t + \epsilon} \ge \sqrt{\mathbf{v}_t} \ge \sqrt{1 - \beta_2} \Delta_t$  and  $\sqrt{\beta_2 \mathbf{v}_{t-1} + \epsilon} \ge 0$ . This concludes the proof. 

**Lemma F.3.** For the variance difference sequence  $\widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2}$ , then  $\sum_{t=1}^{T} \left\| \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right\|_{1} \le \frac{d}{\sqrt{\epsilon}}, \sum_{t=1}^{T} \left\| \widehat{\mathbf{V}}_{t-1}^{-1/2} - \widehat{\mathbf{V}}_{t}^{-1/2} \right\|^{2} \le \frac{d}{\epsilon}.$ (F.3) 

*Proof.* By the definition of variance matrix  $\hat{\mathbf{V}}_t$ , and the non-decreasing update of FedCAMS, i.e.,  $\widehat{\mathbf{v}}_{t-1} \leq \widehat{\mathbf{v}}_t = \max(\widehat{\mathbf{v}}_{t-1}, \mathbf{v}_t, \epsilon), \text{then}$ ш

where the inequality holds by the definition of  $\hat{v}_t \in \mathbb{R}^d$ : For the sum of the variance difference under  $\ell_2$  norm, then

$$\begin{split} \sum_{t=1}^{T} \left\| \hat{\mathbf{V}}_{t-1}^{-1/2} - \hat{\mathbf{V}}_{t}^{-1/2} \right\|^{2} &= \sum_{t=1}^{T} \left\| \frac{1}{\sqrt{\hat{\mathbf{v}}_{t-1}}} - \frac{1}{\sqrt{\hat{\mathbf{v}}_{t}}} \right\|^{2} \\ &= \sum_{t=1}^{T} \left( \frac{1}{\sqrt{\hat{\mathbf{v}}_{t-1}}} - \frac{1}{\sqrt{\hat{\mathbf{v}}_{t}}} \right)^{2} \\ &\leq \sum_{t=1}^{T} \left( \frac{1}{\widehat{\mathbf{v}}_{t-1}} - \frac{1}{\widehat{\mathbf{v}}_{t}} \right) \\ &\leq \frac{1}{\widehat{\mathbf{v}}_{0}} - \frac{1}{\widehat{\mathbf{v}}_{T}} \\ &\leq \frac{d}{\epsilon}, \end{split}$$
(F.5)

where the first inequality holds by the element-wise operation:  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d, \mathbf{0} \leq \mathbf{y} \leq \mathbf{x}$ , we have  $(\mathbf{x} - \mathbf{y})^2 \le (\mathbf{x} - \mathbf{y})(\mathbf{x} + \mathbf{y})) = \mathbf{x}^2 - \mathbf{y}^2$ . It concludes the proof. 

Lemma F.4. The compression error has the following absolute bound 

$$\|\mathbf{e}_t^i\|^2 \le \frac{4q^2}{(1-q^2)^2} \eta_l^2 K^2 G^2, \quad \|\mathbf{e}_t\|^2 \le \frac{4q^2}{(1-q^2)^2} \eta_l^2 K^2 G^2.$$
(F.6)

*Proof.* For all  $t \in [T]$ , by Assumption 3.4 and Young's inequality, then if  $i \notin M_t$ 

$$\begin{aligned} \|\mathbf{e}_{t+1}^{i}\|^{2} &= \|\Delta_{t}^{i} + \mathbf{e}_{t}^{i} - \mathcal{C}(\Delta_{t}^{i} + \mathbf{e}_{t}^{i})\|^{2} \\ &\leq q^{2} \|\Delta_{t}^{i} + \mathbf{e}_{t}^{i}\|^{2} \\ &\leq q^{2}(1+\rho) \|\mathbf{e}_{t}^{i}\|^{2} + q^{2} \left(1 + \frac{1}{\rho}\right) \|\Delta_{t}^{i}\|^{2} \\ &\leq \frac{1+q^{2}}{2} \|\mathbf{e}_{t}^{i}\|^{2} + \frac{2q^{2}}{1-q^{2}} \|\Delta_{t}^{i}\|^{2}, \end{aligned}$$

if  $i \in M_t$ 

$$\begin{split} \|\mathbf{e}_{t+1}^{i}\|^{2} &= \|\Delta_{t}^{i} + \mathbf{e}_{t}^{i} - \mathcal{C}(\Delta_{t-1}^{i} + \mathbf{e}_{t-1}^{i})\|^{2} \\ &= \|\Delta_{t}^{i} + \mathbf{e}_{t}^{i} - \mathcal{C}(\Delta_{t}^{i} + \mathbf{e}_{t}^{i}) + \mathcal{C}(\Delta_{t}^{i} + \mathbf{e}_{t}^{i}) - \mathcal{C}(\Delta_{t-1}^{i} + \mathbf{e}_{t-1}^{i})\|^{2} \\ &\leq 2q^{2} \|\Delta_{t}^{i} + \mathbf{e}_{t}^{i}\|^{2} + 2\|\mathcal{C}(\Delta_{t}^{i} + \mathbf{e}_{t}^{i}) - \mathcal{C}(\Delta_{t-1}^{i} + \mathbf{e}_{t-1}^{i})\|^{2} \\ &\leq 2q^{2}(1+\rho)\|\mathbf{e}_{t}^{i}\|^{2} + \left(2q^{2}(1+\frac{1}{\rho}) + \frac{2C^{2}}{\alpha^{2}m^{2}}\right)\|\Delta_{t}^{i}\|^{2} \\ &\leq \frac{2(1+q^{2})}{2}\|\mathbf{e}_{t}^{i}\|^{2} + \left(\frac{4q^{2}}{1-q^{2}} + \frac{2C^{2}}{\alpha^{2}m^{2}}\right)\|\Delta_{t}^{i}\|^{2}, \end{split}$$

where the last inequality holds by choosing  $\rho = \frac{1-q^2}{2q^2}$ . 

Thus obtain the absolute bound for the error terms 

$$\|\mathbf{e}_t^i\|^2 \le \frac{4q^2}{(1-q^2)^2} \eta_l^2 K^2 G^2$$

$$\|\mathbf{e}_t^i\|^2 \le (\frac{8q^2}{(1-q^2)^2} + \frac{2C^2}{\alpha^2 m^2} \frac{1-q^2}{2q^2})\eta_l^2 K^2 G^2,$$

then

or

$$\|\mathbf{e}_{t}\|^{2} = \left\|\frac{1}{m}\sum_{i=1}^{m}\mathbf{e}_{t}^{i}\right\|^{2} \leq \frac{1}{m}\sum_{i=1}^{m}\|\Delta_{t}^{i} + \mathbf{e}_{t}^{i} - \mathcal{C}(\Delta_{t}^{i} + \mathbf{e}_{t}^{i})\|^{2} + \frac{1}{M_{t}}\sum_{i\in M_{t}}\|\Delta_{t}^{i} + \mathbf{e}_{t}^{i} - \mathcal{C}(\Delta_{t-1}^{i} + \mathbf{e}_{t-1}^{i})\|^{2}$$
$$\leq \frac{4q^{2}}{(1-q^{2})^{2}}\eta_{l}^{2}K^{2}G^{2} + \left(\frac{8q^{2}}{(1-q^{2})^{2}} + \frac{2C^{2}}{\alpha^{2}m^{2}}\frac{1-q^{2}}{2q^{2}}\right)\eta_{l}^{2}K^{2}G^{2}.$$
(F.7)

In the case of partial participation, suppose that client i has the participated time set  $T_i$ , and we rewrite the  $T_i = \{t_0, t_1, ..., t_{p_i}\}$ , where  $t_0 < t_1 < \cdots < t_{p_i}$ . Since when client *i* are not selected to participate local training, the error stay unchanged. Then for  $t_s \in \mathcal{T}_i$ 

thus by the similar recursive approach, since  $\mathbf{e}_{t_0}^i = 0$ , 

$$\mathbb{E}[\|\mathbf{e}_{t_{s+1}}^{i}\|^{2}] \leq \frac{2q^{2}}{1-q^{2}} \sum_{\tau=1}^{s} \left(\frac{1+q^{2}}{2}\right)^{s-\tau} \mathbb{E}[\|\Delta_{t_{\tau}}^{i}\|^{2}].$$

Thus obtain the absolute bound for the error terms,

$$\|\mathbf{e}_t^i\|^2 \le \frac{4q^2}{(1-q^2)^2} \eta_l^2 K^2 G^2$$

or

$$\|\mathbf{e}_t^i\|^2 \le (\frac{8q^2}{(1-q^2)^2} + \frac{2C^2}{\alpha^2 m^2} \frac{1-q^2}{2q^2})\eta_l^2 K^2 G^2,$$

$$\begin{aligned} \|\mathbf{e}_{t}\|^{2} &= \left\|\frac{1}{m}\sum_{i=1}^{m}\mathbf{e}_{t}^{i}\right\|^{2} \leq \frac{1}{m}\sum_{i=1}^{m}\|\Delta_{t}^{i} + \mathbf{e}_{t}^{i} - \mathcal{C}(\Delta_{t}^{i} + \mathbf{e}_{t}^{i})\|^{2} + \frac{1}{M_{t}}\sum_{i\in M_{t}}\|\Delta_{t}^{i} + \mathbf{e}_{t}^{i} - \mathcal{C}(\Delta_{t-1}^{i} + \mathbf{e}_{t-1}^{i})\|^{2} \\ &\leq \frac{4q^{2}}{(1-q^{2})^{2}}\eta_{t}^{2}K^{2}G^{2} + \left(\frac{8q^{2}}{(1-q^{2})^{2}} + \frac{2C^{2}}{\alpha^{2}m^{2}}\frac{1-q^{2}}{2q^{2}}\right)\eta_{t}^{2}K^{2}G^{2}. \end{aligned}$$
(F.8)

This is the end of the proof. 

**Lemma F.5.** Under Assumptions 3.2 and 3.4, for FedAMS, we have  $\|\nabla f(\theta)\| \leq G, \|\tilde{\Delta}_t\| \leq$  $\begin{aligned} &\eta_l KG, \|\mathbf{m}_t\| \leq \eta_t KG \text{ and } \|\mathbf{v}_t\| \leq \eta_t^2 K^2 G^2. \text{ For FedCAMS, we have } \|\nabla f_i(\theta)\| \leq G, \|\widehat{\Delta}_t\|^2 \leq \\ &\frac{4(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2 \|\nabla f(\theta)\| \leq G, \|\widehat{\Delta}_t\|^2 \leq \frac{4(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2, \|\mathbf{m}_t'\| \leq \eta_l KG \text{ and } \|\mathbf{v}_t\| \leq \\ &\frac{4(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2, \text{where } \mathbf{m}_t' = \beta_1 \mathbf{m}_{t-1}' + (1-\beta_1) \widetilde{\Delta}_t. \end{aligned}$ 

*Proof.* Since f has G-bounded stochastic gradients, for any  $\theta$  and  $\xi$ , we have  $\|\nabla f(\theta, \xi)\| \leq G$ , that

$$\|\nabla f(\theta)\| = \|\mathbb{E}_{\xi} \nabla f(\theta, \xi)\| \le \mathbb{E}_{\xi} \|\nabla f(\theta, \xi)\| \le G.$$

For Fed, the model difference  $\tilde{\Delta}_t^i$ , by definition, has the following formula, 

therefore,

$$\tilde{\Delta}_{t}^{i} = \theta_{t,K}^{i} - \theta_{t} = -\eta \sum_{k=1}^{K} \mathbf{g}_{t,k}^{i} \text{ or } \tilde{\Delta}_{t}^{i} = \theta_{t-1,K}^{i} - \theta_{t-1} = -\eta \sum_{k=1}^{K} \mathbf{g}_{t-1,k}^{i},$$

$$\|\tilde{\Delta}_{t}\| \leq \frac{1}{|S_{t}|} \sum_{i \in S_{t}} \|\tilde{\Delta}_{t}^{i}\| \leq \eta_{l} KG.$$

$$3325$$

Thus the bound for momentum  $\mathbf{m}_t$  and variance  $\mathbf{v}_t$  has the formula of 

$$\|\mathbf{m}_t\| = (1 - \beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \|\tilde{\Delta}_t\| \le \eta_l KG,$$

$$\|\mathbf{v}_t\| = (1 - \beta_2) \sum_{\tau=1}^t \beta_2^{t-\tau} \|\tilde{\Delta}_t\|^2 \le \eta_l^2 K^2 G^2$$

For the compressed version, FedCAMS, 

$$\begin{split} \Delta_t^i &= \theta_{t,K}^i - \theta_t = -\eta \sum_{k=1}^K \mathbf{g}_{t,k}^i ,\\ \|\Delta_t\| &\leq \frac{1}{m} \sum_{i=1}^m \|\Delta_t^i\| \leq \eta_l KG. \end{split}$$

Thus the bound for momentum  $\mathfrak{m}_t$  and variance  $v_t$  has the formula of

3343  
3344  
3345 
$$\|\mathbf{m}_t\| = (1 - \beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \|\Delta_t\| \le \eta_l KG,$$

3346  
3347 
$$\|\mathbf{v}_t\| = (1 - \beta_2) \sum_{j=1}^{t} \beta_2^{t-\tau} \|\Delta_j\|$$

3347 
$$\|\mathbf{v}_t\| = (1 - \beta_2) \sum_{\tau=1}^{t} \beta_2^{t-\tau} \|\Delta_t\|^2 \le \eta_l^2 K^2 G^2$$

$$\begin{split} \|\Delta_t^i\|^2 &\leq \|\mathcal{C}(\Delta_t^i + \mathbf{e}_t^i)\|^2 \\ &\leq \|\mathcal{C}(\Delta_t^i + \mathbf{e}_t^i) - (\Delta_t^i + \mathbf{e}_t^i) + (\Delta_t^i + \mathbf{e}_t^i)\|^2 \\ &\leq 2(q^2 + 1)\|\Delta_t^i + \mathbf{e}_t^i\|^2 \\ &\leq 4(q^2 + 1)[\|\Delta_t^i\|^2 + \|\mathbf{e}_t^i\|^2], \end{split}$$

 then

$$\|\widehat{\Delta}_t\|^2 = \left\|\frac{1}{m}\sum_{i=1}^m \widehat{\Delta}_t^i\right\|^2 \le \frac{4(1+q^2)^3}{(1-q^2)^2}\eta^2 K^2 G^2 + \left(\frac{8(1+q^2)^3}{(1-q^2)^2} + \frac{4(1-q^2)(q^2+1)C^2}{\alpha^2 m^2}\right)\eta_l^2 K^2 G^2$$

 $\text{if } i \notin M_t, \|\widehat{\Delta}_t^i\|^2 \leq \frac{4(1+q^2)^3}{(1-q^2)^2} \eta_l^2 K^2 G^2. \text{ if } i \in M_t, \|\widehat{\Delta}_t^i\|^2 \leq \big(\frac{8(1+q^2)^3}{(1-q^2)^2} + \frac{4(1-q^2)(q^2+1)C^2}{\alpha^2 m^2}\big) \eta_l^2 K^2 G^2.$ 

where the third inequality holds due to Assumption 3.4, and the last inequality holds due to LemmaF.4 e virtual momentum sequence  $||\mathbf{m}'_t||$  has the same bound as  $\mathbf{m}_t$  of FedAMS. For the variance sequence of FedCAMS, we have

$$\|\mathbf{v}_t\| = (1 - \beta_2) \sum_{\tau=1}^t \beta_2^{t-\tau} \|\hat{\Delta}_t\|^2 \le \frac{4(1 + q^2)^3}{(1 - q^2)^2} \eta_l^2 K^2 G^2.$$

This concludes the proof.

**Lemma F.6.** The global model difference  $\Delta_t = \sum_{i=1}^m \tilde{\Delta}_t^i$  in full participation cases satisfy

$$\mathbb{E}[\|\tilde{\Delta}_t\|^2] \le \frac{K\eta_l^2}{m} \sigma_l^2 + \frac{\eta_l^2}{m^2} \mathbb{E}\left\|\sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i)\right\|^2 + \frac{1}{m^2} \mathbb{E}\left\|\frac{1}{M_t} \sum_{i \in M_t} q_t^i\right\|^2.$$
(F.9)

*Proof.* For  $\mathbb{E}[\|\tilde{\Delta}_t\|^2]$  in full participation case, then

г

$$\begin{split} \mathbb{E}[\|\tilde{\Delta}_{t}\|^{2}] &= \mathbb{E}\left[\left\|\frac{1}{m}\sum_{i=1}^{m}\sum_{k=0}^{K-1}\eta_{i}\mathbf{g}_{t,k}^{i} - \frac{1}{M_{t}}\sum_{i\in M_{t}}q_{t}^{i}\right\|^{2}\right] \\ &= \frac{\eta_{t}^{2}}{m^{2}}\mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\mathbf{g}_{t,k}^{i} - \frac{1}{\eta_{l}M_{c}}\sum_{i\in M_{c}}q_{t}^{i}\right\|^{2}\right] \\ &= \frac{\eta_{t}^{2}}{m^{2}}\mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}(\mathbf{g}_{t,k}^{i} - \nabla F_{i}(\theta_{t,k}^{i}))\right\|^{2}\right] + \frac{\eta_{t}^{2}}{m^{2}}\mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i}) - \frac{1}{\eta_{l}M_{t}}\sum_{i\in M_{t}}q_{t}^{i}\right\|^{2}\right] \\ &\leq \frac{K\eta_{t}^{2}}{m}\sigma_{t}^{2} + \frac{2\eta_{t}^{2}}{m^{2}}\left[\mathbb{E}\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2} + \mathbb{E}\left\|\frac{1}{\eta_{t}M_{t}}\sum_{i\in M_{t}}q_{t}^{i}\right\|^{2}\right] \\ &= \frac{K\eta_{t}^{2}}{m}\sigma_{t}^{2} + \frac{2\eta_{t}^{2}}{m^{2}}\mathbb{E}\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2} + \frac{2}{m^{2}}\mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{t}^{i}\right\|^{2}, \end{split}$$

where the inequality holds by Assumption 3.2. The end of the proof. 

**Lemma F.7.** The global model difference  $\Delta_t = \sum_{i=1}^m \Delta_t^i$  in full participation cases satisfy  $1^{2}$ 

$$\mathbb{E}[\|\Delta_t\|^2] \leq \frac{K\eta_l^2}{m}\sigma_l^2 + \frac{\eta_l^2}{m^2}\mathbb{E}\left\|\sum_{i=1}^m\sum_{k=0}^{K-1}\nabla F_i(\theta_{t,k}^i)\right\|$$

.

*Proof.* For  $\mathbb{E}[\|\Delta_t\|^2]$  in full participation case, then

$$\mathbb{E}[\|\Delta_t\|^2] = \mathbb{E}\left[\left\|\frac{1}{m}\sum_{i=1}^m\sum_{k=0}^{K-1}\eta_l \mathbf{g}_{t,k}^i\right\|^2\right]$$

$$= \frac{\eta_l^2}{m^2} \mathbb{E}\left[\left\|\sum_{i=1}^m\sum_{k=0}^{K-1}\mathbf{g}_{t,k}^i\right\|^2\right]$$

$$= \frac{\eta_l^2}{m^2} \mathbb{E}\left[\left\|\sum_{i=1}^m\sum_{k=0}^{K-1}(\mathbf{g}_{t,k}^i - \nabla F_i(\theta_{t,k}^i))\right\|^2\right] + \frac{\eta_l^2}{m^2} \mathbb{E}\left[\left\|\sum_{i=1}^m\sum_{k=0}^{K-1}\nabla F_i(\theta_{t,k}^i)\right\|^2\right]$$

$$\leq \frac{K\eta_l^2}{m}\sigma_l^2 + \frac{\eta_l^2}{m^2} \left[\mathbb{E}\left\|\sum_{i=1}^m\sum_{k=0}^{K-1}\nabla F_i(\theta_{t,k}^i)\right\|^2\right]$$

$$= \frac{K\eta_l^2}{m}\sigma_l^2 + \frac{\eta_l^2}{m^2} \mathbb{E}\left\|\sum_{i=1}^m\sum_{k=0}^{K-1}\nabla F_i(\theta_{t,k}^i)\right\|^2, \quad (F.10)$$
here the inequality holds by Assumption 3.2. The end of the proof.

wł Ŋ

**Lemma F.8.** The global model difference 
$$\tilde{\Delta}_t = \sum_{i \in S_t} \tilde{\Delta}_t^i$$
 in partial participation cases satisfy  

$$\mathbb{E}[\|\tilde{\Delta}_t\|^2] = \frac{K\eta_l^2}{n}\sigma_l^2 + \frac{\eta_l^2(m-n)}{mn(m-1)}[15mK^3L^3\eta_l^2(\sigma_l^2 + 6K\sigma_g^2) + 90mK^4L^2\eta_l^2 + 3mK^2\|\nabla f(\theta_t)\|^2 + 3mK^2\sigma_g^2] + \frac{\eta_l^2(n-1)}{mn(m-1)}\mathbb{E}\Big[\Big\|\sum_{i=1}^m\sum_{k=0}^{K-1}\nabla F_i(\theta_{t,k}^i) - \frac{1}{M_t}\sum_{i\in M_t}q_t^i\Big\|^2\Big].$$

Proof. Then

$$\begin{split} \mathbb{E}[\|\tilde{\Delta}_{t}\|^{2}] &= \mathbb{E}\left[\left\|\frac{1}{n}\sum_{i\in\mathcal{S}_{t}}\tilde{\Delta}_{t}^{i}\right\|^{2}\right] \\ &= \frac{1}{n^{2}}\mathbb{E}\left[\left\|\sum_{i=1}^{m}\mathbb{I}\{i\in\mathcal{S}_{t}\}\tilde{\Delta}_{t}^{i}\right\|^{2}\right] \\ &= \frac{1}{n^{2}}\mathbb{E}\left[\left\|\eta_{t}^{2}\sum_{i=1}^{m}\mathbb{I}\{i\in\mathcal{S}_{t}\}\sum_{k=0}^{K-1}\left[\mathbf{g}_{t,k}^{i}-\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}+\left\|\eta_{t}^{2}\sum_{i=1}^{m}\mathbb{I}\{i\in\mathcal{S}_{t}\}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})-\frac{1}{M_{t}}\sum_{i\in\mathcal{M}_{t}}q_{t}^{i}\right\|^{2}\right] \\ &= \frac{1}{n^{2}}\mathbb{E}\left[\left\|\eta_{t}^{2}\sum_{i=1}^{m}\mathbb{P}\{i\in\mathcal{S}_{t}\}\sum_{k=0}^{K-1}\left[\mathbf{g}_{t,k}^{i}-\nabla F_{i}(\theta_{t,k}^{i})\right]\right\|^{2}+\left\|\eta_{t}^{2}\sum_{i=1}^{m}\mathbb{P}\{i\in\mathcal{S}_{t}\}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})-\frac{1}{M_{t}}\sum_{i\in\mathcal{M}_{t}}q_{t}^{i}\right\|^{2}\right] \\ &= \frac{\eta_{t}^{2}}{mn}\mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\left[\mathbf{g}_{t,k}^{i}-\nabla F_{i}(\theta_{t,k}^{i})\right]\right\|^{2}\right]+\frac{1}{n^{2}}\mathbb{E}\left[\left\|\eta_{t}^{2}\sum_{i=1}^{m}\mathbb{P}\{i\in\mathcal{S}_{t}\}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})-\frac{1}{M_{t}}\sum_{i\in\mathcal{M}_{t}}q_{t}^{i}\right\|^{2}\right] \\ &\leq \frac{K\eta_{t}^{2}}{n}\sigma_{t}^{2}+\frac{2\eta_{t}^{2}}{n^{2}}\mathbb{E}\left\|\sum_{i=1}^{m}\mathbb{P}\{i\in\mathcal{S}_{t}\}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}+\frac{2}{n^{2}}\mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in\mathcal{M}_{t}}q_{t}^{i}\right\|^{2}, \quad (F.11) \end{split}$$

 $q_t^i$ 

where the fifth equation holds due to  $\mathbb{P}\{i \in S_t\} = \frac{n}{m}$ . Note that

where the second equation holds due to  $\|\sum_{i=1}^{m} \theta_i\|^2 = \sum_{i=1}^{m} m \|\theta_i\|^2 - \frac{1}{2} \sum_{i \neq j} \|\theta_i - \theta_j\|^2$ . By the sampling strategy (without replacement),  $\mathbb{P}\{i \in S_t\} = \frac{n}{m}$  and  $\mathbb{P}\{i, j \in S_t\} = \frac{n(n-1)}{m(m-1)}$ , thus

$$\begin{aligned} \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \mathbb{P}\{i \in \mathcal{S}_{t}\} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \\ &= \sum_{i=1}^{m} \mathbb{P}\{i \in \mathcal{S}_{t}\} \left\| \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} + \sum_{i \neq j} \mathbb{P}\{i, j \in \mathcal{S}_{t}\} \left\langle \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}), \sum_{k=0}^{K-1} \nabla F_{j}(\theta_{t,k}^{j}) \right\rangle \\ &= \frac{n}{m} \sum_{i=1}^{m} \left\| \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} + \frac{n(n-1)}{m(m-1)} \sum_{i \neq j} \left\langle \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}), \sum_{k=0}^{K-1} \nabla F_{j}(\theta_{t,k}^{j}) \right\rangle \\ &= \frac{n^{2}}{m} \sum_{i=1}^{m} \left\| \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} - \frac{n(n-1)}{2m(m-1)} \sum_{i \neq j} \left\| \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) - \sum_{k=0}^{K-1} \nabla F_{j}(\theta_{t,k}^{j}) \right\|^{2} \\ &= \frac{n(m-n)}{m(m-1)} \sum_{i=1}^{m} \left\| \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} + \frac{n(n-1)}{m(m-1)} \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2}, \end{aligned}$$
(F.13)

where the third equation holds due to  $\langle \mathbf{x}, \mathbf{y} \rangle = \frac{1}{2} [\|\mathbf{x}^2 + \mathbf{y}^2\| - \|\mathbf{x} - \mathbf{y}\|^2]$  and the last equation holds due to  $\frac{1}{2} \sum_{i \neq j} \|\theta_i - \theta_j\|^2 = \sum_{i=1}^m m \|\theta_i\|^2 - \|\sum_{i=1}^m \theta_i\|^2$ . Therefore, for the last term in F.11, then

$$\mathbb{E}[\|\tilde{\Delta}_{t}\|^{2}] = \frac{K\eta_{l}^{2}}{n}\sigma_{l}^{2} + \frac{2\eta_{l}^{2}(m-n)}{mn(m-1)}\sum_{i=1}^{m}\mathbb{E}\left[\left\|\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}\right] + \frac{2\eta_{l}^{2}(n-1)}{mn(m-1)}\mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}\right] + \frac{2}{n^{2}}\mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{t}^{i}\right\|^{2}$$

The second term in (F.13) is bounded partially following Reddi et al. (2020),

$$\begin{split} \sum_{i=1}^{m} \left\| \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} &= \sum_{i=1}^{m} \mathbb{E} \left\| \sum_{k=0}^{K-1} \left[ \nabla F_{i}(\theta_{t,k}^{i}) - \nabla F_{i}(\theta_{t}) + \nabla F_{i}(\theta_{t}) - \nabla f(\theta_{t}) + \nabla f(\theta_{t}) \right] \right\|^{2} \\ &\leq 3 \sum_{i=1}^{m} \mathbb{E} \left\| \sum_{k=0}^{K-1} \left[ \nabla F_{i}(\theta_{t,k}^{i}) - \nabla F_{i}(\theta_{t}) \right] \right\|^{2} + 3mK^{2}\sigma_{g}^{2} + 3mK^{2} \|\nabla f(\theta_{t})\|^{2} \\ &\leq 3KL^{2} \sum_{i=1}^{m} \sum_{k=0}^{K-1} \mathbb{E} [\|\theta_{t,k}^{i} - \theta_{t}\|^{2}] + 3mK^{2}\sigma_{g}^{2} + 3mK^{2} \|\nabla f(\theta_{t})\|^{2} \\ &\leq 15mK^{3}L^{3}\eta_{l}^{2}(\sigma_{l}^{2} + 6K\sigma_{g}^{2}) + (90mK^{4}L^{2}\eta_{l}^{2} + 3mK^{2}) \|\nabla f(\theta_{t})\|^{2} + 3mK^{2}\sigma_{g}^{2} \\ &\quad (F.14) \end{split}$$

where the last inequality holds by applying Lemma C.9 (also follows from Reddi et al. (2020)).
Substituting (F.14) into (F.13), this concludes the proof.

**Lemma F.9.** The global model difference  $\Delta_t = \sum_{i \in S_t} \Delta_t^i$  in partial participation cases satisfy

$$\begin{split} \mathbf{\tilde{s}12} & \quad \mathbb{E}[\|\Delta_t\|^2] = \frac{K\eta_l^2}{n} \sigma_l^2 + \frac{\eta_l^2(m-n)}{mn(m-1)} [15mK^3L^3\eta_l^2(\sigma_l^2 + 6K\sigma_g^2) + 90mK^4L^2\eta_l^2 + 3mK^2\|\nabla f(\theta_t)\|^2 \\ \mathbf{\tilde{s}14} & \quad + 3mK^2\sigma_g^2] + \frac{\eta_l^2(n-1)}{mn(m-1)} \mathbb{E}\Big[\Big\|\sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i)\Big\|^2\Big]. \end{split}$$

*Proof.* we have

$$\begin{aligned}
& \text{3521} \\
& \text{3522} \\
& \text{3523} \\
& \text{E}[\|\Delta_t\|^2] = \mathbb{E}\left[\left\|\frac{1}{n}\sum_{i\in\mathcal{S}_t}\Delta_t^i\right\|^2\right] \\
& \text{3524} \\
& \text{3525} \\
& \text{3526} \\
& \text{3526} \\
& \text{3527} \\
& \text{3528} \\
& \text{3529} \\
& \text{3529} \\
& \text{3529} \\
& \text{3530} \\
& \text{3531} \\
& \text{3532} \\
& \text{3532} \\
& \text{3532} \\
& \text{3533} \\
& \text{3534} \\
& \text{3535} \\
& \text{3536} \\
& \text{3537} \\
& \text{3537} \\
& \text{3538} \\
& \text{3536} \\
& \text{3536} \\
& \text{3537} \\
& \text{3537} \\
& \text{3538} \\
& \text{3538} \\
& \text{3538} \\
& \text{3539} \\
& \text{3540} \\
& \text{3540} \\
& \text{3540} \\
& \text{3559} \\
& \text{3540} \\
& \text{3559} \\
& \text{3560} \\
& \text{3560} \\
& \text{3577} \\
& \text{3577} \\
& \text{3578} \\
& \text{3578} \\
& \text{3577} \\
& \text{3578} \\
& \text{3579} \\
& \text{3577} \\
& \text{3578} \\
& \text{3579} \\
& \text{3577} \\
& \text{3578} \\
& \text{3578} \\
& \text{3577} \\
& \text{3578} \\
& \text{3578} \\
& \text{3577} \\
& \text{3577} \\
& \text{3578} \\
& \text{3577} \\
& \text{3578} \\
& \text{3577} \\
& \text{$$

where the fifth equation holds due to  $\mathbb{P}\{i \in \mathcal{S}_t\} = \frac{n}{m}$ .Note that

$$\begin{aligned} \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} &= \sum_{i=1}^{m} \left\| \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} + \sum_{i \neq j} \left\langle \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}), \sum_{k=0}^{K-1} \nabla F_{j}(\theta_{t,k}^{j}) \right\rangle \\ &= \sum_{i=1}^{m} m \left\| \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} - \frac{1}{2} \sum_{i \neq j} \left\| \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) - \sum_{k=0}^{K-1} \nabla F_{j}(\theta_{t,k}^{j}) \right\|^{2}, \end{aligned}$$
(F.16)

where the second equation holds due to  $\|\sum_{i=1}^{m} \theta_i\|^2 = \sum_{i=1}^{m} m \|\theta_i\|^2 - \frac{1}{2} \sum_{i \neq j} \|\theta_i - \theta_j\|^2$ . By the sampling strategy (without replacement), we have  $\mathbb{P}\{i \in S_t\} = \frac{n}{m}$  and  $\mathbb{P}\{i, j \in S_t\} = \frac{n(n-1)}{m(m-1)}$ , thus

$$\begin{split} & \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \mathbb{P}\{i \in \mathcal{S}_{t}\} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} \\ &= \sum_{i=1}^{m} \mathbb{P}\{i \in \mathcal{S}_{t}\} \left\| \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} + \sum_{i \neq j} \mathbb{P}\{i, j \in \mathcal{S}_{t}\} \left\langle \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}), \sum_{k=0}^{K-1} \nabla F_{j}(\theta_{t,k}^{j}) \right\rangle \\ &= \frac{n}{m} \sum_{i=1}^{m} \left\| \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} + \frac{n(n-1)}{m(m-1)} \sum_{i \neq j} \left\langle \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}), \sum_{k=0}^{K-1} \nabla F_{j}(\theta_{t,k}^{j}) \right\rangle \end{split}$$

where the third equation holds due to  $\langle \mathbf{x}, \mathbf{y} \rangle = \frac{1}{2} [\|\mathbf{x}^2 + \mathbf{y}^2\| - \|\mathbf{x} - \mathbf{y}\|^2]$  and the last equation holds due to  $\frac{1}{2} \sum_{i \neq j} \|\theta_i - \theta_j\|^2 = \sum_{i=1}^m m \|\theta_i\|^2 - \|\sum_{i=1}^m \theta_i\|^2$ . Therefore, for the last term in F.15, then

 $= \frac{n(m-n)}{m(m-1)} \sum_{i=1}^{m} \left\| \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i) \right\|^2 + \frac{n(n-1)}{m(m-1)} \left\| \sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i) \right\|^2,$ 

 $= \frac{n^2}{m} \sum_{i=1}^{m} \left\| \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i) \right\|^2 - \frac{n(n-1)}{2m(m-1)} \sum_{i\neq i} \left\| \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i) - \sum_{k=0}^{K-1} \nabla F_j(\theta_{t,k}^j) \right\|^2$ 

$$\mathbb{E}[\|\Delta_t\|^2] = \frac{K\eta_l^2}{n}\sigma_l^2 + \frac{\eta_l^2(m-n)}{mn(m-1)}\sum_{i=1}^m \mathbb{E}\left[\left\|\sum_{k=0}^{K-1}\nabla F_i(\theta_{t,k}^i)\right\|^2\right] + \frac{\eta_l^2(n-1)}{mn(m-1)}\mathbb{E}\left[\left\|\sum_{i=1}^m\sum_{k=0}^{K-1}\nabla F_i(\theta_{t,k}^i)\right\|^2\right].$$
(F.17)

The second term in (F.17) is bounded partially following (Reddi et al., 2020),

$$\begin{split} \sum_{i=1}^{m} \left\| \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i}) \right\|^{2} &= \sum_{i=1}^{m} \mathbb{E} \left\| \sum_{k=0}^{K-1} \left[ \nabla F_{i}(\theta_{t,k}^{i}) - \nabla F_{i}(\theta_{t}) + \nabla F_{i}(\theta_{t}) - \nabla f(\theta_{t}) + \nabla f(\theta_{t}) \right] \right\|^{2} \\ &\leq 3 \sum_{i=1}^{m} \mathbb{E} \left\| \sum_{k=0}^{K-1} \left[ \nabla F_{i}(\theta_{t,k}^{i}) - \nabla F_{i}(\theta_{t}) \right] \right\|^{2} + 3mK^{2}\sigma_{g}^{2} + 3mK^{2} \|\nabla f(\theta_{t})\|^{2} \\ &\leq 3KL^{2} \sum_{i=1}^{m} \sum_{k=0}^{K-1} \mathbb{E} [\|\theta_{t,k}^{i} - \theta_{t}\|^{2}] + 3mK^{2}\sigma_{g}^{2} + 3mK^{2} \|\nabla f(\theta_{t})\|^{2} \\ &\leq 15mK^{3}L^{3}\eta_{l}^{2}(\sigma_{l}^{2} + 6K\sigma_{g}^{2}) + (90mK^{4}L^{2}\eta_{l}^{2} + 3mK^{2}) \|\nabla f(\theta_{t})\|^{2} + 3mK^{2}\sigma_{g}^{2}, \end{split}$$
(F.18)

where the last inequality holds by applying Lemma C.9 (also follows from Reddi et al. (2020)). Substituting F.18 into F.17, this concludes the proof.  $\Box$ 

**Lemma F.10.** Under Assumptions 3.1-3.4, for the momentum sequence  $\mathbf{m}_t = (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \tilde{\Delta}_{\tau}$ and accumulated error sequence  $\Gamma_t = (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \mathbf{e}_{\tau}$  in full participation settings, then

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathbf{m}_{t}\|^{2}] \leq \frac{TK\eta_{l}^{2}}{m}\sigma_{l}^{2} + \frac{2\eta_{l}^{2}}{m^{2}}\sum_{t=1}^{T} \mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\right] + \frac{2}{m^{2}}\sum_{t=1}^{T} \mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{t}^{i}\right\|^{2}\right]$$

3610 and

$$\sum_{t=1}^{3611} \mathbb{E}[\|\mathbf{\Gamma}_{t}\|^{2}] \leq \frac{4Tq^{2}}{(1-q^{2})^{2}} \frac{K\eta_{i}^{2}}{m} \sigma_{l}^{2} + \frac{\eta_{l}^{2}}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{t=1}^{T} \mathbb{E}\Big\|\sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i})\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{t=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{M_{t}} \sum_{i \in M_{t}} q_{t}^{i}\Big\|^{2} + \frac{1}{m^{2}} \frac{4q^{2}}{(1-q^{2})^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{m^{2}} \sum_{i=1}^{T} \frac{1}{m^{2}} \frac{1}{(1-q^{2})^{2}} \sum_{i=1}^{T} \frac{1}{m^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{m^{2}} \sum_{i=1}^{T} \frac{1}{m^{2}} \sum_{i=1}^{T} \frac{1}{m^{2}} \sum_{i=1}^{T} \frac{1}{m^{2}} \sum_{i=1}^{T} \mathbb{E}\Big\|\frac{1}{m^{2}} \sum_{i=1}^{T} \frac{1}{m^{2}} \sum_{i=1}^{T} \frac{1}{m^{2}} \sum_{i=1}^{T} \frac{1}{m^{2}} \sum_{i=1}^$$

*Proof.* By the updating rule, we get

 $\mathbb{E}[\|\mathbf{m}_t\|^2] = \mathbb{E}\left[\|(1-\beta_1)\sum_{\tau=1}^t \beta_1^{t-\tau} \tilde{\Delta}_{\tau}\|^2\right]$ 

 $\leq (1-\beta_1)^2 \mathbb{E}\left[\left(\sum_{i=1}^t \beta_1^{t-\tau} \tilde{\Delta}_{\tau,i}\right)^2\right]$ 

 $\leq (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \mathbb{E}[\|\tilde{\Delta}_{\tau}\|^2]$ 

 $\leq (1-\beta_1)^2 \mathbb{E}\left[\left(\sum_{j=1}^t \beta_1^{t-\tau}\right) \left(\sum_{j=1}^t \beta_1^{t-\tau} \tilde{\Delta}_{\tau,i}^2\right)\right]$ 

$$\leq \frac{K\eta_l^2}{m}\sigma_l^2 + \frac{\eta_l^2}{m^2}(1-\beta_1)\sum_{\tau=1}^t \beta_1^{t-\tau} \mathbb{E}\left[ \left\| \sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i) - \frac{1}{M_t} \sum_{i\in M_t} q_t^i \right\|^2 \right] \\ - \frac{2\eta_l^2}{m^2}(1-\beta_1)\sum_{\tau=1}^t \beta_1^{t-\tau} \mathbb{E}\left\| \sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i) \right\|^2 + \frac{2}{m^2}(1-\beta_1)\sum_{\tau=1}^t \beta_1^{t-\tau} \mathbb{E}\left\| \frac{1}{M_t} \sum_{i\in M_t} q_t^i \right\|^2$$

where the second inequality holds by applying Cauchy-Schwarz inequality, and the third inequality holds by summation of series. The last inequality holds by Lemma F.6. Hence summing over  $t = 1, \dots, T$ , we get

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathbf{m}_{t}\|^{2}] \leq \frac{TK\eta_{l}^{2}}{m}\sigma_{l}^{2} + \frac{2\eta_{l}^{2}}{m^{2}}\sum_{t=1}^{T} \mathbb{E}\left[\left\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\right] + \frac{2}{m^{2}}\sum_{t=1}^{T} \mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in M_{t}}q_{t}^{i}\right\|^{2}\right]$$

**Lemma F.11.** Under Assumptions 3.1-3.4, for the momentum sequence  $\mathbf{m}_t = (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \Delta_{\tau}$ and accumulated error sequence  $\Gamma_t = (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \mathbf{e}_{\tau}$  in full participation settings, we have

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathbf{m}_{t}\|^{2}] \leq \frac{TK\eta_{l}^{2}}{m}\sigma_{l}^{2} + \frac{\eta_{l}^{2}}{m^{2}}\sum_{t=1}^{T} \mathbb{E}\Big[\Big\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\Big]$$

and

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathbf{\Gamma}_t\|^2] \le \frac{4Tq^2}{(1-q^2)^2} \frac{K\eta_i^2}{m} \sigma_l^2 + \frac{\eta_l^2}{m^2} \frac{4q^2}{(1-q^2)^2} \sum_{t=1}^{T} \mathbb{E}\Big\|\sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i)\Big\|^2.$$

*Proof.* By the updating rule, we get

$$\mathbb{E}[\|\mathbf{m}_t\|^2] = \mathbb{E}\left[\|(1-\beta_1)\sum_{\tau=1}^t \beta_1^{t-\tau} \Delta_{\tau}\|^2\right]$$
$$\leq (1-\beta_1)^2 \mathbb{E}\left[\left(\sum_{\tau=1}^t \beta_1^{t-\tau} \Delta_{\tau,i}\right)^2\right]$$

$$\leq (1-\beta_1) \sum_{\tau=1}^{t} \beta_1^{t-\tau} \mathbb{E}[\|\Delta_{\tau}\|^2]$$

 $\leq (1-\beta_1)^2 \mathbb{E}\left[\left(\sum_{\tau=1}^{t} \beta_1^{t-\tau}\right) \left(\sum_{\tau=1}^{\iota} \beta_1^{t-\tau} \Delta_{\tau,i}^2\right)\right]$ 

$$\leq \frac{K\eta_l^2}{m}\sigma_l^2 + \frac{\eta_l^2}{m^2}(1-\beta_1)\sum_{\tau=1}^t \beta_1^{t-\tau} \mathbb{E}\left[\left\|\sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i)\right\|^2\right]$$

 $= \left\| \frac{1}{m} \sum_{i=1}^{m} [\Delta_t^i + \mathbf{e}_t^i] - \frac{1}{m} \sum_{i=1}^{m} \mathcal{C}(\Delta_t^i + \mathbf{e}_t^i) + \frac{1}{M} \sum_{i=1}^{m} \mathcal{C}(q_t^i) \right\|$ 

$$\frac{\eta_l^2}{m^2} (1 - \beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \mathbb{E} \left\| \sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{t,k}^i) \right\|^2$$
3679

where the second inequality holds by applying Cauchy-Schwarz inequality, and the third inequality holds by summation of series. The last inequality holds by Lemma F.7. Hence summing over  $t = 1, \cdots, T$ , we have

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathbf{m}_{t}\|^{2}] \leq \frac{TK\eta_{l}^{2}}{m}\sigma_{l}^{2} + \frac{\eta_{l}^{2}}{m^{2}}\sum_{t=1}^{T} \mathbb{E}\Big[\Big\|\sum_{i=1}^{m}\sum_{k=0}^{K-1}\nabla F_{i}(\theta_{t,k}^{i})\Big].$$

For the compress error  $e_t$ , by Assumption 3.4-C.1, then

 $\left\|\mathbf{e}_{t+1}\right\| = \left\|\frac{1}{m}\sum_{i=1}^{m}\mathbf{e}_{t+1}^{i}\right\|$ 

$$= \left\| \frac{m}{i-1} - m\frac{m}{i-1} - \mathcal{C}\left(\frac{1}{m}\sum_{i=1}^{m} [\Delta_{t}^{i} + \mathbf{e}_{t}^{i}]\right) - \frac{m}{m}\sum_{i=1}^{m} [\Delta_{t}^{i} + \mathbf{e}_{t}^{i}] - \mathcal{C}\left(\frac{1}{m}\sum_{i=1}^{m} [\Delta_{t}^{i} + \mathbf{e}_{t}^{i}]\right) + \left\| \mathcal{C}\left(\frac{1}{m}\sum_{i=1}^{m} [\Delta_{t}^{i} + \mathbf{e}_{t}^{i}]\right) - \frac{1}{m}\sum_{i=1}^{m} \mathcal{C}(\Delta_{t}^{i} + \mathbf{e}_{t}^{i}) \right\| + \left\| \frac{1}{M_{t}}\sum_{i\in M_{t}} \mathcal{C}(q_{t}^{i}) \right\|$$

$$\leq q \left\| \frac{1}{m} \sum_{i=1}^{m} [\Delta_t^i + \mathbf{e}_t^i] \right\| + \gamma \left\| \frac{1}{m} \sum_{i=1}^{m} \Delta_t^i \right\| + \frac{C}{\alpha m} \left\| \frac{1}{M_t} \sum_{i \in M_t} \Delta_t^i \right\|$$

$$\leq q \|\Delta_t\| + q \|\mathbf{e}_t\| + \gamma \|\Delta_t\| + \frac{\lambda C}{\alpha m} \|\Delta_t\|$$

$$= q \|\mathbf{e}_t\| + (q + \gamma + \frac{\lambda C}{\alpha m}) \|\Delta_t\|,$$

where the first equation holds by the definition for error  $e_{t+1}$ , and the second one holds by the update rule for  $e_{t+1}^i$ . The first inequality holds by  $\|\mathbf{a} + \mathbf{b}\| \leq \|\mathbf{a}\| + \|\mathbf{b}\|$ , and the second one holds by Assumption 3.4-C.1. Thus by Young's inequality, then 

$$\|\mathbf{e}_{t+1}\|^2 \le \left(q\|\mathbf{e}_t\| + (q+\gamma + \frac{\lambda C}{\alpha m})\|\Delta_t\|\right)^2$$

$$\leq q^{2}(1+\rho) \|\mathbf{e}_{t}\|^{2} + (q+\gamma+\frac{\lambda C}{\alpha m})^{2}(1+\rho^{-1}) \|\Delta_{t}\|^{2}$$
$$= \frac{1+q^{2}}{2} \|\mathbf{e}_{t}\|^{2} + \frac{(q+\gamma+\frac{\lambda C}{\alpha m})^{2}}{1-q^{2}} \|\Delta_{t}\|^{2},$$
(F.19)

where the equation holds by letting  $\rho = \frac{1-q^2}{2q^2}$ , and  $1 + \rho^{-1} = \frac{1+q^2}{1-q^2} \le \frac{2}{1-q^2}$ , then by the similar recursive approach in the proof of Lemma F.3, we have 

$$\begin{aligned} & \mathbb{E}[\|\mathbf{e}_{t+1}\|^2] \leq \frac{2(q+\gamma+\frac{\lambda C}{\alpha m})^2}{1-q^2} \sum_{\tau=1}^t \left(\frac{1+q^2}{2}\right)^{t-\tau} \mathbb{E}[\|\Delta_{\tau}\|^2] \\ & \text{3723} \\ & \text{3724} \\ & \text{3725} \\ & \leq \frac{4(q+\gamma+\frac{\lambda C}{\alpha m})^2}{(1-q^2)^2} \frac{K\eta_l^2}{m} \sigma_l^2 + \frac{\eta_l^2}{m^2} \frac{2(q+\gamma+\frac{\lambda C}{\alpha m})^2}{1-q^2} \sum_{\tau=1}^t \left(\frac{1+q^2}{2}\right)^{t-\tau} \mathbb{E}\left[\left\|\sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\mathbf{x}_{\tau,k}^i)\right\|^2\right] \\ & \text{3726} \end{aligned}$$

For the sequence  $\Gamma_t$ , similar as the previous analysis, we have 

$$\begin{array}{l} \mathbf{3728} \\ \mathbf{3729} \\ \mathbf{3730} \\ \mathbf{3731} \\ \mathbf{3732} \end{array} \qquad \mathbb{E}[\|\mathbf{\Gamma}_{\ell}\|^2] = \mathbb{E}\left[\left\| (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \mathbf{e}_{\tau} \right\|^2 \\ \leq (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \mathbb{E}[\|\mathbf{e}_{\tau}\|^2] \\ \end{array} \right]$$

Summing over  $t = 1, \dots, T$ , then

$$\begin{split} \sum_{t=1}^{T} \mathbb{E}[\|\mathbf{\Gamma}_{t}\|^{2}] &\leq \frac{4T(q+\gamma+\frac{\lambda C}{\alpha m})^{2}}{(1-q^{2})^{2}} \frac{K\eta_{i}^{2}}{m} \sigma_{i}^{2} + \frac{\eta_{l}^{2}}{m^{2}} \frac{2(q+\gamma+\frac{\lambda C}{\alpha m})^{2}}{1-q^{2}} \sum_{t=1}^{T} \sum_{\tau=1}^{t} \left(\frac{1+q^{2}}{2}\right)^{t-\tau} \mathbb{E}\left[\left\|\sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{\tau,k}^{i})\right\|^{2}\right] \\ &+ \frac{\eta_{l}^{2}}{m^{2}} \frac{2(q+\gamma+\frac{\lambda C}{\alpha m})^{2}(1-\beta_{1})}{1-q^{2}} \sum_{\tau=1}^{t} \beta_{1}^{t-\tau} \sum_{j=1}^{\tau} \left(\frac{1+q^{2}}{2}\right)^{\tau-j} \mathbb{E}\left\|\frac{1}{M_{t}} \sum_{i\in M_{t}} q_{t}^{i}\right\|^{2} \\ &\leq \frac{4T(q+\gamma+\frac{\lambda C}{\alpha m})^{2}}{(1-q^{2})^{2}} \frac{K\eta_{l}^{2}}{m} \sigma_{l}^{2} + \frac{\eta_{l}^{2}}{m^{2}} \frac{4(q+\gamma+\frac{\lambda C}{\alpha m})^{2}}{(1-q^{2})^{2}} \sum_{t=1}^{T} \mathbb{E}\left[\left\|\sum_{i=1}^{m} \sum_{k=0}^{K-1} \nabla F_{i}(\theta_{t,k}^{i})\right\|^{2}\right]. \end{split}$$
The end of the proof

 $\leq \frac{4(q+\gamma+\frac{\lambda C}{\alpha m})^2}{(1-q^2)^2} \frac{K\eta_l^2}{m} \sigma_l^2 + \frac{\eta_l^2}{m^2} \frac{2(q+\gamma+\frac{\lambda C}{\alpha m})^2(1-\beta_1)}{1-q^2} \sum_{\tau=1}^t \beta_1^{t-\tau} \sum_{j=1}^\tau \left(\frac{1+q^2}{2}\right)^{\tau-j} \mathbb{E}\left[ \left\| \sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{j,k}^i) \right\|^2 \right] + \frac{1}{2} \sum_{j=1}^t \left(\frac{1+q^2}{2}\right)^{\tau-j} \mathbb{E}\left[ \left\| \sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{j,k}^i) \right\|^2 \right] + \frac{1}{2} \sum_{j=1}^t \left(\frac{1+q^2}{2}\right)^{\tau-j} \mathbb{E}\left[ \left\| \sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{j,k}^i) \right\|^2 \right] + \frac{1}{2} \sum_{j=1}^t \left(\frac{1+q^2}{2}\right)^{\tau-j} \mathbb{E}\left[ \left\| \sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{j,k}^i) \right\|^2 \right] + \frac{1}{2} \sum_{j=1}^t \left(\frac{1+q^2}{2}\right)^{\tau-j} \mathbb{E}\left[ \left\| \sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{j,k}^i) \right\|^2 \right] + \frac{1}{2} \sum_{j=1}^t \left(\frac{1+q^2}{2}\right)^{\tau-j} \mathbb{E}\left[ \left\| \sum_{i=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{j,k}^i) \right\|^2 \right] + \frac{1}{2} \sum_{j=1}^t \left(\frac{1+q^2}{2}\right)^{\tau-j} \mathbb{E}\left[ \left\| \sum_{j=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{j,k}^i) \right\|^2 \right] + \frac{1}{2} \sum_{j=1}^t \left(\frac{1+q^2}{2}\right)^{\tau-j} \mathbb{E}\left[ \left\| \sum_{j=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{j,k}^i) \right\|^2 \right] + \frac{1}{2} \sum_{j=1}^t \left(\frac{1+q^2}{2}\right)^{\tau-j} \mathbb{E}\left[ \left\| \sum_{j=1}^m \sum_{k=0}^{K-1} \nabla F_i(\theta_{j,k}^k) \right\|^2 \right] + \frac{1}{2} \sum_{j=1}^t \left(\frac{1+q^2}{2}\right)^{\tau-j} \mathbb{E}\left[ \left(\frac{1+q^2}{2}\right)^{\tau-j} \mathbb{E}\left[ \sum_{j=1}^m \sum_{k=0}^{K-1} \nabla F_j(\theta_{j,k}^k) \right] \right] + \frac{1}{2} \sum_{j=1}^t \left(\frac{1+q^2}{2}\right)^{\tau-j} \mathbb{E}\left[ \sum_{j=1}^t \sum_{k=0}^{K-1} \nabla F_j(\theta_{j,k}^k) \right] + \frac{1}{2} \sum_{j=1}^t \sum_{j=1}^t \sum_{k=0}^t \nabla F_j(\theta_{j,k}^k) \right] + \frac{1}{2} \sum_{j=1}^t \sum_{k=0}^t \nabla F_j(\theta_{j,k}^k) + \frac{1}{2} \sum_{j=1}^t \sum_{j=1}^t \sum_{k=0}^t \nabla F_j(\theta_{j,k}^k) \right] + \frac{1}{2} \sum_{j=1}^t \sum_{j=1}^t$ 

The end of the proof.

**Lemma F.12.** Under Assumptions 3.1-3.4, for the momentum sequence  $\mathbf{m}_t = (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \tilde{\Delta}_{\tau}$ inpartial participation settings, we have 

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathbf{m}_{t}\|^{2}] \leq \frac{KT\eta_{l}^{2}}{n} \sigma_{l}^{2} + \frac{2\eta_{l}^{2}}{n^{2}} \sum_{t=1}^{T} \mathbb{E}\Big[\Big\|\sum_{i\in\mathcal{S}_{t}}\sum_{k=0}^{K-1} \nabla F_{i}(\mathbf{x}_{t,k}^{i})\Big\|^{2}\Big] + \frac{2}{n^{2}} \mathbb{E}\left\|\frac{1}{M_{t}}\sum_{i\in\mathcal{M}_{t}}q_{t}^{i}\right\|^{2}.$$

*Proof.* The proof outline is the same as the proof of Lemma F.10, the main difference is  $E[||\Delta_t||^2]$ has changed, so we need to apply Lemma F.8 instead of Lemma F.6 during the proof.

**Lemma F.13.** Under Assumptions 3.1-3.4, for the momentum sequence  $\mathbf{m}_t = (1-\beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \Delta_{\tau}$ in partial participation settings, then 

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathbf{m}_t\|^2] \le \frac{KT\eta_l^2}{n} \sigma_l^2 + \frac{\eta_l^2}{n^2} \sum_{t=1}^{T} \mathbb{E}\Big[\Big\|\sum_{i \in \mathcal{S}_t} \sum_{k=0}^{K-1} \nabla F_i(\mathbf{x}_{t,k}^i)\Big\|^2\Big].$$

*Proof.* The proof outline is the same as the proof of Lemma F.11, the main difference is  $E[||\Delta_t||^2]$ has changed, so we need to apply Lemma F.9 instead of Lemma F.7 during the proof. 

**Lemma F.14.** (*This lemma directly follows from Lemma 3 in Reddi et al.* (2020)). For local learning rate which satisfying  $\eta_l \leq \frac{1}{8KL}$ , the local model difference after  $k \; (\forall k \in \{0, 1, ..., K - L\})$ 1}) steps local updates satisfies, 

$$\frac{1}{m} \sum_{i=1}^{m} \mathbb{E}[\|\mathbf{x}_{t,k}^{i} - \mathbf{x}_{t}\|^{2}] \le 5K\eta_{l}^{2}(\sigma_{l}^{2} + 6K\sigma_{g}^{2}) + 30K^{2}\eta_{l}^{2}\mathbb{E}[\|\nabla f(\mathbf{x}_{t})\|^{2}].$$

*Proof.* The proof of Lemma F.14 is exactly same as the proof of Lemma 3 in (Reddi et al. (2020)).  $\Box$