
Communication-Efficient Adaptive Federated Learning

Yujia Wang¹ Lu Lin² Jinghui Chen¹

Abstract

Federated learning is a machine learning training paradigm that enables clients to jointly train models without sharing their own localized data. However, the implementation of federated learning in practice still faces numerous challenges, such as the large communication overhead due to the repetitive server-client synchronization and the lack of adaptivity by SGD-based model updates. Despite that various methods have been proposed for reducing the communication cost by gradient compression or quantization, and the federated versions of adaptive optimizers such as FedAdam are proposed to add more adaptivity, the current federated learning framework still cannot solve the aforementioned challenges all at once. In this paper, we propose a novel communication-efficient adaptive federated learning method (FedCAMS) with theoretical convergence guarantees. We show that in the nonconvex stochastic optimization setting, our proposed FedCAMS achieves the same convergence rate of $\mathcal{O}(\frac{1}{\sqrt{TKm}})$ as its non-compressed counterparts. Extensive experiments on various benchmarks verify our theoretical analysis.

1. Introduction

Federated learning (FL) (Konečný et al., 2016; McMahan et al., 2017) has recently become a popular machine learning training paradigm where multiple clients cooperate to jointly learn a machine learning model. In the federated learning setting, training data is distributed across a large number of clients, or edge devices, such as smartphones, personal computers, or IoT devices. These clients own valuable data for training a variety of machine learning

models, yet those raw client data is not allowed to share with the server or other clients due to privacy and regulation concerns. Federated Learning (Konečný et al., 2016; McMahan et al., 2017) works by having each client train the ML model locally based on its own data, while having the clients iteratively exchanging and synchronizing their local ML model parameters with each other through a central server. McMahan et al. (2017) proposed FedAvg algorithm, whose global model is updated by averaging multiple steps of local stochastic gradient descent (SGD) updates, and it has become one of the most popular FL methods.

Despite the ability to jointly train the model without directly sharing the data, the implementation of FL in practice still faces several major challenges such as (1) *large communication overhead* due to the repetitive synchronization between the server and the clients; and (2) *lack of adaptivity* as SGD-based update may not be suitable for heavy-tail stochastic gradient noise distributions, which often arise in training large-scale models such as BERT (Devlin et al., 2018), GPT-3 (Brown et al., 2020), GAN (Goodfellow et al., 2014) or ViT (Dosovitskiy et al., 2021).

Note that various attempts have been made to solve the aforementioned challenges individually but not all of them at once. In terms of reducing communication costs, one can avoid transmitting the complete model updates when synchronizing. Several works, including (Reisizadeh et al., 2020; Jin et al., 2020; Jhunjunwala et al., 2021; Chen et al., 2021b), have studied the compressed and quantized federated learning optimization method based on FedAvg. Another way is to reduce the number of participating clients such that only part of the clients participate in the model training at each round (Yang et al., 2021; Li et al., 2019b; Nishio & Yonetani, 2019; Li & Wang, 2019). Besides, the network resources allocation also plays an important role in communication-efficient federated learning problems (Li et al., 2019a; Yang et al., 2020). For adaptivity concerns, recently, FedAdam (Reddi et al., 2020) and other variants, such as FedYogi (Reddi et al., 2020) and FedAMSGrad (Tong et al., 2020) were proposed to introduce adaptive gradient methods (Kingma & Ba, 2014; Reddi et al., 2018) into federated learning framework and provided provable convergence guarantees. However, it is still an open problem how to achieve communication efficient adaptive federated optimization while still providing rigorous convergence

¹College of Information Sciences and Technology, Pennsylvania State University, State College, PA, United States
²Department of Computer Science, University of Virginia, Charlottesville, VA, United States. Correspondence to: Jinghui Chen <jzc5917@psu.edu>.

guarantees.

In this paper, we aim to develop a new compressed federated adaptive gradient optimization method that is communication-efficient while also provably convergences. Specifically, we first propose FedAMS, a variant of FedAdam, with an improved convergence analysis over the original FedAdam. Based on FedAMS, we propose FedCAMS, a **F**ederated **C**ommunication-compressed **A**MSGrad with **M**ax **S**tabilization (FedCAMS), which addresses both the communication and adaptivity challenges within one training framework. We summarize our contributions as follows:

- We provide an improved analysis on the convergence behaviour of FedAMS, a variant of the existing federated adaptive gradient method FedAdam (Reddi et al., 2020), whose analysis is simplified for only considering the case where no momentum is been used. In particular, we prove that FedAMS (with momentum) can achieve the same convergence rate of $\mathcal{O}(\frac{1}{\sqrt{TKm}})$ w.r.t total iterations T , the number of local updates K , and the number of workers m for both full participation and partial participation schemes.
- We propose a new communication-efficient adaptive federated optimization method, FedCAMS, which to the best of our knowledge, for the first time, achieves both communication efficiency and adaptivity in federated learning with one single learning framework. FedCAMS largely reduces the communication cost by error feedback and compression strategy and it is compatible with various commonly-used compressors in practice. We prove that FedCAMS achieves the same convergence rate of $\mathcal{O}(\frac{1}{\sqrt{TKm}})$, as its uncompressed counterpart FedAMS.
- We conduct experiments on various benchmarks and show that our proposed FedAMS and FedCAMS achieve good adaptivity in training real-world machine learning models. Furthermore, we show that FedCAMS effectively reduced the communication cost (number of bits for communication) by orders of magnitude while sacrificing little in terms of prediction accuracy.

Notation: For vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $\sqrt{\mathbf{x}}; \mathbf{x}^2; \mathbf{x}=\mathbf{y}$ denote the element-wise square root, square, and division of the vectors. For vector \mathbf{x} and matrix A , $\|\cdot\|$ denotes the ℓ_2 norm of vector/matrix, i.e., $\|\mathbf{x}\| = \|\mathbf{x}\|_2$ and $\|A\| = \|A\|_2$.

2. Related Work

SGD and Adaptive Gradient Methods: Stochastic gradient descent (SGD) (Robbins & Monro, 1951) has been widely applied in training machine learning models for decades. Although SGD is straightforward to implement, it is known to be sensitive to parameters and relatively slow to converge when facing heavy-tail stochastic gradi-

ent noise. Adaptive gradient methods were proposed to overcome these issues of SGD, including AdaGrad (Duchi et al., 2011), RMSProp (Tieleman et al., 2012), AdaDelta (Zeiler, 2012). Adam (Kingma & Ba, 2014) and its variant AMSGrad (Reddi et al., 2018), are tremendously used in training deep neural networks, and other variants (Luo et al., 2019; Loshchilov & Hutter, 2017; Chen et al., 2020a) also play important roles in improving adaptive gradient methods through different aspects.

Federated Learning: As the demand of locally data storing and training models at edge devices, Federated Learning (Konečný et al., 2016; Li et al., 2020) rapidly attracts growing interest in recent years. Federated Averaging method (FedAvg) (McMahan et al., 2017) works by periodically averaging local SGD updates. Stich (2018) provided a concise theoretical convergence guarantee for local SGD. Lin et al. (2018) proposed a variant of local SGD with empirical improvements. There are many works based on FedAvg such as FedProx (Li et al., 2020), FedNova (Wang et al., 2020), SCAFFOLD (Karimireddy et al., 2020), and other work discussed the variants of FedAvg (Yang et al., 2021; Li et al., 2019b; Hsu et al., 2019; Wang et al., 2019). Reddi et al. (2020) recently proposed several adaptive federated optimization methods including FedAdagrad, FedYogi and FedAdam to overcome the existing convergence issues of FedAvg. Chen et al. (2020b) proposed Local AMSGrad and Tong et al. (2020) proposed a family of federated adaptive gradient methods with calibrations. Another line of research focused on addressing data heterogeneity issues or the network resource allocation issues (Ghosh et al., 2019; Li & Wang, 2019; Yang et al., 2020).

Communication-Compressed Federated Learning: Various strategies have been proposed for reducing communication costs in distributed learning for SGD based algorithms (Bernstein et al., 2018; Seide et al., 2014; Alistarh et al., 2017; Basu et al., 2019; Stich et al., 2018; Stich & Karimireddy, 2019; Karimireddy et al., 2019) and also adaptive gradient methods (Tang et al., 2021; Wang et al., 2022). In terms of federated learning, many studies have tried to apply the aforementioned methods to FedAvg and have attracted growing interest recently, e.g., FedPAQ (Reisizadeh et al., 2020), FedCOM (Haddadpour et al., 2021), sign SGD in federated learning (Jin et al., 2020), communication-efficient federated learning (Chen et al., 2021b), AdaQuantFL (Jhunjhunwala et al., 2021). However, there are fewer attempts to develop communication-efficient adaptive gradient methods in federated learning, which is our key focus in this work.

3. Proposed Method

In this paper, we aim to study the following federated learning nonconvex optimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m F_i(\mathbf{x}); \quad (3.1)$$

where m is the total amount of local clients, d denotes the dimension of the model parameters, $F_i(\mathbf{x}) = \mathbb{E}_{\xi \sim \mathcal{D}_i} F_i(\mathbf{x}; \xi)$ is the local nonconvex loss function on client i associated with a local distribution \mathcal{D}_i . In the stochastic setting, we can only obtain the unbiased estimator of $F_i(\mathbf{x})$, i.e., the stochastic gradient $\mathbf{g}_i^j = \nabla F_i(\mathbf{x}; \xi_j)$. In the non i.i.d setting, distributions $\mathcal{D}_i, \mathcal{D}_j$ can vary from each other, i.e., $\mathcal{D}_i \neq \mathcal{D}_j, \forall i \neq j$.

FedAvg (McMahan et al., 2017) is a commonly used optimization approach to solve (3.1). Let \mathbf{x}_t denotes the global model parameters before the t -th iteration. Now at iteration t , the participating client i from the selected subset \mathcal{S}_t (with size n) receives the model \mathbf{x}_t from the server, conducts K steps of local SGD updates with local learning rate η , obtains the local model $\mathbf{x}_{t,K}^i$. Client i then sends the model difference $\Delta_t^i = \mathbf{x}_{t,K}^i - \mathbf{x}_t$ to the server. And the server updates the global model difference Δ_t by simply averaging the local model differences Δ_t^i . The server then updates the global model \mathbf{x}_{t+1} by $\mathbf{x}_{t+1} = \mathbf{x}_t + \Delta_t$, which is the same¹ as directly averaging the local model $\mathbf{x}_{t,K}^i$, i.e., $\mathbf{x}_{t+1} = \mathbf{x}_t + \frac{1}{n} \sum_{i \in \mathcal{S}_t} (\mathbf{x}_{t,K}^i - \mathbf{x}_t) = \frac{1}{n} \sum_{i \in \mathcal{S}_t} \mathbf{x}_{t,K}^i$.

FedAdam was then proposed among several adaptive optimization methods in federated learning (Reddi et al., 2020). FedAdam changes the global update rule of FedAvg from one-step SGD to one-step adaptive gradient optimization. Specifically, after gathering local differences Δ_t^i and averaging to Δ_t , the server updates the global model by Adam optimizer:

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \Delta_t; \quad (3.2)$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \Delta_t^2; \quad (3.3)$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t}}; \quad (3.4)$$

where Δ_t acts as pseudo gradient, and the global update can be viewed as one step Adam update using Δ_t . Several variants were also proposed will slight changes in the variance term \mathbf{v}_t , such as FedAdagrad and FedYogi (Reddi et al., 2020) and FedAMSGrad (Tong et al., 2020). Note that the \mathbf{v}_t in (3.4) is used for numerical stabilization purpose as the \mathbf{v}_t term can be quite small and cause unstable optimization behaviours.

3.1. Federated AMSGrad with Max Stabilization

In this section, we propose a general adaptive federated optimization framework, **Federated AMSGrad with Max Stabilization** (FedAMS), where the server conducts one additional max stabilization step before the final update.

Algorithm 1 summarize the details of general FedAMS framework. At the beginning of global round t , we first

¹The global update of FedAvg is equivalent to perform one step SGD update with the pseudo gradient Δ_t and learning rate $\eta = 1$.

select a subset of clients \mathcal{S}_t , each participating client $i \in \mathcal{S}_t$ obtains the local model $\mathbf{x}_{t,K}^i$ after K steps of local SGD updates with learning rate η . The model difference Δ_t^i is the difference between the local updated model $\mathbf{x}_{t,K}^i$ and the current global model \mathbf{x}_t , i.e., $\Delta_t^i = \mathbf{x}_{t,K}^i - \mathbf{x}_t$. The server aggregates Δ_t^i and gets the global difference Δ_t , this Δ_t acts as a pseudo gradient to calculate momentum \mathbf{m}_t and variance \mathbf{v}_t following (3.2) and (3.3). Now for updating \mathbf{x}_{t+1} , our general FedAMS framework provides two options for max stabilization:

$$\text{Option 1: } \mathbf{v}_t = \max(\mathbf{v}_{t-1}, \mathbf{v}_t); \mathbf{x}_{t+1} = \mathbf{x}_t + \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t}};$$

$$\text{Option 2: } \mathbf{v}_t = \max(\mathbf{v}_{t-1}, \mathbf{v}_t); \mathbf{x}_{t+1} = \mathbf{x}_t + \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t + \epsilon}};$$

Note that Option 2 is the same as the AMSGrad (Reddi et al., 2018) update rule, which brings a non-decreasing \mathbf{v}_t to solve a non-convergence issue in Adam (Kingma & Ba, 2014). For Option 1, FedAMS directly adopts $\sqrt{\mathbf{v}_t}$ as the denominator where ϵ is taken as the part of the max operation in \mathbf{v}_t . Intuitively, the unstable behaviour of the denominator (small value in the \mathbf{v}_t) usually only happens for a small set of dimensions. Therefore, the max stabilization strategy in Option 1 only affects those dimensions with small \mathbf{v}_t values, while the traditional adding strategy as in Option 2 will affect the accuracy on all dimensions.

Moreover, we want to emphasize that although the theoretical analysis in Reddi et al. (2020) assumes $\beta_1 = 0$ and only considers the impact of variance \mathbf{v}_t , thus the non-decreasing variance is not necessary for the analysis in Reddi et al. (2020). While the non-decreasing variance is indeed necessary for us to obtain the complete proof with a positive β_1 (see Appendix for details).

3.2. Federated Communication-Compressed AMSGrad

In order to reduce the communication costs between synchronization, we propose **Federated Communication-compressed AMSGrad with Max Stabilization** (FedCAMS), which is summarized in Algorithm 2. The main difference lies in that after the client i obtains the model differences Δ_t^i via local SGD, FedCAMS will compress Δ_t^i to $\hat{\Delta}_t^i$ via error feedback compression strategy, and then send $\hat{\Delta}_t^i$ to the central server. In details, at round t , the client i will apply the compressor on the summation of model differences Δ_t^i together with the cumulative compression error \mathbf{e}_t^i to obtain $\hat{\Delta}_t^i$. After that, the client will update term \mathbf{e}_{t+1}^i by calculating the new cumulative compression error, i.e., $\mathbf{e}_{t+1}^i = \Delta_t^i + \mathbf{e}_t^i - \hat{\Delta}_t^i$, which will be useful for next round's computation. The rest part of FedCAMS is similar to FedAMS: the server aggregates $\hat{\Delta}_t^i$ and obtains $\hat{\Delta}_t$, which will participate in the global update.

To summarize, FedCAMS is indeed a communication-

Algorithm 1 FedAMS

Input: initial point x_1 , local step size η , global stepsize, $\beta_1; \beta_2; \dots$.

- 1: $m_0 = 0, v_0 = 0$
- 2: for $t = 1$ to T do
- 3: Random sample a subset of clients
- 4: Server sends s_t to the subset S_t of clients
- 5: $x_{t,0}^i = x_t$
- 6: for each client $i \in S_t$ in parallel do
- 7: for $k = 0; \dots; K - 1$ do
- 8: Compute local stochastic gradient $g_{t,k}^i = \frac{1}{r} F_i(x_{t,k}^i; \xi_{t,k}^i)$
- 9: $x_{t,k+1}^i = x_{t,k}^i - \eta g_{t,k}^i$
- 10: end for
- 11: $\hat{x}_t = x_{t,K}^i - x_t$
- 12: end for
- 13: Server aggregates local update: $\hat{g}_t = \frac{1}{|S_t|} \sum_{i \in S_t} \hat{x}_t^i$
- 14: Update $m_t = \beta_1 m_{t-1} + (1 - \beta_1) \hat{g}_t$
- 15: Update $v_t = \beta_2 v_{t-1} + (1 - \beta_2) \hat{g}_t^2$
- //Option 1:
- 16: $b_t = \max(b_{t-1}; v_t)$, update $x_{t+1} = x_t + \frac{m_t}{b_t}$
- //Option 2:
- 17: $b_t = \max(b_{t-1}; v_t)$, update $x_{t+1} = x_t + \frac{m_t}{b_t + \epsilon}$
- 18: end for

efficient with the following features.

Error Feedback Compression: Although error-feedback strategy (Karimireddy et al., 2019; Stich et al., 2018; Stich & Karimireddy, 2019) has been widely used in various distributed learning settings, there is much less use of error-feedback in the federated settings, especially for adaptive federated optimization. Note that combining error-feedback with adaptive federated optimization is not a trivial task at all, instead, it is actually quite complicated. Specifically, the theoretical analysis of the adaptive gradient method in the typical nonconvex setting relies on the construction of the Lyapunov function. Compared to directly analyzing the model parameter, this causes extra difficulty as applying the error feedback strategy on the smoothness-expanded terms from the Lyapunov function will result in an accumulation of the compression error which leads to divergence. In our theoretical analysis, we have to modify the original construction of the Lyapunov function and introduce a new auxiliary sequence about the compression error which eliminates the accumulation of compression error. Unlike those direct compression strategies such as simple quantization or direct compression strategies (Haddadpour et al., 2021; Rejzizadeh et al., 2020; Jin et al., 2020) which usually require

²Similar divergence issue has also been discussed in Tang et al. (2021); Wang et al. (2022) in the distributed setting.

an unbiased compressor to work, error feedback allows for various biased compressors such as commonly used scaled sign compressor or top-k-compressors. Furthermore, the design of the error feedback strategy is well-known for reducing unnecessary compression error, which leads to a more precise model update.

Support for Partial Participation: Here we also make error feedback compatible with partial participation settings by keeping the stale cumulative compression error for clients who were not selected for the current round training (see Lines 14-16 in Algorithm 2). Such design makes FedCAMS more practical and communication efficient. Note that the default client sampling strategy in FedCAMS is to randomly select the participating clients (without replacement) in each round, i.e. $p_i = \frac{|S_t|}{n}$. This can be easily extended to the weighted sampling strategy with probability $p_i = w_i$, even with varying numbers of participating workers.

Algorithm 2 FedCAMS

Input: initial point x_1 , local step size η , global stepsize, $\beta_1; \beta_2; \dots$, compressor $C(\cdot)$.

- 1: $m_0 = 0, v_0 = 0; e_1 = 0$
- 2: for $t = 1$ to T do
- 3: Random sample a subset of clients
- 4: Server sends s_t to the subset S_t of clients
- 5: $x_{t,0}^i = x_t$
- 6: for each client $i \in S_t$ in parallel do
- 7: for $k = 0; \dots; K - 1$ do
- 8: Compute local stochastic gradient $g_{t,k}^i = \frac{1}{r} F_i(x_{t,k}^i; \xi_{t,k}^i)$
- 9: $x_{t,k+1}^i = x_{t,k}^i - \eta g_{t,k}^i$
- 10: end for
- 11: $\hat{x}_t = x_{t,K}^i - x_t$
- 12: Compress $b_t^i = C(\hat{x}_t + e_t^i)$, send b_t^i to the server and update $e_{t+1}^i = \hat{x}_t + e_t^i - b_t^i$
- 13: end for
- 14: for each client $j \notin S_t$ in parallel do
- 15: client j maintains the stale compression error $e_{t+1}^j = e_t^j$
- 16: end for
- 17: Server aggregates local update: $\hat{g}_t = \frac{1}{|S_t|} \sum_{i \in S_t} b_t^i$
- 18: Server updates s_{t+1} using b_t in the same way as in Algorithm 1 (Line 14-17)
- 19: end for

4. Convergence Analysis

In this section, we present the theoretical convergence results of our proposed FedAMS and FedCAMS in Algorithm 1 and 2. We first introduce some assumptions needed for the proof.

Assumption 4.1 (Smoothness) Each loss function on the i -th worker $F_i(x)$ is L -smooth, i.e., $\|F_i(x) - F_i(y)\| \leq L\|x - y\|$.

$$\|F_i(x) - F_i(y)\| \leq L\|x - y\| + \frac{L}{2}\|x - y\|^2$$

This also implies the L -gradient Lipschitz condition, i.e., $\|F_i(x) - F_i(y) - L(x - y)\| \leq \frac{L}{2}\|x - y\|^2$. Assumption 4.1 is a standard assumption in nonconvex optimization problems, which has been also adopted in Kingma & Ba (2014); Reddi et al. (2018); Li et al. (2019b); Yang et al. (2021).

Assumption 4.2 (Bounded Gradient) Each loss function on the i -th worker $F_i(x)$ has G -bounded stochastic gradient, i.e., for all x , we have $\|g_i(x)\| \leq G$.

The assumption of bounded gradient is usually adopted in adaptive gradient methods (Kingma & Ba, 2014; Reddi et al., 2018; Zhou et al., 2018; Chen et al., 2020a)

Assumption 4.3 (Bounded Variance) Each stochastic gradient on the i -th worker has a bounded local variance, i.e., for all x ; $i \in [m]$, we have $\mathbb{E} \|g_i(x) - \nabla F_i(x)\|^2 \leq \frac{\sigma_i^2}{m}$, and the loss function on each worker has a global variance bound, $\frac{1}{m} \sum_{i=1}^m \mathbb{E} \|g_i(x) - \nabla f(x)\|^2 \leq \frac{\sigma_g^2}{m}$.

Assumption 4.3 is widely used in federated optimization problems (Li et al., 2019b; Reddi et al., 2020; Yang et al., 2021). The bounded local variance represents the randomness of stochastic gradients, and the bounded global variance represents data heterogeneity between clients. Note that $\sigma_g = 0$ corresponds to the i.i.d setting, in which datasets from each client have the same distribution.

In the following, we will show the convergence results of FedAMS³ and FedCAMS.

4.1. Convergence Analysis for FedAMS

Full Participation: For the full participation scheme, all workers participate in the communication rounds and model update, i.e., $|S_t| = m$; $\forall t \in [1; T]$.

Theorem 4.4. Under Assumptions 4.1-4.3, if the local learning rate η satisfies the following condition:

$$\eta \leq \min \left\{ \frac{1}{8KL}, \frac{1}{K} \frac{1}{2K^2G^2 + [(3 + C_f^2)L + 2]G} \right\},$$

then the iterates of FedAMS in Algorithm 1 under full participation scheme satisfy

$$\min_{t \in [2; T]} \mathbb{E} \|g(x_t)\|^2 \leq \frac{q}{4} \frac{1}{2K^2G^2 + \frac{f_0 - f}{\eta KT}} + \frac{1}{T} + \frac{\sigma_g^2}{m}; \quad (4.1)$$

where $q = \frac{C_f G^2 d}{m} + \frac{2C_f^2 \eta KLG^2 d}{m}$, $\sigma_g^2 = \frac{5}{m} \frac{K^2 L^2}{2} (\eta^2 + \frac{\sigma_g^2}{m})$

³For simplicity, we will only present the convergence guarantee with Option 1. Note that the theoretical analysis can be easily extended to Option 2 with constant-only changes.

$$\frac{1}{1 - \eta} \left[\frac{1}{2} (3 + C_f^2)L + 2G \right] \frac{1}{2m} \frac{1}{\eta} \frac{1}{1 - \eta} \text{ and } C_1 = \frac{1}{1 - \eta}.$$

Remark 4.5. The upper bound $\min_{t \in [2; T]} \mathbb{E} \|g(x_t)\|^2$ contains three parts: the first two items that directly related to the total number of steps T are vanishing as $T \rightarrow \infty$. The last term in (4.1) relates to the local stochastic variance σ_i and global variance σ_g . In the i.i.d setting where each worker has the same data distribution, we have zero global variance, i.e., $\sigma_g = 0$, and the variance term will be smaller and less dependent on the number of local steps.

Corollary 4.6. Suppose we choose the global learning rate $\eta = \frac{1}{K}$ and local learning rate $\eta_i = \frac{1}{TK}$, when T is sufficient large, i.e., $T \gg Km$, the convergence rate for FedAMS in Algorithm 1 under full participation scheme satisfies

$$\min_{t \in [2; T]} \mathbb{E} \|g(x_t)\|^2 = O \left(\frac{1}{TKm} \right); \quad (4.2)$$

Remark 4.7. Corollary 4.6 suggests that with sufficient large T , FedAMS achieves a convergence rate $O(\frac{1}{TKm})$, which matches the result for general federated non-convex optimization methods such as SCAFFOLD (Karimireddy et al., 2020), FedAdam (Reddi et al., 2020).

Remark 4.8. Note that compared with FedAdam (Reddi et al., 2020), our theoretical analysis on FedAMS makes improvements in completing the proof. The analysis in Reddi et al. (2020) can only consider the case where $\eta = 0$ which largely simplifies the proof in their theoretical analysis, while we provide a full analysis on FedAMS with non-zero momentum term.

Partial Participation: In the partial participation scheme, we assume that only m workers participate the local update and communicate with the central server on each step, i.e., $|S_t| = m$; $\forall t \in [1; T]$. 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 m workers for local updating, for any two workers $i, j \in S_t$, the probability of being sampled to participate in model update are $\mathbb{P}(i \in S_t) = \frac{m}{n}$ and $\mathbb{P}(i, j \in S_t) = \frac{m(m-1)}{n(n-1)}$.

Theorem 4.9. Under Assumption 4.1-4.3, if the local learning rate η satisfies: $\eta \leq \min \left\{ \frac{1}{8KL}, \frac{1}{K} \frac{1}{2K^2G^2 + [(3 + C_f^2)L + 2]G} \right\}$, then the iterates of FedAMS in Algorithm 1 under partial participation scheme satisfy

$$\min_{t \in [2; T]} \mathbb{E} \|g(x_t)\|^2 \leq \frac{q}{8} \frac{1}{2K^2G^2 + \frac{f_0 - f}{\eta KT}} + \frac{1}{T} + \frac{\sigma_g^2}{m}; \quad (4.3)$$

where $\alpha = \frac{C_1 G^2 d}{p} + \frac{2C_1^2 \kappa L G^2 d}{p}$ and $\beta = \frac{5^2 \kappa L^2}{2} \left(\frac{1}{p} + 6K \frac{2}{g} \right) + [(3 + C_1^2) L + 2 \frac{p}{2(1-\alpha)} G] \frac{1}{2n} \frac{1}{p} + [(3 + C_1^2) L + 2 \frac{p}{2(1-\alpha)} G] \frac{1}{2n} \frac{(m-n)}{(m-1)} [15K^2 L^2 \frac{2}{2} \left(\frac{1}{p} + 6K \frac{2}{g} \right) + 3K \frac{2}{g}]$ and $C_1 = \frac{1}{1-\alpha}$.

Remark 4.10. The upper bound for $\min_{t \in [T]} E[kr f(x_t)k^2]$ of partial participation is similar to full participation case but with a larger variance term. This is due to the fact that random sampling of participating workers introduces an additional variance during sampling. In the setting where we have zero global variance, i.e. $\sigma_g = 0$, the variance term will get smaller and less dependent on the number of local steps K as well.

Corollary 4.11. Suppose we choose the global learning rate $\eta = \left(\frac{1}{\sqrt{Kn}} \right)$ and local learning rate $\eta_l = \left(\frac{1}{p \sqrt{TK}} \right)$, the convergence rate for FedAMS in Algorithm 1 under partial participation scheme without replacement sampling is

$$\min_{t \in [T]} E[kr f(x_t)k^2] = O\left(\frac{p}{\sqrt{TK}}\right) \quad (4.4)$$

Remark 4.12. Note that Corollary 4.11 suggests that the dominant term in (4.3) is $O\left(\frac{p}{\sqrt{TK}}\right)$, which directly relates to the global variance σ_g^2 . Such convergence rate is consistent with the partial participation result of FedAvg in the non i.i.d case in (Yang et al., 2021). It shows that the global variance has more impact on convergence behaviour in partial participation cases, especially in highly non i.i.d cases where σ_g is large. Corollary 4.11 also suggests that larger number of participating clients would accelerate the convergence. Note that although FedAdam (Reddi et al., 2020) did not provide explicit conclusions on the partial participation setting, its appendix introduced the necessary steps for analyzing such setting, which cannot imply such desired relationship with.

Remark 4.13. The impact of the number of local updates K is complicated. In partial participation settings, it shows that larger K slows down the convergence while full participation suggests the opposite. Similar slow-down result has also been mentioned in Li et al. (2019b), while some others (Stich, 2018; McMahan et al., 2017) showed that larger K would increase the convergence rate. We will leave this as future work.

4.2. Convergence Analysis for FedCAMS

Let us first introduce the assumption for the compressor.

Assumption 4.14 (Biased Compressor). Consider a biased operator $C: \mathbb{R}^d \rightarrow \mathbb{R}^d$: for $\delta \times 2 \mathbb{R}^d$, there exists constant $0 < q < 1$ such that

$$E[kC(x) - xk] \leq qkxk; \delta \times 2 \mathbb{R}^d:$$

Note that $q = 0$ leads to $C(x) = x$ which means no compression. Assumption 4.14 is a standard assumption for

biased compressors (Karimireddy et al., 2019; Stich et al., 2018; Tang et al., 2021). There are several widely used compressors satisfying 4.14 such as scaled-sign compressor and topk compressor.

Top-k (Stich et al., 2018): For $k \leq d$ and $\delta \times 2 \mathbb{R}^d$, the coordinate of x is ordered by the magnitude $|x_{(1)}| \geq |x_{(2)}| \geq \dots \geq |x_{(d)}|$. Denote e_1, e_2, \dots, e_d as standard unit basis vectors in \mathbb{R}^d . The compressor $C_{top}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ is defined as: $C_{top}(x) = \sum_{i=d-k+1}^d x_{(i)} e_{(i)}$.

Remark 4.15. Let us define the compression ratio as $r = k/d$. It can be shown that $C_{top}(x) - xk^2 \leq \frac{1}{k} \sum_{i=1}^{d-k} |x_{(i)}|^2$, and thus we have $q = 1 - r$.

Scaled sign (Karimireddy et al., 2019): For $k \leq d$ and $\delta \times 2 \mathbb{R}^d$, the compressor $C_{sign}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ is defined as

$$C_{sign}(x) = \frac{xk}{k_1} \text{sign}(x)$$

Remark 4.16. For scaled sign compressor, we have $kC_{sign}(x) - xk^2 = (1 - k \frac{xk^2}{d k_1^2}) kxk^2$, thus we have $q = 1 - k \frac{xk^2}{d k_1^2}$.

Next, we show the convergence analysis for FedCAMS. Due to the space limit, we only show the full participation setting and leave the partial participation setting in Appendix B.3.

Theorem 4.17. Under Assumptions 4.1-4.3 and 4.14, if the local learning rate η_l satisfies the following condition:

$$\eta_l \leq \min\left\{ \frac{1}{8KL}; \frac{p}{4 \frac{2(1+q^2)^3(1-q^2)^2 K^2 G^2 + \dots} \right\}, \text{ where } C_{i,q} = \frac{p}{4 \frac{2(1+q^2)^3(1-q^2)^2 K^2 G^2 + \dots}$$

$$\min_{t \in [T]} E[kr f(x_t)k^2] \leq \frac{f_0 - f^*}{\eta_l K T} + \frac{1}{T} + \dots$$

$$4 \frac{4(1+q^2)^3}{2(1-q^2)^2} \frac{1}{p} K^2 G^2 + \frac{f_0 - f^*}{\eta_l K T} + \frac{1}{T} + \dots \quad (4.5)$$

where $\alpha = \frac{C_1 G^2 d}{p} + \frac{2C_1^2 \kappa L G^2 d}{p}$ and $\beta = \frac{5^2 \kappa L^2}{2} \left(\frac{1}{p} + 6K \frac{2}{g} \right) + [(3 + 2C_2) L + 2 \frac{p}{2(1-\alpha)} G] \frac{1}{2m} \frac{1}{p} + [(3 + 2C_2) L + 2 \frac{p}{2(1-\alpha)} G] \frac{1}{2m} \frac{2q}{(1-\alpha)^2} + \frac{4q^2}{(1-q^2)^2}$.

Remark 4.18. The convergence rate in Theorem 4.17 contains three parts as well, the first two parts are related to total iterates T , and they vanish as T increases. The last term shows no direct dependency on T , but on local and global variances. In the i.i.d case where $\sigma_g = 0$, the variance will decrease and show less dependency on the number of local steps K .

Corollary 4.19. Suppose we choose the global learning rate $\eta = \left(\frac{1}{\sqrt{Km}} \right)$ and local learning rate $\eta_l = \left(\frac{1}{p \sqrt{TK}} \right)$, when T is sufficient large, i.e. $T \gg Km$, the convergence rate for FedCAMS in Algorithm 2 under full participation

scheme satisfies

$$\min_{\{x_t\}} \mathbb{E} \|\sum_{k=1}^K f(x_t) - \sum_{k=1}^K \bar{f}_k\|^2 = O\left(\frac{1}{TKm}\right) \quad (4.6)$$

Remark 4.20. Corollary 4.19 suggests that with sufficient large T , FedCAMS achieves the desired $O\left(\frac{1}{TKm}\right)$ convergence rate which matches the result for its uncompressed counterpart, FedAMS. This suggests that FedCAMS indeed achieves better communication efficiency without sacrificing much on the accuracy.

Remark 4.21. The constants C_1 and C_2 in Theorem 4.17 are related to the compression constant β . Specifically, if we track the dependency on the convergence rate, we have $O\left(\frac{1}{TKm}\right)$ under the full participation scheme.

A larger q ($q \neq 1$) corresponds to a stronger compression we applied, leading to worse convergence due to heavier information losses. Note that this dependency is common for the adaptive gradient method since the convergence proof of adaptive gradient methods heavily relies on the bounded gradient assumption and thus the compressed gradient bound is related to q . Similar type of dependency also occurs in other communication compressed distributed Adam methods such as Chen et al. (2021a).

5. Experiments

In this section, we present empirical validations toward the effectiveness of our proposed algorithms. Firstly, we provide comparisons between FedAMS and other first-order federated optimization baselines. Secondly, we provide experimental results of our proposed communication-efficient adaptive federated learning method, FedCAMS, to show its effectiveness in achieving communication-efficient adaptive federated learning.

Experimental Setup: We test all federated learning baselines, including ours on CIFAR10 and CIFAR100 datasets (Krizhevsky et al., 2009) using the following two models: (1) ResNet-18 (He et al., 2016), a widely used convolutional neural network model which is commonly trained by SGD; and (2) ConvMixer model (Trockman & Kolter, 2022), which shares similar ideas to vision transformer (Dosovitskiy et al., 2021) to use patch embeddings to preserve locality and similarly is trained via adaptive gradient methods by default. We set in total 100 clients for all federated training experiments. We set the partial participation ratio as α , i.e., in each round, the server picks α out of 100 clients to participate in the communication and model update. In each round, the client will perform β local epochs of local training with batch size 20. We search for the best training

hyper-parameters for each baseline, including ours. Due to the space limit, we leave all the hyper-parameter details as well as the CIFAR-100 experiments in the Appendix.

5.1. FedAMS and Adaptive Federated Optimization

We compare two options of the FedAMS framework with several state-of-the-art adaptive federated learning optimization methods, including: (1) FedAdam (Reddi et al., 2020) (2) FedYogi (Reddi et al., 2020) as well as standard federated baselines: (3) FedAvg (McMahan et al., 2017). Note that the Option 2 for FedAMS is same as FedAMSGrad (Tong et al., 2020). Thus in this section, we denote FedAMS for Option 1 and FedAMSGrad for Option 2 in the general FedAMS framework.

Figure 1 shows the convergence result of FedAMS and other federated learning baselines on training CIFAR-10 dataset with ResNet-18 model and ConvMixer-256-8 model. We compare the training loss and test accuracy against global rounds for each model. For the ResNet-18 model, FedAMS and FedYogi achieve quite similar performances, which are significantly better than the other three baselines. In particular, FedAMS performs the best in terms of the final training loss and test accuracy. On the other hand, FedAMSGrad and FedAdam obtain quite similar results on test accuracy and training loss. FedAvg achieves a slightly better training loss to FedAdam and FedAMSGrad but much higher test accuracy which is close to FedYogi and FedAMS. For the ConvMixer-256-8 model, which is typically trained via adaptive gradient methods, we observe that all adaptive federated optimization methods (FedAdam, FedYogi, FedAMSGrad and FedAMS) achieve much better performance in terms of both training loss and test accuracy than FedAvg. In detail, FedAMS again achieves a significantly better result than other baselines. Other adaptive methods, including FedAdam, FedYogi, and FedAMSGrad, have similar convergence behaviour when training the ConvMixer-256-8 model. Such results empirically show the effectiveness of our proposed FedAMS method with max stabilization.

Figure 2 shows the effect of parameter α on the convergence rate by choosing different number of α from $\{5, 10, 20\}$. From Figure 2 we can observe that a larger number of participating clients in general achieves a faster convergence rate. This verified our theoretical results in Section 4.1.

Figure 3 shows the ablation study with different local epochs by choosing different number of local epochs from $\{3, 10, 30, 100\}$ when training CIFAR-10 data on ResNet-18 with FedAMS optimizer. We observe that large β leads to faster convergence, but large β does not show a significant advantage in achieving a higher test accuracy. We follow FedAvg (McMahan et al., 2017) and FedAdam (Reddi et al., 2020) and set local epoch $\beta = 3$ by default unless otherwise specified.

⁴Note that the q -dependency in 1-bit Adam (Tang et al., 2021) is actually different due to the use of variance-freedom Adam update, i.e., freeze the variance term of Adam update as a constant after a few epochs, which make it resembles momentum SGD.

(a) ResNet-18 (b) ResNet-18

(c) ConvMixer-256-8 (d) ConvMixer-256-8

Figure 1. The learning curves for FedAMS and other federated learning baselines on training CIFAR-10 data (a)(b) show the results for the ResNet-18 model and (c)(d) show the results for the ConvMixer-256-8 model.

(a) ResNet-18 (b) ConvMixer-256-8

Figure 2. The learning curves for FedAMS with different participating number of clients in training CIFAR-10 data on the ResNet-18 and the ConvMixer-256-8 models.

(a) Training Loss (b) Test Accuracy

Figure 3. The learning curves for FedAMS with different number of local epochs in training CIFAR-10 data on the ResNet-18 model.

5.2. Communication-Eficient FedCAMS

Figure 4 shows the convergence results of FedAMS and FedCAMS with different compression strategies on training CIFAR-10 dataset with the ResNet-18 model. It includes comparisons between scaled sign compressor and top-k compressor with compression ratio $r = 1/64, 1/128, 1/256$. We compare the training loss and test accuracy against global rounds and the (pseudo) gradient communication bits. FedAMS, who does not conduct any communication compression, performs the best in terms of training loss yet requires a large volume of communication costs. For our FedCAMS compression methods, sign compressor and top-k compressor with ratio $r = 1/64$ achieve similar performance in terms of test accuracy against the training rounds and obtain the best trade-off between communication efficiency and model accuracy. Figure 4 (b)(d) show the direct comparison against the communication bits of training ResNet-18 on CIFAR-10. In particular, we can observe that the top-k compressor with a smaller r (i.e., a heavier compression with more information lost), obtains better communication efficiency but a slower convergence rate. Note that for a d -dimensional vector, the overall cost of a scaled sign compressor is d bits, and it is roughly the same communication costs as a top-k compressor with a ratio $r = 1/64$. This verifies our theoretical results in Section 4.2.

Figure 5 shows the convergence results of FedAMS and FedCAMS with the same compression strategies as in Figure 4 on training CIFAR-10 dataset with the ConvMixer-256-8 model. We notice that FedCAMS with the scaled sign compressor achieves roughly the same training loss and test accuracy as FedAMS but with a few orders of magnitude less in communication costs, while other top-k compression models have significantly worse performance. Among the top-k compressor trained models, the one with compression ratio $r = 1/64$ still obtains better training loss and test accuracy but higher communication costs. These results suggest that our proposed FedCAMS is communication-efficient while maintaining high accuracy.

6. Conclusions and Future Work

In this paper, we propose a communication-efficient compressed federated adaptive gradient optimization framework, FedCAMS, which largely reduces the communication overhead and addresses the adaptivity issue in federated optimization methods. FedCAMS is based on our proposed gen-

⁵Here FedCAMS adopt Option 1 for the local update step for fairness comparisons (same as FedAMS).

⁶Note that here we only count the client-to-server one-way communications compression.

⁷Top-k compressor also needs to communicate about the chosen locations which roughly double the costs.

(a) Training Loss (b) Training Loss

(c) Test Accuracy (d) Test Accuracy

Figure 4. The learning curves for FedCAMS and uncompressed FedAMS on training CIFAR-10 data on the ResNet-18 model.

(a) Training Loss (b) Training Loss

(c) Test Accuracy (d) Test Accuracy

Figure 5. The learning curves for FedCAMS and uncompressed FedAMS on training CIFAR-10 data on ConvMixer-256-8 model.

eral adaptive federated optimization framework, FedAMS, which contains variants of FedAdam feature max stabilization mechanisms. We present an improved theoretical convergence analysis of adaptive federated optimization, based on which we prove that in the nonconvex stochastic optimization setting, our proposed FedCAMS achieves the same convergence rate as its uncompressed counterpart FedAMS with a few orders of magnitude less communication cost. Experiments on various benchmarks verified our theoretical

results.

Our current analysis is limited to one-way communication compression from clients to the central server. However, extending our current analysis to two-way communication compression is highly non-trivial as it can be hard to guarantee the distributed global model from server to clients to stay synchronized due to error feedback and the biased compressor, especially in the partial participation setting. We leave it as future work.

Acknowledgements

We thank the anonymous reviewers for their helpful comments. This research was supported in part by a Seed Grant award from the Institute for Computational and Data Sciences at the Pennsylvania State University as well as Dell Technology AI Infrastructure Level Technologies Grant. The views and conclusions contained in this paper are those of the authors and should not be interpreted as representing any funding agencies.

References

Alistarh, D., Grubic, D., Li, J., Tomioka, R., and Vojnovic, M. Qsgd: Communication-efficient sgd via gradient quantization and encoding. *Advances in Neural Information Processing Systems*, 30:1709–1720, 2017.

Basu, D., Data, D., Karakus, C., and Diggavi, S. Qsparse-local-sgd: Distributed sgd with quantization, sparsification, and local computations. *arXiv preprint arXiv:1906.02367* 2019.

Bernstein, J., Wang, Y.-X., Azzadenesheli, K., and Anandkumar, A. signsgd: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning* pp. 560–569. PMLR, 2018.

Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165* 2020.

Chen, C., Shen, L., Huang, H., and Liu, W. Quantized adam with error feedback. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2(5):1–26, 2021a.

Chen, J., Zhou, D., Tang, Y., Yang, Z., Cao, Y., and Gu, Q. Closing the generalization gap of adaptive gradient methods in training deep neural networks. *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)* 2020a.

Chen, M., Shlezinger, N., Poor, H. V., Eldar, Y. C., and Cui, S. Communication-efficient federated learning.

- ceedings of the National Academy of Sciences (17), 2021b.
- Chen, X., Liu, S., Sun, R., and Hong, M. On the convergence of a class of adam-type algorithms for non-convex optimization. *arXiv preprint arXiv:1808.02941*, 2018.
- Chen, X., Li, X., and Li, P. Toward communication efficient adaptive gradient method. *Proceedings of the 2020 ACM-IMS on Foundations of Data Science Conference*, pp. 119–128, 2020b.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. An image is worth 16x16 words: Transformers for image recognition at scale. *International Conference on Learning Representations*, 2021.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Ghosh, A., Hong, J., Yin, D., and Ramchandran, K. Robust federated learning in a heterogeneous environment. *arXiv preprint arXiv:1906.06629*, 2019.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial networks. *Advances in neural information processing systems*, 27, 2014.
- Haddadpour, F., Kamani, M. M., Mokhtari, A., and Mahdavi, M. Federated learning with compression: Unified analysis and sharp guarantees. *International Conference on Artificial Intelligence and Statistics*, pp. 2350–2358. PMLR, 2021.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Hsu, T.-M. H., Qi, H., and Brown, M. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.
- Jhunjhunwala, D., Gadhikar, A., Joshi, G., and Eldar, Y. C. Adaptive quantization of model updates for communication-efficient federated learning. *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 3110–3114. IEEE, 2021.
- Jin, R., Huang, Y., He, X., Dai, H., and Wu, T. Stochastic-sign sgd for federated learning with theoretical guarantees. *arXiv preprint arXiv:2002.10940*, 2020.
- Karimireddy, S. P., Rebjock, Q., Stich, S., and Jaggi, M. Error feedback ressignsgd and other gradient compression schemes. *International Conference on Machine Learning*, pp. 3252–3261. PMLR, 2019.
- Karimireddy, S. P., Kale, S., Mohri, M., Reddi, S., Stich, S., and Suresh, A. T. Scaffold: Stochastic controlled averaging for federated learning. *International Conference on Machine Learning*, pp. 5132–5143. PMLR, 2020.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Konecny, J., McMahan, H. B., Yu, F. X., Ricătk, P., Suresh, A. T., and Bacon, D. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Li, D. and Wang, J. Fedmd: Heterogeneous federated learning via model distillation. *arXiv preprint arXiv:1910.03581*, 2019.
- Li, T., Sanjabi, M., Beirami, A., and Smith, V. Fair resource allocation in federated learning. *arXiv preprint arXiv:1905.10497*, 2019a.
- Li, T., Sahu, A. K., Talwalkar, A., and Smith, V. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020.
- Li, X., Huang, K., Yang, W., Wang, S., and Zhang, Z. On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*, 2019b.
- Lin, T., Stich, S. U., Patel, K. K., and Jaggi, M. Don't use large mini-batches, use local sgd. *arXiv preprint arXiv:1808.07217*, 2018.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Luo, L., Xiong, Y., Liu, Y., and Sun, X. Adaptive gradient methods with dynamic bound of learning rate. *arXiv preprint arXiv:1902.09843*, 2019.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. Communication-efficient learning of deep networks from decentralized data. *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.

- Nishio, T. and Yonetani, R. Client selection for federated learning with heterogeneous resources in mobile edge. In ICC 2019-2019 IEEE International Conference on Communications (ICC) pp. 1–7. IEEE, 2019.
- Wang, J., Tantia, V., Ballas, N., and Rabbat, M. Slowmo: Improving communication-efficient distributed sgd with slow momentum. arXiv preprint arXiv:1910.00643 2019.
- Reddi, S., Charles, Z., Zaheer, M., Garrett, Z., Rush, K., Konecny, J., Kumar, S., and McMahan, H. B. Adaptive federated optimization arXiv preprint arXiv:2003.00295 2020.
- Wang, J., Liu, Q., Liang, H., Joshi, G., and Poor, H. V. Tackling the objective inconsistency problem in heterogeneous federated optimization arXiv preprint arXiv:2007.07481 2020.
- Reddi, S. J., Kale, S., and Kumar, S. On the convergence of adam and beyond. In International Conference on Learning Representation 2018.
- Wang, Y., Lin, L., and Chen, J. Communication-compressed adaptive gradient method for distributed nonconvex optimization. In International Conference on Artificial Intelligence and Statistics pp. 6292–6320. PMLR, 2022.
- Reisizadeh, A., Mokhtari, A., Hassani, H., Jadbabaie, A., and Pedarsani, R. Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization. In International Conference on Artificial Intelligence and Statistics pp. 2021–2031. PMLR, 2020.
- Yang, H., Fang, M., and Liu, J. Achieving linear speedup with partial worker participation in non-iid federated learning. arXiv preprint arXiv:2101.11203 2021.
- Robbins, H. and Monro, S. A stochastic approximation method. The annals of mathematical statistics pp. 400–407, 1951.
- Yang, Z., Chen, M., Saad, W., Hong, C. S., and Shikh-Bahaei, M. Energy efficient federated learning over wireless communication networks IEEE Transactions on Wireless Communications 20(3):1935–1949, 2020.
- Seide, F., Fu, H., Droppo, J., Li, G., and Yu, D. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns. Fifteenth annual conference of the international speech communication association Citeseer, 2014.
- Zeiler, M. D. Adadelta: an adaptive learning rate method. arXiv preprint arXiv:1212.5701 2012.
- Stich, S. U. Local sgd converges fast and communicates little. arXiv preprint arXiv:1805.09767 2018.
- Zhou, D., Chen, J., Cao, Y., Tang, Y., Yang, Z., and Gu, Q. On the convergence of adaptive gradient methods for non-convex optimization arXiv preprint arXiv:1808.05671 2018.
- Stich, S. U. and Karimireddy, S. P. The error-feedback framework: Better rates for sgd with delayed gradients and compressed communication. arXiv preprint arXiv:1909.05350 2019.
- Stich, S. U., Cordonnier, J.-B., and Jaggi, M. Sparse sgd with memory. arXiv preprint arXiv:1809.07599 2018.
- Tang, H., Gan, S., Awan, A. A., Rajbhandari, S., Li, C., Lian, X., Liu, J., Zhang, C., and He, Y. 1-bit adam: Communication efficient large-scale training with adam's convergence speed. arXiv preprint arXiv:2102.02888 2021.
- Tieleman, T., Hinton, G., et al. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. COURSE 6: Neural networks for machine learning (2):26–31, 2012.
- Tong, Q., Liang, G., and Bi, J. Effective federated adaptive gradient methods with non-iid decentralized data. arXiv preprint arXiv:2009.06557 2020.
- Trockman, A. and Kolter, J. Z. Patches are all you need? arXiv preprint arXiv:2201.09792 2022.

A. Proof in Section 4.1

A.1. Proof of Theorem 4.4

Similar to previous works studied adaptive methods (Chen et al., 2018; Zhou et al., 2018; Chen et al., 2020a), we introduce a Lyapunov sequence $\{z_t\}$: assume $x_0 = x_1$, for each $t \geq 1$, we have

$$z_t = x_t + \frac{1}{1+\nu_t}(x_t - x_{t-1}) = \frac{1}{1+\nu_t}x_t + \frac{\nu_t}{1+\nu_t}x_{t-1}. \quad (\text{A.1})$$

For the difference of sequence $\{z_t\}$, we have

$$\begin{aligned} z_{t+1} - z_t &= \frac{1}{1+\nu_{t+1}}(x_{t+1} - x_t) - \frac{1}{1+\nu_t}(x_t - x_{t-1}) \\ &= \frac{1}{1+\nu_{t+1}}(\nu_{t+1}^{1=2} m_t) - \frac{1}{1+\nu_t} \nu_t^{1=2} m_{t-1} \\ &= \frac{1}{1+\nu_{t+1}} \nu_{t+1}^{1=2} m_t + (1 - \frac{\nu_{t+1}}{1+\nu_{t+1}}) \frac{1}{1+\nu_t} \nu_t^{1=2} m_{t-1} \\ &= \nu_{t+1}^{1=2} m_t + \frac{1}{1+\nu_{t+1}} \nu_t^{1=2} \nu_t^{1=2} m_{t-1}. \end{aligned}$$

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

$$\begin{aligned} &E[f(z_{t+1})] - f(z_t) \\ &= E[\langle \nabla f(z_t); z_{t+1} - z_t \rangle] + \frac{L}{2} E[\|z_{t+1} - z_t\|^2] \\ &= E[\langle \nabla f(z_t); \nu_{t+1}^{1=2} m_t \rangle] + E[\langle \nabla f(z_t); \frac{1}{1+\nu_t} \nu_t^{1=2} \nu_t^{1=2} m_{t-1} \rangle] \\ &\quad + \frac{2L}{2} E[\nu_{t+1}^{1=2} m_t \langle \nabla f(z_t); \frac{1}{1+\nu_t} \nu_t^{1=2} \nu_t^{1=2} m_{t-1} \rangle] \\ &= E[\langle \nabla f(x_t); \nu_{t+1}^{1=2} m_t \rangle] + E[\langle \nabla f(z_t); \frac{1}{1+\nu_t} \nu_t^{1=2} \nu_t^{1=2} m_{t-1} \rangle] \\ &\quad + \frac{2L}{2} E[\nu_{t+1}^{1=2} m_t \langle \nabla f(z_t); \frac{1}{1+\nu_t} \nu_t^{1=2} \nu_t^{1=2} m_{t-1} \rangle] + E[\langle \nabla f(z_t) - \nabla f(x_t); \nu_{t+1}^{1=2} m_t \rangle]; \quad (\text{A.2}) \end{aligned}$$

here we recall the notation $\nu_t = \text{diag}(\mathbf{v}_t) = \text{diag}(\max(\mathbf{v}_{t-1}; \mathbf{v}_t))$.

Bounding I_1 : We have

$$\begin{aligned} I_1 &= E[\langle \nabla f(x_t); \nu_{t+1}^{1=2} m_t \rangle] \\ &= E[\langle \nabla f(x_t); \nu_{t+1}^{1=2} m_t \rangle] + E[\langle \nabla f(x_t); \nu_{t+1}^{1=2} m_t \rangle] \\ &= \frac{p}{2} E[\langle \nabla f(x_t); \nu_{t+1}^{1=2} m_t \rangle] + \frac{p}{2} E[\langle \nabla f(x_t); \nu_{t+1}^{1=2} m_t \rangle]; \quad (\text{A.3}) \end{aligned}$$

where the first inequality follows by the fact that $\frac{v_t+1}{2} \geq \frac{v_t}{2}$. For the second term in (A.3), we have

$$\begin{aligned} & \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \leq \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \\ & \leq \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \\ & \leq \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \\ & \leq \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \end{aligned} \quad (\text{A.4})$$

where the second inequality follows from Lemma C.1 and C.4, and we will further apply the bound $\mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right]$ following from Lemma C.5. For the first term in (A.3), we have

$$\begin{aligned} & \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \\ & = \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \\ & = \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \\ & = \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \\ & = \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \end{aligned} \quad (\text{A.5})$$

where the third equality follows the local update rule. For the last term in (A.5), we have

$$\begin{aligned} & \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \\ & = \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \\ & = \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \\ & = \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \\ & = \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \\ & = \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] \end{aligned} \quad (\text{A.6})$$

where the second equation follows from $\mathbb{E} \left[\sum_{t=0}^T \left\| \nabla f(x_t) \right\|^2 \right] = \frac{1}{2} \mathbb{E} \left[\sum_{t=0}^T \left(\left\| \nabla f(x_t) \right\|^2 + \left\| \nabla f(x_t) \right\|^2 \right) \right]$, the first inequality holds by applying Cauchy-Schwarz inequality, the second inequality follows from Assumption 4.1.

Hence by applying Lemma C.9 with the local learning rate condition: $\frac{1}{8KL}$, we have

$$\begin{aligned} & \frac{\rho}{2} \mathbb{E} \left[\frac{r f(x_t)}{2v_{t-1} + 1} \right]; \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^n \sum_{k=0}^{K-1} g_{t;k}^i + \frac{1}{m} \sum_{i=1}^n r F_i(x_t) \right] \\ & \frac{3\rho}{4} \frac{1}{K} \mathbb{E} \left[\frac{r f(x_t)}{2v_{t-1} + 1} \right]^2 + \frac{5}{2} \frac{3K^2 L^2}{\rho} \left(\frac{1}{m} + 6K \frac{1}{g} \right) \frac{\rho}{2} \frac{1}{K} \mathbb{E} \left[\frac{1}{2v_{t-1} + 1} \sum_{i=1}^n \sum_{k=0}^{K-1} r F_i(x_{t;k}^i) \right]^2; \end{aligned} \quad (\text{A.7})$$

Then merging pieces together, we have

$$\begin{aligned} I_1 & \frac{\rho}{4} \frac{1}{K} \mathbb{E} \left[\frac{r f(x_t)}{2v_{t-1} + 1} \right]^2 + \frac{5}{2} \frac{3K^2 L^2}{\rho} \left(\frac{1}{m} + 6K \frac{1}{g} \right) \\ & \frac{\rho}{2} \frac{1}{K} \mathbb{E} \left[\frac{1}{2v_{t-1} + 1} \sum_{i=1}^n \sum_{k=0}^{K-1} r F_i(x_{t;k}^i) \right]^2 + \frac{\rho}{2(1-\alpha)} G \mathbb{E}[k_t^2] \\ & \frac{1}{4} \mathbb{E} \left[\frac{r f(x_t)}{2v_{t-1} + 1} \right]^2 + \frac{5}{2} \frac{3K^2 L^2}{\rho} \left(\frac{1}{m} + 6K \frac{1}{g} \right) \\ & \frac{1}{2Km^2} \mathbb{E} \left[\frac{1}{2v_{t-1} + 1} \sum_{i=1}^n \sum_{k=0}^{K-1} r F_i(x_{t;k}^i) \right]^2 + \frac{\rho}{2(1-\alpha)} G \mathbb{E}[k_t^2]; \end{aligned} \quad (\text{A.8})$$

Bounding I_2 : The bound for I_2 mainly follows by the update rule and definition of virtual sequence

$$\begin{aligned} I_2 & \mathbb{E} \left[r f(z_t); \frac{1}{1-\alpha} \psi_t^{1=2} \psi_t^{1=2} m_{t-1} \right] \\ & \mathbb{E} \left[r f(z_t) - f(x_t) + r f(x_t); \frac{1}{1-\alpha} \psi_t^{1=2} \psi_t^{1=2} m_{t-1} \right] \\ & \mathbb{E} \left[\|r f(x_t)\| \frac{1}{1-\alpha} \psi_t^{1=2} \psi_t^{1=2} m_{t-1} \right] \\ & + 2LE \frac{1}{\psi_t^{1=2}} \frac{1}{1-\alpha} m_{t-1} \frac{1}{1-\alpha} \psi_t^{1=2} \psi_t^{1=2} m_{t-1} \\ & \frac{1}{1-\alpha} \frac{1}{1-\alpha} \frac{1}{K} G^2 \mathbb{E} \left[\psi_t^{1=2} \psi_t^{1=2} \right] + 2 \frac{1}{(1-\alpha)^2} L^2 K^2 G^2 \mathbb{E} \left[\psi_t^{1=2} \psi_t^{1=2} \right]; \end{aligned} \quad (\text{A.9})$$

where the last inequality holds by applying Lemma C.4 and the fact of $\psi_t^{1=2}$.

Bounding I_3 : It can be bounded as follows:

$$\begin{aligned} I_3 & \frac{2L}{2} \mathbb{E} \left[\psi_t^{1=2} m_{t-1} + \frac{1}{1-\alpha} \psi_t^{1=2} \psi_t^{1=2} m_{t-1} \right]^2 \\ & 2LE \mathbb{E} \left[\psi_t^{1=2} m_{t-1} \right]^2 + 2LE \frac{1}{1-\alpha} \mathbb{E} \left[\psi_t^{1=2} \psi_t^{1=2} m_{t-1} \right]^2 \\ & 2LE \mathbb{E} \left[\psi_t^{1=2} m_{t-1} \right]^2 + 2L \frac{1}{(1-\alpha)^2} K^2 G^2 \mathbb{E} \left[\psi_t^{1=2} \psi_t^{1=2} \right]^2; \end{aligned} \quad (\text{A.10})$$

where the first inequality follows by Cauchy-Schwarz inequality, and the second one follows by Lemma C.4.

Bounding I_4 :

$$\begin{aligned}
 I_4 &= \mathbb{E} \left[\sum_{t=1}^T \left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] \\
 &\leq \mathbb{E} \left[\sum_{t=1}^T \left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] \\
 &\leq \mathbb{E} \left[\sum_{t=1}^T \left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] \\
 &\leq \frac{2L}{2} \mathbb{E} \left[\sum_{t=1}^T \left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] + \frac{2L}{2} \mathbb{E} \left[\sum_{t=1}^T \frac{1}{1} \left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right];
 \end{aligned}$$

where the first inequality holds by the fact $\|a\|_{\Psi^{-1}} \leq \|a\|_{\Psi^{-1}}$; the second one follows from Assumption 4.1 and the third one holds by the definition of virtual sequence and the fact $\|a\|_{\Psi^{-1}} \leq \frac{1}{2} \|a\|_{\Psi^{-1}} + \frac{1}{2} \|a\|_{\Psi^{-1}}$. Then summing I_4 over $t = 1; \dots; T$, we have

$$\begin{aligned}
 \sum_{t=1}^T I_4 &\leq \frac{2L}{2} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] + \frac{2L}{2} \sum_{t=1}^T \mathbb{E} \left[\frac{1}{1} \left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] \\
 &\leq \frac{2L}{2} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] + \frac{2L}{2} \frac{1}{(1-1)^2} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right]; \tag{A.11}
 \end{aligned}$$

By Lemma C.7, we have

$$\sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] \leq \frac{TK}{m} \sum_{i=1}^n \frac{1}{m} \sum_{k=0}^{X-1} \mathbb{E} \left[\left\| \nabla F_i(x_{t;k}^i) \right\|_{\Psi_t^{-1}}^2 \right];$$

Therefore, the summation I_4 term is bounded by

$$\begin{aligned}
 \sum_{t=1}^T I_4 &\leq \frac{2L}{2} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] + \frac{1}{(1-1)^2} \frac{2L}{2} \frac{1}{m^2} \sum_{t=1}^T \mathbb{E} \left[\sum_{i=1}^n \sum_{k=0}^{X-1} \left\| \nabla F_i(x_{t;k}^i) \right\|_{\Psi_t^{-1}}^2 \right] \\
 &\quad + \frac{1}{(1-1)^2} \frac{2L}{2} \frac{TK}{m} \sum_{i=1}^n \frac{1}{m} \sum_{k=0}^{X-1} \mathbb{E} \left[\left\| \nabla F_i(x_{t;k}^i) \right\|_{\Psi_t^{-1}}^2 \right]. \tag{A.12}
 \end{aligned}$$

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

$$\begin{aligned}
 \mathbb{E} \left[f(z_{T+1}) \right] - f(z_1) &= \sum_{t=1}^T \left[I_1 + I_2 + I_3 + I_4 \right] \\
 &\leq \frac{1}{4} \sum_{t=1}^T \mathbb{E} \left[\frac{\left\| \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2}{2V_{t-1}} \right] + \frac{5}{2} \frac{3K^2 L^2 T}{2} \left(\frac{1}{1} + 6K \frac{1}{9} \right) + \frac{1}{2(1-1)^2} \frac{G}{G} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] \\
 &\quad + \frac{1}{2Km^2} \sum_{t=1}^T \mathbb{E} \left[\frac{1}{2V_{t-1}} \sum_{i=1}^n \sum_{k=0}^{X-1} \left\| \nabla F_i(x_{t;k}^i) \right\|_{\Psi_t^{-1}}^2 \right] \\
 &\quad + \frac{1}{1-1} \frac{1}{1} \frac{1}{1} \frac{1}{1} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] + \frac{1}{(1-1)^2} \frac{2}{2} \frac{2K^2 G^2}{2} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] \\
 &\quad + \frac{1}{(1-1)^2} \frac{2}{2} \frac{2K^2 LG^2}{2} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] + \frac{2L}{2} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] \\
 &\quad + \frac{2L}{2} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right] + \frac{2L}{2} \frac{1}{(1-1)^2} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla f(z_t) - \nabla f(x_t) \right\|_{\Psi_t^{-1}}^2 \right]; \tag{A.13}
 \end{aligned}$$

A.2. Proof of Corollary 4.6

If we pick $\eta = (\frac{p}{TK})$ and $\beta = (\frac{p}{Km})$, we have $\min_{t \in [T]} E[kr f(x_t)k^2] = O(\frac{1}{TKm})$.

A.3. Proof of Theorem 4.9

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, i.e., $\bar{w}_t = \frac{1}{n} \sum_{i \in S_t} w_i^t$. Denote $\bar{w}_t = \frac{1}{m} \sum_{i=1}^m w_i^t$, and for convenience, we follow the previous notation of $\bar{w}_t = \text{diag}(\bar{w}_t + \cdot)$. Next we show that the global model difference is an unbiased estimator of \bar{w}_t :

$$E_{S_t}[\bar{w}_t] = \frac{1}{n} E_{S_t}[\sum_{i=1}^n w_i^t] = E_{S_t}[\bar{w}_t] = \frac{1}{m} \sum_{i=1}^m w_i^t = \bar{w}_t: \quad (\text{A.17})$$

Define the virtual sequence as same as previous: assume $w_i = x_i$, for each $i = 1$, we have

$$z_t = x_t + \frac{1}{1-\beta} (x_t - x_{t-1}) = \frac{1}{1-\beta} x_t - \frac{\beta}{1-\beta} x_{t-1}; \quad (\text{A.18})$$

$$z_{t+1} - z_t = \bar{w}_t - \frac{1}{1-\beta} \bar{w}_t^{1=2} - \bar{w}_t^{1=2} - m_{t-1}; \quad (\text{A.19})$$

By Assumption 4.1, we have

$$\begin{aligned} & E[f(z_{t+1})] - f(z_t) \\ &= E[r f(z_t); \bar{w}_t^{1=2} - \frac{1}{1-\beta} \bar{w}_t^{1=2} - \bar{w}_t^{1=2} - m_{t-1}] \\ &+ \frac{2L}{2} E[\bar{w}_t^{1=2} - \frac{1}{1-\beta} \bar{w}_t^{1=2} - \bar{w}_t^{1=2} - m_{t-1}]^2 \\ &= E[r f(x_t); \bar{w}_t^{1=2} - \frac{1}{1-\beta} \bar{w}_t^{1=2} - \bar{w}_t^{1=2} - m_{t-1}] \\ &+ \frac{2L}{2} E[\bar{w}_t^{1=2} - \frac{1}{1-\beta} \bar{w}_t^{1=2} - \bar{w}_t^{1=2} - m_{t-1}]^2 + E[r f(z_t) - r f(x_t); \bar{w}_t^{1=2} - \frac{1}{1-\beta} \bar{w}_t^{1=2} - \bar{w}_t^{1=2} - m_{t-1}]; \quad (\text{A.20}) \end{aligned}$$

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

Note that the bound for I_2^0 is exactly the same as the bound for I_3^0 and I_4^0 which include the second-order momentum estimate of \bar{w}_t . For I_1^0 , we have

$$\begin{aligned} I_1^0 &= E[r f(x_t); \bar{w}_t - \frac{1}{1-\beta} \bar{w}_t - \bar{w}_t - m_{t-1}] \\ &= E[r f(x_t); \bar{w}_t - \frac{1}{1-\beta} \bar{w}_t - \bar{w}_t - m_{t-1}] \\ &= E[r f(x_t); \bar{w}_t - \frac{1}{1-\beta} \bar{w}_t - \bar{w}_t - m_{t-1}] + E[r f(x_t); \bar{w}_t - \frac{1}{1-\beta} \bar{w}_t - \bar{w}_t - m_{t-1}]; \quad (\text{A.21}) \end{aligned}$$

The first term in (A.21) does not change in partial participation scheme. The second term is changed due to the variance of \bar{w}_t changes. For the second term, we have

$$\frac{p}{2} E[r f(x_t); \bar{w}_t - \frac{1}{1-\beta} \bar{w}_t - \bar{w}_t - m_{t-1}] = \frac{p}{2(1-\beta)} G E[k\bar{w}_t k^2]; \quad (\text{A.22})$$

For I_3^0 , we have

$$\sum_{t=1}^T I_3^0 = \frac{2L}{2} \sum_{t=1}^T \mathbb{E}[k_t^2] + \frac{2L}{2} \frac{1}{(1-\rho)^2} \sum_{t=1}^T \mathbb{E}[\|\mathbf{v}_t^{1=2}\|_{\mathbf{K}^2 \mathbf{G}^2}^2]; \quad (\text{A.23})$$

and for I_4^0 , similar to (A.11), we have

$$\sum_{t=1}^T I_4^0 = \frac{2L}{2} \sum_{t=1}^T \mathbb{E}[k_t^2] + \frac{2L}{2} \frac{1}{(1-\rho)^2} \sum_{t=1}^T \mathbb{E}[k m_t k^2]; \quad (\text{A.24})$$

From Lemma C.8, we have

$$\sum_{t=1}^T \mathbb{E}[k m_t k^2] = \frac{KT}{n} \sum_{i=1}^n \sum_{k=0}^{S_t} \mathbb{E}[\|\mathbf{x}_{t,k}^i\|_{\mathbf{F}_i}^2]; \quad (\text{A.25})$$

Then substituting (A.25) into (A.24), we have

$$\begin{aligned} \sum_{t=1}^T I_4^0 &= \frac{2L}{2} \sum_{t=1}^T \mathbb{E}[k_t^2] + \frac{2L}{2} \frac{1}{(1-\rho)^2} \sum_{t=1}^T \mathbb{E}[\|\mathbf{x}_{t,k}^i\|_{\mathbf{F}_i}^2] + \frac{2L}{2} \frac{1}{(1-\rho)^2} \sum_{t=1}^T \frac{2L}{2n} \sum_{i=1}^n \sum_{k=0}^{S_t} \mathbb{E}[\|\mathbf{x}_{t,k}^i\|_{\mathbf{F}_i}^2] \\ &= \frac{2L}{2} \sum_{t=1}^T \mathbb{E}[k_t^2] + \frac{2L}{2} \frac{1}{(1-\rho)^2} \sum_{t=1}^T \mathbb{E}[\|\mathbf{x}_{t,k}^i\|_{\mathbf{F}_i}^2] + \frac{2L}{2} \frac{1}{(1-\rho)^2} \sum_{t=1}^T \frac{2L}{2n} \sum_{i=1}^n \sum_{k=0}^{S_t} \mathbb{E}[\|\mathbf{x}_{t,k}^i\|_{\mathbf{F}_i}^2] \\ &\quad + \frac{2L}{2} \frac{1}{(1-\rho)^2} \sum_{t=1}^T \frac{2L}{2n} \sum_{i=1}^n \sum_{k=0}^{S_t} \mathbb{E}[\|\mathbf{x}_{t,k}^i\|_{\mathbf{F}_i}^2]; \end{aligned} \quad (\text{A.26})$$

where we will further apply the bound for $\mathbb{E}[k_t^2]$ following by Lemma C.6. The second term in (A.26) can be bounded from (C.10). Therefore, summing up (A.22), (A.23) and (A.9), summing over from $t = 1$ to T , then adding (A.24), we have

$$\begin{aligned} \mathbb{E}[f(z_{T+1})] - f(z_1) &= \sum_{t=1}^T [I_1^0 + I_2 + I_3^0 + I_4^0] \\ &= \frac{1}{4} \sum_{t=1}^T \mathbb{E}[\|\mathbf{r}_t\|_{\mathbf{F}_t}^2] + \frac{5}{2} \sum_{t=1}^T \frac{3K^2 L^2 T}{2} (\rho^t + 6K^2 \rho^t) + \frac{\rho}{2(1-\rho)^2} \sum_{t=1}^T \mathbb{E}[k_t^2] \\ &\quad + \frac{1}{2Km^2} \sum_{t=1}^T \mathbb{E}[\|\mathbf{x}_{t,k}^i\|_{\mathbf{F}_i}^2] \\ &\quad + \frac{1}{1-\rho} \sum_{t=1}^T \mathbf{K} \mathbf{G}^2 \mathbb{E}[\|\mathbf{v}_t^{1=2}\|_{\mathbf{K}^2 \mathbf{G}^2}^2] + \frac{2L}{2} \frac{1}{(1-\rho)^2} \sum_{t=1}^T \mathbb{E}[\|\mathbf{v}_t^{1=2}\|_{\mathbf{K}^2 \mathbf{G}^2}^2] \\ &\quad + \frac{2L}{2} \frac{1}{(1-\rho)^2} \sum_{t=1}^T \mathbf{K}^2 \mathbf{L} \mathbf{G}^2 \mathbb{E}[\|\mathbf{v}_t^{1=2}\|_{\mathbf{K}^2 \mathbf{G}^2}^2] + \frac{2L}{2} \sum_{t=1}^T \mathbb{E}[\|\mathbf{v}_t^{1=2}\|_{\mathbf{K}^2 \mathbf{G}^2}^2] \\ &\quad + \frac{2L}{2} \sum_{t=1}^T \mathbb{E}[k_t^2] + \frac{2L}{2} \frac{1}{(1-\rho)^2} \sum_{t=1}^T \mathbb{E}[k m_t k^2]; \end{aligned} \quad (\text{A.27})$$

By applying Lemma C.5 into all terms containing the second moment estimate of model difference in (A.27), using the fact that $(\frac{1}{2} \mathbf{K}^2 \mathbf{G}^2 + \frac{1}{2} \mathbf{K}^2 \mathbf{G}^2)^{-1} \mathbf{K} \mathbf{K} \mathbf{K} = \frac{\rho}{2V_{t+1}} \mathbf{K} \mathbf{K} \mathbf{K}$, and applying Lemma C.8, we

have

$$\begin{aligned}
 & E[f(z_{T+1})] - f(z_1) \\
 & \leq \frac{\rho}{4} \frac{1}{2} \frac{K}{\rho^2 K^2 G^2 + K m^2} \sum_{t=1}^T E[\text{tr} f(x_t) k^2] + \frac{5}{\rho} \frac{3K^2 L^2 T}{2} (\frac{1}{\rho} + 6K \frac{2}{\rho}) \\
 & + \frac{1}{1} \frac{1}{1} \frac{1}{\rho} \frac{K G^2 d}{\rho} + \frac{2}{(1-\frac{1}{\rho})^2} \frac{2}{\rho} \frac{2}{\rho} \frac{2}{\rho} K^2 L G^2 d \\
 & + \frac{3}{2} \frac{2L}{2} + \frac{2}{2(1-\frac{1}{\rho})^2} 2L + \frac{\rho}{2(1-\frac{1}{\rho})^2} G \frac{K T}{n} \frac{2}{\rho} \sum_{t=1}^T E \sum_{i=1}^n \sum_{k=0}^{X^1} F_i(x_{t;k})^2 \\
 & \frac{\rho}{2} \frac{1}{2} \frac{1}{2} \frac{K^2 G^2 + K m^2} {2} \frac{3}{2} \frac{2L}{2} + \frac{2}{2(1-\frac{1}{\rho})^2} 2L + \frac{\rho}{2(1-\frac{1}{\rho})^2} G \frac{2}{\rho} \frac{(n-1)}{mn(m-1)} \\
 & + \frac{3}{2} \frac{2L}{2} + \frac{2}{2(1-\frac{1}{\rho})^2} 2L + \frac{\rho}{2(1-\frac{1}{\rho})^2} G \frac{2}{\rho} \frac{(m-n)}{mn(m-1)} 15mK^3 L^3 \frac{2}{\rho} (\frac{1}{\rho} + 6K \frac{2}{\rho}) T \\
 & + (90mK^4 L^2 \frac{2}{\rho} + 3mK^2) \sum_{t=1}^T E[\text{tr} f(x_t) k^2] + 3mK^2 T \frac{2}{\rho}; \tag{A.28}
 \end{aligned}$$

then we have

$$\begin{aligned}
 & E[f(z_{T+1})] - f(z_1) \\
 & \leq \frac{\rho}{4} \frac{1}{2} \frac{K}{\rho^2 K^2 G^2 + K m^2} \sum_{t=1}^T E[\text{tr} f(x_t) k^2] + \frac{5}{\rho} \frac{3K^2 L^2 T}{2} (\frac{1}{\rho} + 6K \frac{2}{\rho}) + \frac{1}{1} \frac{1}{1} \frac{1}{\rho} \frac{K G^2 d}{\rho} \\
 & + \frac{2}{(1-\frac{1}{\rho})^2} \frac{2}{\rho} \frac{2}{\rho} \frac{2}{\rho} K^2 L G^2 d + \frac{3}{2} \frac{2L}{2} + \frac{2}{2(1-\frac{1}{\rho})^2} 2L + \frac{\rho}{2(1-\frac{1}{\rho})^2} G \frac{K T}{n} \frac{2}{\rho} \sum_{t=1}^T \\
 & + \frac{3}{2} \frac{2L}{2} + \frac{2}{2(1-\frac{1}{\rho})^2} 2L + \frac{\rho}{2(1-\frac{1}{\rho})^2} G \frac{2}{\rho} \frac{(m-n)}{mn(m-1)} 15mK^3 L^3 \frac{2}{\rho} (\frac{1}{\rho} + 6K \frac{2}{\rho}) T \\
 & + (90mK^4 L^2 \frac{2}{\rho} + 3mK^2) \sum_{t=1}^T E[\text{tr} f(x_t) k^2] + 3mK^2 T \frac{2}{\rho}; \tag{A.29}
 \end{aligned}$$

By adopting additional constraint of local learning rate with the inequality $\frac{3}{2} \frac{2L}{2} + \frac{2}{2(1-\frac{1}{\rho})^2} 2L + \frac{\rho}{2(1-\frac{1}{\rho})^2} G \frac{2}{\rho} \frac{(n-1)}{mn(m-1)} \frac{\rho}{2} \frac{1}{2} \frac{K}{\rho^2 K^2 G^2 + K m^2} = 0$, thus we obtain the constraint ρ

$\frac{n(m-1)}{m(n-1)} \frac{\rho}{2} \frac{1}{2} \frac{K}{\rho^2 K^2 G^2 + K m^2} \frac{\rho}{2(1-\frac{1}{\rho})^2} G$, and we further need satisfies $\frac{\rho}{4} \frac{1}{2} \frac{K}{\rho^2 K^2 G^2 + K m^2} \frac{3}{2} \frac{2L}{2} + \frac{2}{2(1-\frac{1}{\rho})^2} 2L + \frac{\rho}{2(1-\frac{1}{\rho})^2} G \frac{2}{\rho} \frac{(m-n)}{mn(m-1)} (90mK^4 L^2 \frac{2}{\rho} + 3mK^2) \frac{\rho}{8} \frac{1}{2} \frac{K}{\rho^2 K^2 G^2 + K m^2}$. Hence we have the following condition on local learning rate ρ ,

$$\rho \leq \frac{n(m-1)}{48m(n-1)} K \frac{\rho}{2K^2 G^2 + K m^2} \frac{3L}{2} + \frac{2}{2(1-\frac{1}{\rho})^2} L + \frac{\rho}{2(1-\frac{1}{\rho})^2} G \frac{1}{n}; \tag{A.30}$$

then we have

$$\begin{aligned}
 & \frac{\rho}{8} \frac{1}{2} \frac{K}{\rho^2 K^2 G^2 + K m^2} \sum_{i=1}^T E[\text{tr} f(x_t) k^2] \\
 & \frac{f(z_0)}{T} E[f(z_T)] + \frac{C_1}{T} \frac{1}{\rho} \frac{K G^2 d}{\rho} + \frac{2C_1^2}{T} \frac{2}{\rho} \frac{2}{\rho} \frac{2}{\rho} K^2 L G^2 d \\
 & + \frac{5}{\rho} \frac{3K^2 L^2}{2} (\frac{1}{\rho} + 6K \frac{2}{\rho}) + \frac{3}{2} \frac{2L}{2} + \frac{2}{2(1-\frac{1}{\rho})^2} 2L + \frac{\rho}{2(1-\frac{1}{\rho})^2} G \frac{K}{n} \frac{2}{\rho} \sum_{t=1}^T \\
 & + \frac{3}{2} \frac{2L}{2} + \frac{2}{2(1-\frac{1}{\rho})^2} 2L + \frac{\rho}{2(1-\frac{1}{\rho})^2} G \frac{2}{\rho} \frac{(m-n)}{mn(m-1)} [15mK^3 L^2 \frac{2}{\rho} (\frac{1}{\rho} + 6K \frac{2}{\rho}) + 3mK^2 \frac{2}{\rho}]; \tag{A.31}
 \end{aligned}$$

Therefore,

$$\min E[\text{kr f}(x_t)k^2] \leq \frac{q}{8} \frac{1}{2} \frac{K^2 G^2}{L} + \frac{f_0}{K} \frac{1}{T} + \frac{1}{T} ; \quad (\text{A.32})$$

where $\frac{1}{T} = \frac{C_1 G^2 d + 2C_1^2 \frac{1}{L} K L G^2 d}{2} ; \frac{1}{T} = \frac{5}{2} \frac{K L^2}{L} (\frac{1}{L} + 6K \frac{2}{g}) + [(3 + C_1^2) L + 2 \frac{P}{2(1 - \frac{1}{2})} G] \frac{1}{2n} \frac{1}{L} + [(3 + C_1^2) L + 2 \frac{P}{2(1 - \frac{1}{2})} G] \frac{1}{2n} \frac{(m - n)}{(m - 1)} [15K^2 L^2 \frac{2}{L} (\frac{1}{L} + 6K \frac{2}{g}) + 3K \frac{2}{g}]$ and $C_1 = \frac{1}{1 - \frac{1}{2}}$.

A.4. Proof of Corollary 4.11

If we choose $\frac{1}{T} = (\frac{1}{P} \frac{1}{K})$ and $\frac{1}{T} = (\frac{P}{K n})$, we have $\min_{t \in [T]} E[\text{kr f}(x_t)k^2] = O(\frac{P}{K n})$.

B. Proof of Theorems in Section 4.2 and Partial Participation Setting for FedCAMS

B.1. Proof of Theorem 4.17

Notations and equations. From the update rule of Algorithm 2, we have $e_t = \frac{1}{m} \sum_{i=1}^m e_t^i$ and $m_t = (1 - \frac{1}{m}) \sum_{i=1}^m m_t^i + \frac{1}{m} \sum_{i=1}^m b_i$. Denote a global uncompressed difference $e_t = \frac{1}{m} \sum_{i=1}^m e_t^i$. Denote a virtual momentum sequence: $m_t^0 = (1 - \frac{1}{m}) m_{t-1}^0 + \frac{1}{m} \sum_{i=1}^m m_{t-1}^i$. By the aforementioned definition and notation, we have

$$b_t = \frac{1}{m} \sum_{i=1}^m (b_t^i - m_t^i) = \frac{1}{m} \sum_{i=1}^m (e_t^i - e_{t+1}^i) = e_t - e_{t+1}; \quad (\text{B.1})$$

Denote the weighted averaging error sequence $e_t = (1 - \frac{1}{m}) \sum_{i=1}^m e_t^i$, with the impute $e_1 = 0$, we obtain the relation between e_t and m_t as follows

$$m_t - m_t^0 = (1 - \frac{1}{m}) \sum_{i=1}^m (b_t^i - m_t^i) = (1 - \frac{1}{m}) \sum_{i=1}^m (e_t^i - e_{t+1}^i) = e_t - e_{t+1}; \quad (\text{B.2})$$

where the last step holds due to $m_{t+1} = (1 - \frac{1}{m}) \sum_{i=1}^m m_{t+1}^i + \frac{1}{m} \sum_{i=1}^m b_{t+1}^i = (1 - \frac{1}{m}) \sum_{i=1}^m m_{t+1}^i + \frac{1}{m} \sum_{i=1}^m e_{t+1}^i + \frac{1}{m} \sum_{i=1}^m e_{t+1}^i$.

Similar to previous works studied adaptive methods (Chen et al., 2020a; Zhou et al., 2018; Chen et al., 2018), we introduce a Lyapunov sequence ψ_t : assume $\psi_0 = x_1$, for each $t \geq 1$, we have

$$y_t = x_t + \frac{1}{1 - \frac{1}{m}} (x_t - x_{t-1}) = \frac{1}{1 - \frac{1}{m}} x_t - \frac{1}{1 - \frac{1}{m}} x_{t+1};$$

Therefore, by the update rule of (2), we have

$$\begin{aligned} y_{t+1} &= x_{t+1} + \frac{1}{1 - \frac{1}{m}} \psi_t^{1=2} m_t \\ &= x_{t+1} + \frac{1}{1 - \frac{1}{m}} \psi_t^{1=2} [m_t^0 + e_t - e_{t+1}] \\ &= x_{t+1} + \frac{1}{1 - \frac{1}{m}} \psi_t^{1=2} m_t^0 + \frac{1}{1 - \frac{1}{m}} \psi_t^{1=2} \frac{(1 - \frac{1}{m}) e_{t+1}}{1} - e_{t+1} \\ &= x_{t+1} + \frac{1}{1 - \frac{1}{m}} \psi_t^{1=2} m_t^0 + \psi_t^{1=2} e_{t+1} - \psi_t^{1=2} e_{t+1}; \end{aligned} \quad (\text{B.3})$$

The third equation holds due to the fact that $m_{t+1} = (1 - \frac{1}{m}) m_t + \frac{1}{m} \sum_{i=1}^m b_{t+1}^i$. We then introduce a new sequence based on the previous Lyapunov sequence ψ_t as follows

$$z_{t+1} = y_{t+1} + \psi_t^{1=2} e_{t+1} = x_{t+1} + \frac{1}{1 - \frac{1}{m}} \psi_t^{1=2} m_t^0 + \psi_t^{1=2} e_{t+1}; \quad (\text{B.4})$$

The sequence difference $z_{t+1} - z_t$ can be represented by

$$\begin{aligned} z_{t+1} - z_t &= x_{t+1} - x_t + \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} m_t^0 - \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} m_{t-1}^0 + \mathbf{v}_t^{1-2} x_{t+1} - \mathbf{v}_t^{1-2} x_t \\ &= \mathbf{v}_t^{1-2} m_t^0 + \mathbf{v}_t^{1-2} x_{t+1} + \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} m_t^0 - \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} m_{t-1}^0 - \mathbf{v}_t^{1-2} x_t; \end{aligned} \quad (\text{B.5})$$

where the second equation follows the update rule of x_t . Following (B.2), then combining likely terms and applying the definition of m_t^0 , we have

$$\begin{aligned} z_{t+1} - z_t &= \mathbf{v}_t^{1-2} m_t^0 + \mathbf{v}_t^{1-2} x_{t+1} + \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} m_t^0 - \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} m_{t-1}^0 - \mathbf{v}_t^{1-2} x_t \\ &= \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} m_t^0 - \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} m_{t-1}^0 + \mathbf{v}_t^{1-2} x_{t+1} - \mathbf{v}_t^{1-2} x_t \\ &= \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} [m_t^0 - m_{t-1}^0 + (1-\alpha) x_t] + \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} m_{t-1}^0 + \mathbf{v}_t^{1-2} x_{t+1} - \mathbf{v}_t^{1-2} x_t \\ &= \mathbf{v}_t^{1-2} x_t + \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} [\mathbf{v}_t^{1-2} m_{t-1}^0 - \mathbf{v}_t^{1-2} m_{t-1}^0] + \mathbf{v}_t^{1-2} x_{t+1} - \mathbf{v}_t^{1-2} x_t; \end{aligned} \quad (\text{B.6})$$

Therefore, we obtain a helpful Lyapunov sequence for our proof of FedCAMS. The proof of FedCAMS in full participation settings has a similar outline with the proof of FedAMS. By Assumption 4.1, we have

$$\begin{aligned} E[f(z_{t+1})] - f(z_t) &= E[\langle \nabla f(z_t); z_{t+1} - z_t \rangle] + \frac{L}{2} E[\|z_{t+1} - z_t\|^2] \\ &= E[\langle \nabla f(z_t); \mathbf{v}_t^{1-2} x_t \rangle] \\ &\quad + E[\langle \nabla f(z_t); \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} [\mathbf{v}_t^{1-2} m_{t-1}^0 - \mathbf{v}_t^{1-2} m_{t-1}^0] + \mathbf{v}_t^{1-2} x_{t+1} - \mathbf{v}_t^{1-2} x_t \rangle] \\ &+ \frac{2L}{2} E[\|\mathbf{v}_t^{1-2} x_t + \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} [\mathbf{v}_t^{1-2} m_{t-1}^0 - \mathbf{v}_t^{1-2} m_{t-1}^0] + \mathbf{v}_t^{1-2} x_{t+1} - \mathbf{v}_t^{1-2} x_t\|^2] \\ &= E[\langle \nabla f(x_t); \mathbf{v}_t^{1-2} x_t \rangle] + E[\langle \nabla f(z_t); \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} [\mathbf{v}_t^{1-2} m_{t-1}^0 - \mathbf{v}_t^{1-2} m_{t-1}^0] + \mathbf{v}_t^{1-2} x_{t+1} - \mathbf{v}_t^{1-2} x_t \rangle] \\ &\quad + \frac{2L}{2} E[\|\mathbf{v}_t^{1-2} x_t + \frac{1}{1-\alpha} \mathbf{v}_t^{1-2} [\mathbf{v}_t^{1-2} m_{t-1}^0 - \mathbf{v}_t^{1-2} m_{t-1}^0] + \mathbf{v}_t^{1-2} x_{t+1} - \mathbf{v}_t^{1-2} x_t\|^2] \\ &\quad + E[\langle \nabla f(z_t) - \nabla f(x_t); \mathbf{v}_t^{1-2} x_t \rangle]; \end{aligned} \quad (\text{B.7})$$

here we recall the notation $\mathbf{h}_t = \text{diag}(\mathbf{v}_t) = \text{diag}(\max(\mathbf{v}_{t-1}; v_t; \dots))$.

Bounding T_1 : We have

$$\begin{aligned} T_1 &= E[\langle \nabla f(x_t); \mathbf{p} \frac{\mathbf{v}_t}{\mathbf{v}_t + \mathbf{v}_t} \rangle] \\ &= \mathbf{p} \frac{1}{2} E[\langle \nabla f(x_t); \mathbf{p} \frac{\mathbf{v}_t}{2\mathbf{v}_t + \mathbf{v}_t} \rangle] + \mathbf{p} \frac{1}{2} E[\langle \nabla f(x_t); \mathbf{p} \frac{\mathbf{v}_t}{\mathbf{v}_t + \mathbf{v}_t} \rangle]; \end{aligned} \quad (\text{B.8})$$

where the first inequality follows by the fact that $\frac{v_t+1}{2}$. For the second term in (B.8), we have

$$\begin{aligned} & \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\|; \rho \frac{t}{v_t+1} \right] \leq \frac{\rho}{2v_{t-1}+1} \\ & \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \mathbb{E} \left[\frac{1}{v_t+1} \right] \right] \leq \frac{\rho}{2v_{t-1}+1} \mathbb{E} \left[\|\mathbf{k}_t\| \right] \\ & \leq \frac{\rho}{2(1-\beta)G} \mathbb{E} \left[\|\mathbf{k}_t\|^2 \right]; \end{aligned} \quad (\text{B.9})$$

where the second inequality follows from Lemma C.1 and C.4, and we will further apply the bound $\mathbb{E}[\|\mathbf{k}_t\|^2]$ by applying Lemma C.5. For the first term in (B.8), we have

$$\begin{aligned} & \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\|; \rho \frac{t}{2v_{t-1}+1} \right] \\ &= \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\|; \frac{t}{2v_{t-1}+1} \right] \leq \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \\ &= \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \\ &= \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \sum_{k=0}^{X-1} g_{t,k}^i \right] + \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \\ &= \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \sum_{k=0}^{X-1} g_{t,k}^i \right] + \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m r F_i(\mathbf{x}_t) \right] : \end{aligned} \quad (\text{B.10})$$

For the last term in (B.10), we have

$$\begin{aligned} & \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \sum_{k=0}^{X-1} g_{t,k}^i \right] + \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m r F_i(\mathbf{x}_t) \right] \\ &= \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \sum_{k=0}^{X-1} g_{t,k}^i \right] + \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \sum_{k=0}^{X-1} (r F_i(\mathbf{x}_{t,k}^i) - r F_i(\mathbf{x}_t)) \right] \\ &= \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \sum_{k=0}^{X-1} (r F_i(\mathbf{x}_{t,k}^i) - r F_i(\mathbf{x}_t)) \right]^2 \\ & \quad + \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \sum_{k=0}^{X-1} r F_i(\mathbf{x}_{t,k}^i) \right]^2 \\ &= \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\|^2 \right] + \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \sum_{k=0}^{X-1} (r F_i(\mathbf{x}_{t,k}^i) - r F_i(\mathbf{x}_t)) \right]^2 \\ & \quad + \frac{\rho}{2} \mathbb{E} \left[\|\mathbf{r}_f(\mathbf{x}_t)\| \right] \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \sum_{k=0}^{X-1} r F_i(\mathbf{x}_{t,k}^i) \right]^2 ; \end{aligned} \quad (\text{B.11})$$

where the second equation follows from $\|\mathbf{y}_i\| = \frac{1}{2} [\|\mathbf{k}_x\|^2 + \|\mathbf{k}_y\|^2 - \|\mathbf{k}_x - \mathbf{k}_y\|^2]$, and the inequality holds by applying

Cauchy-Schwarz inequality. Then by Assumption 4.1, we have

$$\begin{aligned}
 & \frac{\rho}{2} \mathbb{E} \left[\frac{\|r f(x_t)\|^2}{2V_{t-1}} \right]; \mathbb{E} \left[\frac{1}{m} \sum_{i=1}^m \sum_{k=0}^{K-1} g_{t;k}^i + \frac{1}{m} \sum_{i=1}^m r F_i(x_t) \right. \\
 & \left. \frac{\rho}{2} \frac{1}{K} \mathbb{E} \left[\frac{\|r f(x_t)\|^2}{2V_{t-1}} \right] + \frac{\rho}{2} \frac{1}{2m} \sum_{i=1}^m \sum_{k=0}^{K-1} \mathbb{E} \left[\frac{\|x_{t;k}^i - x_t\|^2}{2V_{t-1}} \right] \right. \\
 & \left. \frac{\rho}{2} \frac{1}{2Km^2} \mathbb{E} \left[\frac{1}{2V_{t-1}} \sum_{i=1}^m \sum_{k=0}^{K-1} r F_i(x_{t;k}^i) \right]^2 \right. \\
 & \left. \frac{3\rho}{4} \frac{1}{K} \mathbb{E} \left[\frac{\|r f(x_t)\|^2}{2V_{t-1}} \right] + \frac{5}{2} \frac{3K^2L^2}{\rho} (\frac{1}{L} + 6K \frac{2}{g}) \right. \\
 & \left. \frac{\rho}{2} \frac{1}{2Km^2} \mathbb{E} \left[\frac{1}{2V_{t-1}} \sum_{i=1}^m \sum_{k=0}^{K-1} r F_i(x_{t;k}^i) \right]^2 \right]; \tag{B.12}
 \end{aligned}$$

where the last inequality holds by applying Lemma C.9 and the constraint of local learning rate $\frac{1}{8KL}$. Then we have

$$\begin{aligned}
 T_1 & \frac{\rho}{4} \frac{1}{K} \mathbb{E} \left[\frac{\|r f(x_t)\|^2}{2V_{t-1}} \right] + \frac{5}{2} \frac{3K^2L^2}{\rho} (\frac{1}{L} + 6K \frac{2}{g}) \\
 & \frac{\rho}{2} \frac{1}{2Km^2} \mathbb{E} \left[\frac{1}{2V_{t-1}} \sum_{i=1}^m \sum_{k=0}^{K-1} r F_i(x_{t;k}^i) \right]^2 + \frac{\rho}{2(1-\alpha)} G \mathbb{E}[k_{t-1}^2] \\
 & \frac{1}{4} \mathbb{E} \left[\frac{\|r f(x_t)\|^2}{2V_{t-1}} \right] + \frac{5}{2} \frac{3K^2L^2}{\rho} (\frac{1}{L} + 6K \frac{2}{g}) \\
 & \frac{1}{2Km^2} \mathbb{E} \left[\frac{1}{2V_{t-1}} \sum_{i=1}^m \sum_{k=0}^{K-1} r F_i(x_{t;k}^i) \right]^2 + \frac{\rho}{2(1-\alpha)} G \mathbb{E}[k_{t-1}^2]; \tag{B.13}
 \end{aligned}$$

Bounding T_2 : The bound for T_2 mainly follows by the update rule and definition of virtual sequence

$$\begin{aligned}
 T_2 & = \mathbb{E} \left[r f(z_t); \frac{1}{1-\alpha} \psi_t^{1=2} \psi_t^{1=2} m_{t-1}^0 + \psi_t^{1=2} \psi_t^{1=2} t \right] \\
 & = \mathbb{E} \left[r f(z_t) + r f(z_t) - r f(x_t); \psi_t^{1=2} \psi_t^{1=2} \frac{1}{1-\alpha} m_{t-1}^0 + t \right] \\
 & \quad \mathbb{E} \left[\|r f(x_t)\| \psi_t^{1=2} \psi_t^{1=2} \frac{1}{1-\alpha} m_{t-1}^0 + t \right] \\
 & + 2LE \frac{1}{\psi_{t-1}^{1=2}} \frac{1}{1-\alpha} m_{t-1}^0 + t \psi_t^{1=2} \psi_t^{1=2} \frac{1}{1-\alpha} m_{t-1}^0 + t \\
 & C_1 \frac{1}{K} G^2 \mathbb{E} \left[\psi_t^{1=2} \psi_t^{1=2} \frac{1}{1-\alpha} \right] + 2C_1^2 L^2 K^2 G^2 \frac{1=2}{1-\alpha} \mathbb{E} \left[\psi_t^{1=2} \psi_t^{1=2} \frac{1}{1-\alpha} \right]; \tag{B.14}
 \end{aligned}$$

where the last inequality holds by Lemma C.4, here $\frac{1}{1-\alpha} + \frac{2q}{1-\alpha^2}$.

Bounding T_3 : It can be bounded as follows:

$$\begin{aligned}
 T_3 & = \frac{2L}{2} \mathbb{E} \left[\psi_t^{1=2} t + \frac{1}{1-\alpha} \psi_t^{1=2} \psi_t^{1=2} m_{t-1} \right]^2 \\
 & 2LE \psi_t^{1=2} t^2 + 2LE \frac{1}{1-\alpha} \psi_t^{1=2} \psi_t^{1=2} m_{t-1}^0 + \psi_t^{1=2} \psi_t^{1=2} t^2 \\
 & 2LE \psi_t^{1=2} t^2 + 2LC_1^2 \frac{2}{1-\alpha} K^2 G^2 \mathbb{E} \left[\psi_t^{1=2} \psi_t^{1=2} \right]; \tag{B.15}
 \end{aligned}$$

where the first inequality follows by Cauchy-Schwarz inequality, and the second one follows by Lemma C.4 here $\frac{1}{1-\epsilon} + \frac{2q}{1-q^2}$.

Bounding T_4 :

$$\begin{aligned} T_4 &= \mathbb{E} \sum_{t=1}^T \langle \nabla f(z_t) - \nabla f(x_t), \psi_t \rangle \\ &\leq \mathbb{E} \sum_{t=1}^T \|\nabla f(z_t) - \nabla f(x_t)\| \|\psi_t\| \\ &\leq L \mathbb{E} \sum_{t=1}^T \|z_t - x_t\| \|\psi_t\| \\ &\leq \frac{2L}{2} \mathbb{E} \sum_{t=1}^T \|\psi_t\|^2 + \frac{2L}{2} \mathbb{E} \sum_{t=1}^T \frac{1}{1-\epsilon} m_t + \|\psi_t\|^2; \end{aligned}$$

where the first inequality holds by the fact $\langle a, b \rangle \leq \|a\| \|b\|$, the second one follows from Assumption 4.1 and the third one holds by the definition of virtual sequence and the fact $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$. Then summing T_4 over $t = 1; \dots; T$, we have

$$\begin{aligned} \sum_{t=1}^T T_4 &\leq \frac{2L}{2} \sum_{t=1}^T \mathbb{E} \|\psi_t\|^2 + \frac{2L}{2} \sum_{t=1}^T \mathbb{E} \left[\frac{1}{1-\epsilon} m_t + \|\psi_t\|^2 \right] \\ &\leq \frac{2L}{2} \sum_{t=1}^T \mathbb{E}[k_t^2] + \frac{2L}{(1-\epsilon)^2} \sum_{t=1}^T \mathbb{E}[m_t k^2] + \sum_{t=1}^T \mathbb{E}[k_t^2]; \end{aligned} \quad (B.16)$$

By Lemma C.7, we have

$$\sum_{t=1}^T \mathbb{E}[k_t k^2] \leq \frac{TK}{m} \sum_{i=1}^2 \mathbb{E} \sum_{k=0}^{X_i} \mathbb{E} \|\nabla F_i(x_{t;k}^i)\|^2;$$

and

$$\sum_{t=1}^T \mathbb{E}[k_t^2] \leq \frac{4Tq^2}{(1-q^2)^2} \frac{K}{m} \sum_{i=1}^2 \mathbb{E} \sum_{k=0}^{X_i} \mathbb{E} \|\nabla F_i(x_{t;k}^i)\|^2;$$

Therefore, the T_4 term is bounded by

$$\sum_{t=1}^T T_4 \leq \frac{2L}{2} \sum_{t=1}^T \mathbb{E}[k_t^2] + \frac{C_2}{m^2} \sum_{t=1}^T \mathbb{E} \sum_{i=1}^2 \mathbb{E} \|\nabla F_i(x_{t;k}^i)\|^2 + \frac{C_2}{m} \sum_{i=1}^2 \mathbb{E} \sum_{k=0}^{X_i} \mathbb{E} \|\nabla F_i(x_{t;k}^i)\|^2; \quad (B.17)$$

where $C_2 = \frac{4q^2}{(1-q^2)^2} + \frac{2}{(1-\epsilon)^2}$.

Merging pieces together: Substituting (B.13), (B.14) and (B.15) into (B.7), summing over $t = 1$ to T and then adding

(B.17), we have

$$\begin{aligned}
 E[f(z_{T+1})] - f(z_1) &= \sum_{t=1}^T [T_1 + T_2 + T_3 + T_4] \\
 &= \frac{1}{4} K \sum_{t=1}^T E \left[\frac{r f(x_t)}{2V_{t-1}} \right]^2 + \frac{5}{2} \frac{3K^2 L^2 T}{p} \left(\frac{1}{l} + 6K \frac{2}{g} \right) + \frac{p}{2(1-\alpha)} G \sum_{t=1}^T E[k_t k^2] \\
 &+ \frac{1}{2Km^2} \sum_{t=1}^T E \left[\frac{1}{2V_{t-1}} \sum_{i=1}^n \sum_{k=0}^{K-1} r F_i(x_{t;k}^i) \right]^2 + C_1 \frac{1}{p} K G^2 d \sum_{t=1}^T E \left[\psi_t^{1=2} \psi_t^{1=2} \right] \\
 &+ \frac{C_1^2}{p} \frac{2}{l} K^2 G^2 \sum_{t=1}^T E \left[\psi_t^{1=2} \psi_t^{1=2} \right] + C_1^2 \frac{2}{l} K^2 L G^2 \sum_{t=1}^T E \left[\psi_t^{1=2} \psi_t^{1=2} \right] \\
 &+ \frac{2L}{2} \sum_{t=1}^T E \left[\psi_t^{1=2} \right]^2 + \frac{2L}{2} \sum_{t=1}^T E \left[\psi_t^{1=2} \right]^2 + \frac{2L}{2} \frac{1}{(1-\alpha)^2} \sum_{t=1}^T E[k_t^0 k^2] + \frac{2L}{2} \sum_{t=1}^T E[k_t k^2].
 \end{aligned} \tag{B.18}$$

Hence by organizing and applying Lemmas, we have

$$\begin{aligned}
 E[f(z_{T+1})] - f(z_1) &= \frac{1}{4} K \sum_{t=1}^T E \left[\frac{r f(x_t)}{2V_{t-1}} \right]^2 + \frac{5}{2} \frac{3K^2 L^2 T}{p} \left(\frac{1}{l} + 6K \frac{2}{g} \right) \\
 &+ \frac{1}{2Km^2} \sum_{t=1}^T E \left[\frac{1}{2V_{t-1}} \sum_{i=1}^n \sum_{k=0}^{K-1} r F_i(x_{t;k}^i) \right]^2 + \frac{C_1}{p} K G^2 d + \frac{2C_1^2}{l} \frac{2}{l} K^2 L G^2 d \\
 &+ \frac{2L}{2} + \frac{2L}{2} + \frac{p}{2(1-\alpha)} G \frac{KT}{m} \frac{1}{l} + \frac{1}{m^2} \sum_{t=1}^T E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} r F_i(x_{t;k}^i) \right]^2 \\
 &+ \frac{2L}{m^2} \frac{2C_1}{l} \sum_{t=1}^T E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} r F_i(x_{t;k}^i) \right]^2 + \frac{2L}{m} \frac{TK}{l} \frac{2}{l} C_2 \frac{1}{l};
 \end{aligned} \tag{B.19}$$

by applying Lemma C.5 into all terms containing the second moment estimate of model difference (B.18), using the fact that $\frac{(1+q^2)^3}{2(1-q^2)^2} K^2 G^2 + \frac{1}{K} k k$ $\frac{(1+q^2)^3}{2(1-q^2)^2} \frac{1}{l} K^2 G^2 + \frac{1}{K} k k$ $\frac{1}{2V_{t-1}} \frac{1}{l} k k$, and applying Lemma C.2 and C.8, we have

$$\begin{aligned}
 E[f(z_{T+1})] - f(z_1) &= \frac{1}{4} \frac{K}{4} \frac{(1+q^2)^3}{2(1-q^2)^2} \frac{1}{l} K^2 G^2 + \sum_{t=1}^T E[kr f(x_t)k^2] + \frac{5}{2} \frac{3K^2 L^2 T}{p} \left(\frac{1}{l} + 6K \frac{2}{g} \right) \\
 &+ \frac{C_1}{p} K G^2 d + \frac{2C_1^2}{l} \frac{2}{l} K^2 L G^2 d + \frac{3}{2} \frac{2L}{2} + C_2 \frac{2L}{2} + \frac{p}{2(1-\alpha)} G \frac{KT}{m} \frac{1}{l} \\
 &+ \frac{1}{2} \sum_{t=1}^T E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} r F_i(x_{t;k}^i) \right]^2 + \frac{1}{2} \frac{1}{4} \frac{(1+q^2)^3}{2(1-q^2)^2} \frac{1}{l} K^2 G^2 + \frac{1}{Km^2} \frac{3}{2} \frac{2L}{2} + C_2 \frac{2L}{2} + \frac{p}{2(1-\alpha)} G \frac{1}{m^2} \\
 &+ \frac{1}{4C_0} \sum_{t=1}^T E[kr f(x_t)k^2] + \frac{5}{2} \frac{3K^2 L^2 T}{p} \left(\frac{1}{l} + 6K \frac{2}{g} \right) \\
 &+ \frac{C_1}{p} K G^2 d + \frac{2C_1^2}{l} \frac{2}{l} K^2 L G^2 d + \frac{3}{2} \frac{2L}{2} + C_2 \frac{2L}{2} + \frac{p}{2(1-\alpha)} G \frac{KT}{m} \frac{1}{l};
 \end{aligned} \tag{B.20}$$

where the last inequality holds by $\frac{p}{4} \frac{K}{2(1+q^2)^3(1-q^2)^2 K^2 G^2 + K(3L+2C_2L+2\frac{p}{2(1-q^2)}G)}$. Hence we have

$$\begin{aligned} & -q \frac{K}{4} \frac{1}{2(1+q^2)^3(1-q^2)^2 K^2 G^2 + K(3L+2C_2L+2\frac{p}{2(1-q^2)}G)} \sum_{t=1}^T E[kr f(x_t)k^2] \\ & \frac{f(z_0) - E[f(z_T)]}{T} + \frac{5}{2} \frac{K^2 L^2}{p} (\frac{1}{2} + 6K \frac{2}{g}) + \frac{C_1}{T} \frac{K G^2 d}{p} + \frac{2C_1^2}{T} \frac{K^2 L G^2 d}{p} \\ & + 3 \frac{2L+2C_2}{2} L + 2 \frac{p}{2(1-q^2)} G \frac{K}{2m} \frac{1}{l^2}; \end{aligned} \quad (B.21)$$

where $C_1 = \frac{1}{1-q^2} + \frac{2q}{1-q^2}$ and $C_2 = \frac{1}{(1-q^2)^2} + \frac{4q^2}{(1-q^2)^2}$. (B.21) also implies,

$$\min E[kr f(x_t)k^2] \leq \frac{1}{4} \frac{1}{2(1+q^2)^3(1-q^2)^2 K^2 G^2 + K(3L+2C_2L+2\frac{p}{2(1-q^2)}G)} \left[\frac{f_0 - f}{KT} + \frac{1}{T} + \dots \right]; \quad (B.22)$$

where $\dots = \frac{C_1 G^2 d}{p} + \frac{2C_1^2}{T} \frac{K L G^2 d}{p}$; $\dots = \frac{5}{2} \frac{K^2 L^2}{p} (\frac{1}{2} + 6K \frac{2}{g}) + [3(2C_2)L + 2\frac{p}{2(1-q^2)}G] \frac{1}{2m} \frac{1}{l^2}$, $C_1 = \frac{1}{1-q^2} + \frac{2q}{1-q^2}$ and $C_2 = \frac{1}{(1-q^2)^2} + \frac{4q^2}{(1-q^2)^2}$.

B.2. Proof of Corollary 4.19

Let $\frac{1}{l} = (\frac{p}{TK})$ and $\frac{1}{m} = (\frac{p}{Km})$, the convergence rate under full participation scheme is $\frac{1}{TKm}$.

B.3. Analysis on the Partial Participation Setting for FedCAMS

Let us present the theoretical analysis of the partial participation scheme of FedCAMS (Algorithm 2). Similar to partial participation scheme in Section 4.1, we have the following convergence analysis.

Theorem B.1. Under Assumption 4.1-4.3 and 4.14, if the local learning rate is as the following condition: $\min \frac{1}{8KL}; \frac{n(m-1)}{48m(n-1)} [K \frac{p}{4} \frac{1}{2(1+q^2)^3(1-q^2)^2 K^2 G^2 + (L + \frac{p}{2(1-q^2)}G)]^{-1}$, then the iterates of Algorithm 2 under partial participation scheme satisfy

$$\min_{t \in [T]} Ekr f(x_t)k^2 \leq \frac{1}{8} \frac{1}{4} \frac{1}{2(1+q^2)^3(1-q^2)^2 K^2 G^2 + K(3L+2C_2L+2\frac{p}{2(1-q^2)}G)} \left[\frac{f_0 - f}{KT} + \frac{1}{T} + \dots \right]; \quad (B.23)$$

where $\dots = \frac{C_1 G^2 d}{p} + \frac{2C_1^2}{T} \frac{K L G^2 d}{p}$, $\dots = \frac{C_1}{n} \frac{K L G^2}{p} + \frac{5}{2} \frac{K^2 L^2}{p} (\frac{1}{2} + 6K \frac{2}{g}) + [L + \frac{p}{2(1-q^2)}G] \frac{1}{n} \frac{1}{l^2} + [L + \frac{p}{2(1-q^2)}G] \frac{1}{n(m-1)} [15K^2 L^2 (\frac{1}{2} + 6K \frac{2}{g}) + 3K \frac{2}{g}]$ and $C_1 = \frac{1}{1-q^2} + \frac{2q}{1-q^2}$.

Remark B.2. The upper bound for $\min_{t \in [T]} Ekr f(x_t)k^2$ of partial participation is similar to full participation case but with a larger variance term. This is due to the fact that random sampling of participating workers introduces an additional variance during sampling.

Remark B.3. From Theorem B.1, constant C_1 is related to the compressor constant. The stronger compression we apply to the model difference $\frac{1}{l}$ corresponding to larger q ($q \rightarrow 1$) leads to worse convergence due to larger information losses.

Next, we provide theoretical proofs for the partial participation analysis for FedCAMS.

Proof of Theorem B.1:

Notations and equations. From the update rule of Algorithm 2, we have $e_t = 0$, $e_t = \frac{1}{m} \sum_{i=1}^m e_t^i$ and $m_t = (1 - \frac{1}{m}) \sum_{i=1}^m b_t^i$. Denote a global uncompressed difference $e_t = \frac{1}{p} \sum_{i \in S_t} e_t^i$. Denote a virtual momentum sequence: $m_t^0 = \frac{1}{p} \sum_{i=1}^m m_t^0 + (1 - \frac{1}{m}) \sum_{i=1}^m m_t^0$, hence we have $m_t^0 = (1 - \frac{1}{m}) \sum_{i=1}^m m_t^0$. Define additional two virtual sequences $\frac{1}{l} = \frac{1}{n} \sum_{i=1}^m \frac{1}{l}$ and $b_t^0 = \frac{1}{n} \sum_{i=1}^m b_t^i$. Note that when the client does not take part in the round of participation at step t , we have $\frac{1}{l} = b_t^i = 0$, therefore, $\frac{1}{l} = \frac{1}{l}$ and $b_t^0 = b_t^i$.

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

$$b_t = \frac{1}{|S_t|} \sum_{i \in S_t} (b_t^i - \bar{b}_t) = \frac{1}{n} \sum_{i=1}^n (b_t^i - \bar{b}_t) = \frac{1}{n} \sum_{i=1}^n (e_t^i - e_{t+1}^i) = e_t^0 - e_{t+1}^0; \quad (\text{B.24})$$

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

$$z_{t+1} := (1 - \alpha) z_t + \alpha \sum_{i=1}^n w_i^{t+1} + e_{t+1}^0;$$

and keep using the Lyapunov function from (B.4). For the expectation of model difference, we have

$$E_{S_t} [z_t] = \frac{1}{n} E_{S_t} \sum_{i=1}^n w_i^t = E_{S_t} [w_t] = \frac{1}{m} \sum_{i=1}^m w_i^t = \bar{w}_t; \quad (\text{B.25})$$

The proof of FedCAMS in partial participation settings has a similar outline combining the proof of partial participation in FedAMS and full participation in FedCAMS. By Assumption 4.1, we have

$$\begin{aligned} E[f(z_{t+1})] - f(z_t) &= E \left[r f(x_t); \underbrace{\psi_t^{1=2}}_{T_1^0} \right] \\ &= E \left[r f(z_t); \underbrace{\frac{1}{1} \psi_t^{1=2} \psi_t^{1=2} m_t^0 + \psi_t^{1=2} \psi_t^{1=2}}_{T_2^0} \right] \\ &+ \frac{2L}{2} E \left[\psi_t^{1=2} \psi_t^{1=2} \underbrace{\frac{1}{1} \psi_t^{1=2} \psi_t^{1=2} m_t^0 + \psi_t^{1=2} \psi_t^{1=2}}_{T_3^0} \right] \\ &+ E \left[r f(z_t) - r f(x_t); \underbrace{\psi_t^{1=2}}_{T_4^0} \right]; \end{aligned} \quad (\text{B.26})$$

Note that the bound for T_2^0 is exactly the same as the bound T_1 . For the three corresponding terms T_1^0 , T_3^0 and T_4^0 which include the second-order momentum estimate of w_t , similar to the full participation settings, we have

$$T_1^0 \leq \frac{\rho}{2} E \left[r f(x_t); \underbrace{\psi_t^{1=2}}_{T_1^0} \right] + \frac{\rho}{2} E \left[r f(x_t); \underbrace{\psi_t^{1=2} \psi_t^{1=2}}_{T_3^0} \right] + \frac{\rho}{2} E \left[r f(x_t); \underbrace{\psi_t^{1=2}}_{T_4^0} \right]; \quad (\text{B.27})$$

The first term in (B.27) does not change in partial participation scheme. The second term is changed due to the variance of w_t changes. For the second term T_3^0 , we have

$$\frac{\rho}{2} E \left[r f(x_t); \underbrace{\psi_t^{1=2} \psi_t^{1=2}}_{T_3^0} \right] \leq \frac{\rho}{2} \frac{2(1-\alpha)G}{\alpha} E[k_t^2]; \quad (\text{B.28})$$

For T_3^0 , similar to the proof of T_3 , we have

$$\sum_{t=1}^T \frac{2L}{2} E \left[k_t^2 \right] + \sum_{t=1}^T 2LC_1^2 \alpha^2 K^2 G^2 E \left[\psi_t^{1=2} \psi_t^{1=2} \right]; \quad (\text{B.29})$$

where $C_1 = \frac{1}{1-\alpha} + \frac{2q}{1-q^2}$. For T_4^0 in partial participation, we have

$$\begin{aligned} T_4^0 &= \mathbb{E} \left[f(z_t) - f(x_t); \psi_t^{1=2} \right] \\ &\leq \mathbb{E} \left[k f(z_t) - f(x_t) k \psi_t^{1=2} \right] \\ &\leq \frac{1}{1-\alpha} \mathbb{E} \left[\psi_t^{1=2} m_t^0 + \psi_t^{1=2} \right] \psi_t^{1=2} \\ &\leq \frac{1}{1-\alpha} \mathbb{E} \left[\psi_t^{1=2} \right] \psi_t^{1=2}. \end{aligned} \quad (\text{B.30})$$

Hence, the summation from T_1^0 to T_4^0 over total iteration T is:

$$\begin{aligned} \mathbb{E}[f(z_{T+1})] - f(z_1) &= \sum_{t=1}^T [T_1^0 + T_2^0 + T_3^0 + T_4^0] \\ &\leq \frac{1}{4} \frac{K}{m^2} \sum_{t=1}^T \mathbb{E} \left[\frac{r f(x_t)}{2V_{t+1}} \right]^2 + \frac{5}{1-\alpha} \frac{3K^2 L^2 T}{2} \left(\frac{1}{1-\alpha} + 6K \right) \frac{p}{2(1-\alpha)} \frac{G}{m} \sum_{t=1}^T \mathbb{E}[k_t k^2] \\ &\quad + \frac{1}{2Km^2} \sum_{t=1}^T \mathbb{E} \left[\frac{1}{2V_{t+1}} \sum_{i=1}^n \sum_{k=0}^{K-1} r F_i(x_t) \right]^2 + C_1 \frac{1}{1-\alpha} \frac{KG^2}{m} \sum_{t=1}^T \mathbb{E} \left[\psi_t^{1=2} \right] \psi_t^{1=2} \\ &\quad + C_1^2 \frac{1}{1-\alpha} \frac{2}{1-\alpha} \frac{K^2 LG^2}{m} \sum_{t=1}^T \mathbb{E} \left[\psi_t^{1=2} \right] \psi_t^{1=2} + C_1^2 \frac{1}{1-\alpha} \frac{2}{1-\alpha} \frac{K^2 LG^2}{m} \sum_{t=1}^T \mathbb{E} \left[\psi_t^{1=2} \right] \psi_t^{1=2} \\ &\quad + \frac{2L}{m} \sum_{t=1}^T \mathbb{E}[k_t k^2] + \frac{C_1 T}{1-\alpha} \frac{2}{1-\alpha} \frac{K^2 LG^2}{m} \\ &\quad + \frac{1}{4} \frac{K}{m^2} \frac{1}{2(1-\alpha)^2} \sum_{t=1}^T \mathbb{E}[k r f(x_t) k^2] + \frac{5}{1-\alpha} \frac{3K^2 L^2 T}{2} \left(\frac{1}{1-\alpha} + 6K \right) \frac{p}{2(1-\alpha)} \frac{G}{m} \\ &\quad + \frac{2C_1^2}{1-\alpha} \frac{2}{1-\alpha} \frac{K^2 LG^2}{m} \frac{1}{2} \frac{1}{4} \frac{1}{2(1-\alpha)^2} \frac{1}{K m^2} \sum_{t=1}^T \mathbb{E} \left[\sum_{i=1}^n \sum_{k=0}^{K-1} r F_i(x_t) \right]^2 \\ &\quad + \frac{2}{n} \frac{2}{1-\alpha} \frac{LKT}{m} + \frac{p}{2(1-\alpha)} \frac{2}{1-\alpha} \frac{KTG}{m} \frac{1}{1-\alpha} + \frac{C_1 T}{1-\alpha} \frac{2}{1-\alpha} \frac{K^2 LG^2}{m} \\ &\quad + \frac{2}{1-\alpha} \frac{2}{1-\alpha} \frac{L}{m} + \frac{p}{2(1-\alpha)} \frac{2}{1-\alpha} \frac{G}{m} \frac{m}{mn(m-1)} 15mK^3 L^3 \frac{1}{1-\alpha} \left(\frac{1}{1-\alpha} + 6K \right) T \\ &\quad + (90mK^4 L^2 \frac{1}{1-\alpha} + 3mK^2) \sum_{t=1}^T \mathbb{E}[k r f(x_t) k^2] + 3mK^2 T \frac{p}{2(1-\alpha)} \\ &\quad + \frac{2}{1-\alpha} \frac{2}{1-\alpha} \frac{L}{m} + \frac{p}{2(1-\alpha)} \frac{2}{1-\alpha} \frac{G}{m} \frac{n-1}{mn(m-1)} \sum_{t=1}^T \mathbb{E} \left[\sum_{i=1}^n \sum_{k=0}^{K-1} r F_i(x_t) \right]^2. \end{aligned} \quad (\text{B.31})$$

The proof outline is similar with previous proof. We take the use of Lemma C.2, C.6, C.8 for corresponding terms. By additional constraints of local learning rate with the inequality $\frac{2L}{m} +$

$\frac{p}{2(1-\alpha)} \frac{G}{m} \frac{1}{mn(m-1)} \frac{1}{2Km^2} \frac{1}{4} \frac{1}{2(1-\alpha)^2} \frac{1}{K^2 G^2} + \dots$, we obtain the constraint $\frac{1}{1-\alpha}$
 $\frac{n(m-1)}{m(n-1)} \frac{p}{2(1-\alpha)} \frac{G}{m} \frac{1}{mn(m-1)} \frac{1}{2Km^2} \frac{1}{4} \frac{1}{2(1-\alpha)^2} \frac{1}{K^2 G^2} + [L + \frac{p}{2(1-\alpha)} \frac{G}{m}]$, and we further need satisfies $\frac{1}{1-\alpha} \frac{K}{4} \frac{1}{4} \frac{1}{2(1-\alpha)^2} \frac{1}{K^2 G^2} +$
 $(\frac{2L}{m} + \frac{p}{2(1-\alpha)} \frac{G}{m}) \frac{1}{mn(m-1)} (90mK^4 L^2 \frac{1}{1-\alpha} + 3mK^2) \frac{p}{2(1-\alpha)} \frac{G}{m} \frac{1}{mn(m-1)} \frac{1}{2Km^2} \frac{1}{4} \frac{1}{2(1-\alpha)^2} \frac{1}{K^2 G^2} + \dots$. Hence for the convergence

rate, we have

$$\begin{aligned}
 & \frac{q}{8} \frac{K}{4} \frac{(1+q^2)^3}{2(1-q^2)^2} \frac{1}{T} \sum_{i=1}^T \mathbb{E}[\|k_r f(x_t)k^2] \\
 & \frac{f(z_0) - \mathbb{E}[f(z_T)]}{T} + \frac{5}{p} \frac{3K^2L^2}{2} \left(\frac{1}{n} + 6K \frac{2}{g} \right) + L + \frac{p}{2(1-q^2)} G \frac{K^2}{n} \frac{1}{n} \\
 & + \frac{C_1}{T} \frac{K^2G^2d}{p} + \frac{2C_1^2}{T} \frac{K^2LG^2d}{p} + \frac{C_1}{T} \frac{K^2LG^2}{p} \\
 & + \frac{1}{n} \frac{L}{p} + \frac{p}{2(1-q^2)} \frac{1}{n} G \frac{m}{mn(m-1)} \frac{n}{1} [15mK^3L^2 \left(\frac{1}{n} + 6K \frac{2}{g} \right) + 3mK^2 \frac{2}{g}]: \quad (\text{B.32})
 \end{aligned}$$

Therefore,

$$\min \mathbb{E}[\|k_r f(x_t)k^2] \leq \frac{q}{8} \frac{K}{4} \frac{(1+q^2)^3}{2(1-q^2)^2} \frac{1}{T} \sum_{i=1}^T \mathbb{E}[\|k_r f(x_t)k^2] + \frac{f_0 - f}{KT} + \frac{1}{T} + \dots; \quad (\text{B.33})$$

where $\frac{q}{8} \frac{K}{4} \frac{(1+q^2)^3}{2(1-q^2)^2} \frac{1}{T} \sum_{i=1}^T \mathbb{E}[\|k_r f(x_t)k^2]$; $\frac{f_0 - f}{KT} + \frac{1}{T} + \dots$; $\frac{C_1}{T} \frac{K^2G^2d}{p} + \frac{2C_1^2}{T} \frac{K^2LG^2d}{p} + \frac{C_1}{T} \frac{K^2LG^2}{p}$; $\frac{5}{p} \frac{3K^2L^2}{2} \left(\frac{1}{n} + 6K \frac{2}{g} \right) + [L + \frac{p}{2(1-q^2)} G] \frac{K^2}{n} \frac{1}{n} + [L + \frac{p}{2(1-q^2)} G] \frac{1}{n} \frac{(m-n)}{(m-1)} [15K^2L^2 \left(\frac{1}{n} + 6K \frac{2}{g} \right) + 3K^2 \frac{2}{g}]$ and $C_1 = \frac{1}{1-q} + \frac{2q}{1-q^2}$.

C. Supporting Lemmas

Lemma C.1. For the element-wise difference $\mathbf{w}_t = \frac{1}{v_t} - \frac{1}{2v_{t-1}}$, we have $\|k\mathbf{w}_t\| \leq \frac{p}{1-2} k_t k$.

Proof. Note that we have:

$$\begin{aligned}
 \|k\mathbf{w}_t\| &= \left| \frac{1}{v_t} - \frac{1}{2v_{t-1}} \right| \\
 &= \left| \frac{p}{2v_{t-1}} - \frac{p}{v_t} \right| = \frac{p}{2v_{t-1}} \left(\frac{p}{2v_{t-1}} + \frac{p}{v_t} \right) \\
 &= \frac{p}{v_t} \frac{p}{2v_{t-1}} \frac{1}{\left(\frac{p}{2v_{t-1}} + \frac{p}{v_t} \right)} \\
 &= \frac{p}{v_t} \frac{p}{2v_{t-1}} \frac{(1-2)^{\frac{t}{2}}}{\left(\frac{p}{2v_{t-1}} + \frac{p}{v_t} \right)} \\
 &= \frac{p}{v_t} \frac{p}{2v_{t-1}} \frac{(1-2)^{\frac{t}{2}}}{1 - \frac{2}{t}} \\
 &= \frac{p}{1-2} k_t k; \quad (\text{C.1})
 \end{aligned}$$

where the forth equation holds by the update rule of $v_t = 2v_{t-1} + (1-2)^{\frac{t}{2}}$, and the first inequality holds due to $\frac{p}{v_t} - \frac{p}{v_t} \frac{1}{1-2} \frac{2}{t}$ and $\frac{p}{2v_{t-1}} > 0$. This concludes the proof. \square

Lemma C.2. For the variance difference sequence $\mathbf{v}_t^{1=2} - \mathbf{v}_t^{1=2}$, we have

$$\sum_{t=1}^T \mathbf{v}_t^{1=2} - \mathbf{v}_t^{1=2} \leq \frac{d}{p}; \quad (\text{C.2})$$

Proof. By the definition of variance matrix \mathbf{v}_t , and the non-decreasing update of FedCAMS, we have $\mathbf{v}_t =$

$\max(\mathbf{b}_{t-1}; \mathbf{v}_t)$, we have

$$\begin{aligned}
 \sum_{t=1}^T \mathbf{v}_t^{1:2} - \mathbf{v}_t^{1:2} &= \sum_{t=1}^T \left(\frac{1}{\mathbf{b}_{t-1}} - \frac{1}{\mathbf{b}_t} \right) \\
 &= \sum_{t=1}^T \left(\frac{1}{\mathbf{b}_{t-1}} - \frac{1}{\mathbf{b}_t} \right) \\
 &= \frac{1}{\mathbf{b}_0} - \frac{1}{\mathbf{b}_T} \\
 &\leq d;
 \end{aligned} \tag{C.3}$$

where the inequality holds by the definition of $\mathbf{b} \in \mathbb{R}^d$. For the sum of the variance difference under the norm, we have

$$\begin{aligned}
 \sum_{t=1}^T \mathbf{v}_t^{1:2} - \mathbf{v}_t^{1:2} &= \sum_{t=1}^T \left(\frac{1}{\mathbf{b}_{t-1}} - \frac{1}{\mathbf{b}_t} \right)^2 \\
 &= \sum_{t=1}^T \left(\frac{1}{\mathbf{b}_{t-1}} - \frac{1}{\mathbf{b}_t} \right)^2 \\
 &= \sum_{t=1}^T \left(\frac{1}{\mathbf{b}_{t-1}} - \frac{1}{\mathbf{b}_t} \right) \\
 &= \frac{1}{\mathbf{b}_0} - \frac{1}{\mathbf{b}_T} \\
 &\leq d;
 \end{aligned} \tag{C.4}$$

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

Lemma C.3. The compression error has the following absolute bound

$$\mathbf{ke}_t^i k^2 \leq \frac{4q^2}{(1 - q^2)^2} \mathbf{K}^2 \mathbf{G}^2; \quad \mathbf{ke}_t k^2 \leq \frac{4q^2}{(1 - q^2)^2} \mathbf{K}^2 \mathbf{G}^2; \tag{C.5}$$

Proof. For all $t \in [T]$, by Assumption 4.14 and Young's inequality, we have

$$\begin{aligned}
 \mathbf{ke}_{t+1}^i k^2 &= \mathbf{k} \left(\frac{1}{\mathbf{b}_t} + \mathbf{e}_t^i \right) \mathbf{C} \left(\frac{1}{\mathbf{b}_t} + \mathbf{e}_t^i \right) k^2 \\
 &\leq q^2 \mathbf{k} \left(\frac{1}{\mathbf{b}_t} + \mathbf{e}_t^i \right) k^2 \\
 &\leq q^2 (1 + \epsilon) \mathbf{ke}_t^i k^2 + q^2 \left(1 + \frac{1}{\epsilon} \right) \mathbf{k} \frac{1}{\mathbf{b}_t} k^2 \\
 &\leq \frac{1 + q^2}{2} \mathbf{ke}_t^i k^2 + \frac{2q^2}{1 - q^2} \mathbf{k} \frac{1}{\mathbf{b}_t} k^2;
 \end{aligned}$$

where the last inequality holds by choosing $\epsilon = \frac{1 - q^2}{2q^2}$. Thus we obtain the absolute bound for the error terms,

$$\begin{aligned}
 \mathbf{ke}_t^i k^2 &\leq \frac{4q^2}{(1 - q^2)^2} \mathbf{K}^2 \mathbf{G}^2; \\
 \mathbf{ke}_t k^2 &= \frac{1}{m} \sum_{i=1}^m \mathbf{e}_t^i k^2 \leq \frac{1}{m} \sum_{i=1}^m \mathbf{ke}_t^i k^2 \leq \frac{4q^2}{(1 - q^2)^2} \mathbf{K}^2 \mathbf{G}^2;
 \end{aligned} \tag{C.6}$$

It concludes the proof. \square

Lemma C.4. Under Assumptions 4.2 and 4.14, for FedAMS, we have $\|f(x)\| \leq G$, $\|k_t\| \leq \eta KG$, $\|m_t\| \leq \eta KG$ and $\|v_t\| \leq \frac{1}{2} K^2 G^2$. For FedCAMS, we have $\|f(x)\| \leq G$, $\|b_t\| \leq \frac{4(1+q^2)^3}{(1-q^2)^2} \frac{1}{2} K^2 G^2$, $\|m_t^0\| \leq \eta KG$ and $\|v_t\| \leq \frac{4(1+q^2)^3}{(1-q^2)^2} \frac{1}{2} K^2 G^2$, where $m_t^0 = \eta m_{t-1}^0 + (1-\eta) v_t$.

Proof. Since f has G -bounded stochastic gradients, for any x and η , we have $\|f(x; \eta)\| \leq G$, we have

$$\|k_t f(x)\| = \eta E \|r f(x; \eta)\| \leq E \|k_t f(x; \eta)\| \leq G.$$

For FedAMS, the model difference x_t , by definition, has the following formula,

$$x_t = x_{t,K}^i \quad x_t = \frac{1}{K} \sum_{k=1}^K g_{t,k}^i;$$

therefore,

$$\|k_t\| \leq \eta K \|g_{t,k}^i\| \leq \eta KG;$$

Thus the bound for momentum m_t and variance v_t has the formula of

$$\begin{aligned} \|m_t\| &= (1-\eta) \sum_{i=1}^t \|k_i\| \leq \eta KG; \\ \|v_t\| &= (1-\eta^2) \sum_{i=1}^t \|k_i\|^2 \leq \frac{1}{2} K^2 G^2. \end{aligned}$$

For the compressed version, FedCAMS, we have

$$\begin{aligned} \|b_t\|^2 &\leq K C (\|v_t + e_t\|)^2 \\ &\leq K C (\|v_t + e_t\| + (\|v_t + e_t\|)^2) \\ &\leq 2(q^2 + 1) \|v_t + e_t\|^2 \\ &\leq 4(q^2 + 1) [\|v_t\|^2 + \|e_t\|^2] \\ &\leq \frac{4(1+q^2)^3}{(1-q^2)^2} \frac{1}{2} K^2 G^2; \end{aligned}$$

where the third inequality holds due to Assumption 4.14, and the last inequality holds due to Lemma C.3. The virtual momentum sequence m_t^0 has the same bound as of FedAMS. For the variance sequence of FedCAMS, we have

$$\|v_t\| = (1-\eta^2) \sum_{i=1}^t \|b_i\|^2 \leq \frac{4(1+q^2)^3}{(1-q^2)^2} \frac{1}{2} K^2 G^2.$$

This concludes the proof. □

Lemma C.5. The global model difference $x_t = \frac{1}{m} \sum_{i=1}^m x_t^i$ in full participation cases satisfy

$$E[\|x_t\|^2] \leq \frac{K}{m} \frac{1}{m} + \frac{1}{m^2} E \sum_{i=1}^m \sum_{k=0}^{X-1} \|r F_i(x_{t,k}^i)\|^2;$$

Proof. For $E[k_t^2]$ in full participation case, we have

$$\begin{aligned}
 E[k_t^2] &= E \left[\frac{1}{m} \sum_{i=1}^n \sum_{k=0}^{K-1} \|g_{t,k}^i\|^2 \right] \\
 &= \frac{1}{m^2} E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} \|g_{t,k}^i\|^2 \right] \\
 &= \frac{1}{m^2} E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} (\|g_{t,k}^i - r F_i(x_{t,k}^i)\|^2) \right] + \frac{1}{m^2} E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} \|r F_i(x_{t,k}^i)\|^2 \right] \\
 &= \frac{K}{m} \sum_{i=1}^n \sigma_i^2 + \frac{1}{m^2} E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} \|r F_i(x_{t,k}^i)\|^2 \right]; \tag{C.7}
 \end{aligned}$$

where the inequality holds by Assumption 4.2. This concludes the proof. \square

Lemma C.6. The global model difference $\sum_{i \in S_t} \mathbf{x}_t^i$ in partial participation cases satisfy

$$\begin{aligned}
 E[k_t^2] &= \frac{K}{n} \sum_{i=1}^n \sigma_i^2 + \frac{\sum_{i=1}^n (m-n)}{mn(m-1)} [15mK^3L^3 \sum_{i=1}^n (\sigma_i^2 + 6K \frac{\sigma_i^2}{g}) + 90mK^4L^2 \sum_{i=1}^n \sigma_i^2 + 3mK^2kr f(x_t)k^2 \\
 &\quad + 3mK^2 \frac{\sigma_i^2}{g}] + \frac{\sum_{i=1}^n (n-1)}{mn(m-1)} E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} \|r F_i(x_{t,k}^i)\|^2 \right];
 \end{aligned}$$

Proof. We have

$$\begin{aligned}
 E[k_t^2] &= E \left[\frac{1}{n} \sum_{i \in S_t} \|\mathbf{x}_t^i\|^2 \right] \\
 &= \frac{1}{n^2} E \left[\sum_{i=1}^n \|\mathbf{x}_t^i\|^2 \right] \\
 &= \frac{1}{n^2} E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} \|g_{t,k}^i - r F_i(x_{t,k}^i)\|^2 + \sum_{i=1}^n \sum_{k=0}^{K-1} \|r F_i(x_{t,k}^i)\|^2 \right] \\
 &= \frac{1}{n^2} E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} \|g_{t,k}^i - r F_i(x_{t,k}^i)\|^2 + \sum_{i=1}^n \sum_{k=0}^{K-1} \|r F_i(x_{t,k}^i)\|^2 \right] \\
 &= \frac{1}{mn} E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} \|g_{t,k}^i - r F_i(x_{t,k}^i)\|^2 \right] + \frac{1}{n^2} E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} \|r F_i(x_{t,k}^i)\|^2 \right] \\
 &= \frac{K}{n} \sum_{i=1}^n \sigma_i^2 + \frac{1}{n^2} E \left[\sum_{i=1}^n \sum_{k=0}^{K-1} \|r F_i(x_{t,k}^i)\|^2 \right]; \tag{C.8}
 \end{aligned}$$

where the first equation holds due to $\sum_{i \in S_t} \mathbf{x}_t^i = \frac{n}{m} \sum_{i=1}^n \mathbf{x}_t^i$. Note that we have

$$\begin{aligned}
 \sum_{i=1}^n \sum_{k=0}^{K-1} \|r F_i(x_{t,k}^i)\|^2 &= \sum_{i=1}^n \sum_{k=0}^{K-1} \|r F_i(x_{t,k}^i)\|^2 + \sum_{i \neq j} \sum_{k=0}^{K-1} \|r F_i(x_{t,k}^i) - r F_j(x_{t,k}^j)\|^2 \\
 &= \sum_{i=1}^n \sum_{k=0}^{K-1} \|r F_i(x_{t,k}^i)\|^2 + \frac{1}{2} \sum_{i \neq j} \sum_{k=0}^{K-1} \|r F_i(x_{t,k}^i) - r F_j(x_{t,k}^j)\|^2; \tag{C.9}
 \end{aligned}$$

where the second equation holds due to $\sum_{i=1}^m \mathbf{x}_i k^2 = \sum_{i=1}^m m k x_i k^2 = \frac{1}{2} \sum_{i \neq j} k x_i - x_j k^2$. By the sampling strategy

Proof. By the updating rule, we have

$$\begin{aligned}
 E[k_t k^2] &= E \left[k \left(1 - \frac{1}{m}\right)^{X_t} \sum_{i=1}^t k^2 \right] \\
 &\leq \left(1 - \frac{1}{m}\right)^2 E \sum_{i=1}^{X_t} E \sum_{j=1}^2 k^2 \\
 &\leq \left(1 - \frac{1}{m}\right)^2 E \sum_{i=1}^{X_t} E \sum_{j=1}^2 k^2 \\
 &\leq \left(1 - \frac{1}{m}\right)^2 E \sum_{i=1}^{X_t} E[k^2] \\
 &\leq \frac{K^2}{m} \sum_{i=1}^2 + \frac{1}{m^2} \left(1 - \frac{1}{m}\right)^{X_t} \sum_{i=1}^{X_t} E \sum_{k=0}^{X-1} F_i(x_{t;k}^i)^2; \tag{C.13}
 \end{aligned}$$

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 C.5. Hence summing over $t = 1; \dots; T$, we have

$$\sum_{t=1}^T E[k_t k^2] \leq \frac{TK^2}{m} \sum_{i=1}^2 + \frac{1}{m^2} \sum_{t=1}^T E \sum_{i=1}^{X_t} \sum_{k=0}^{X-1} F_i(x_{t;k}^i)^2; \tag{C.14}$$

For the compression error, following the proof from Lemma C.3, for t and each local worker i , we have the induction

$$\begin{aligned}
 E[k_{t+1}^i k^2] &\leq \frac{2q^2}{1 - q^2} \sum_{i=1}^{X_t} \frac{1 + q^2}{2} E[k^i k^2] \\
 &\leq \frac{4q^2}{(1 - q^2)^2} \frac{K^2}{m} \sum_{i=1}^2 + \frac{1}{m^2} \sum_{i=1}^{X_t} \frac{1 + q^2}{2} E \sum_{k=0}^{X-1} F_i(x_{t;k}^i)^2; \tag{C.15}
 \end{aligned}$$

For the sequence e_t , similar as the previous analysis, we have

$$\begin{aligned}
 E[k_{t+1} k^2] &= E \left[\left(1 - \frac{1}{m}\right)^{X_t} \sum_{i=1}^t e^2 \right] \\
 &\leq \left(1 - \frac{1}{m}\right)^{X_t} \sum_{i=1}^t E[k e^2] \\
 &\leq \left(1 - \frac{1}{m}\right)^{X_t} \sum_{i=1}^t \frac{1}{m} \sum_{i=1}^{X^n} E[k e^2] \\
 &\leq \frac{4q^2}{(1 - q^2)^2} \frac{K^2}{m} \sum_{i=1}^2 + \frac{2q^2}{m^2} \left(1 - \frac{1}{m}\right)^{X_t} \sum_{i=1}^{X_t} \sum_{j=1}^X \frac{1 + q^2}{2} E \sum_{k=0}^{X-1} F_i(x_{j;k}^i)^2; \tag{C.16}
 \end{aligned}$$

Summing over $t = 1; \dots; T$, we have

$$\begin{aligned}
 \sum_{t=1}^T E[k_{t+1} k^2] &\leq \frac{4Tq^2}{(1 - q^2)^2} \frac{K^2}{m} \sum_{i=1}^2 + \frac{2q^2}{m^2} \sum_{t=1}^T \sum_{i=1}^{X_t} \frac{1 + q^2}{2} E \sum_{k=0}^{X-1} F_i(x_{j;k}^i)^2 \\
 &\leq \frac{4Tq^2}{(1 - q^2)^2} \frac{K^2}{m} \sum_{i=1}^2 + \frac{4q^2}{m^2} \sum_{t=1}^T E \sum_{i=1}^{X_t} \sum_{k=0}^{X-1} F_i(x_{t;k}^i)^2; \tag{C.17}
 \end{aligned}$$

□

Lemma C.8. Under Assumptions 4.1-4.3 and Assumption 4.14, for the momentum sequence $\mathbf{m}_t = (1 - \alpha) \mathbf{m}_{t-1} + \alpha \Delta_t$ in partial participation settings, we have

$$\sum_{t=1}^T \mathbb{E}[\|\mathbf{m}_t\|^2] \leq \frac{KT}{n} \sum_{i=1}^2 \frac{1}{i} + \frac{2}{n^2} \sum_{t=1}^T \mathbb{E} \left[\sum_{i \in S_t} \sum_{k=0}^{K-1} \|\nabla F_i(\mathbf{x}_{t,k}^i)\|^2 \right].$$

Proof. The proof outline is the same as the proof of Lemma C.7, the main difference is $E[\|\Delta_t\|^2]$ has changed, so we need to apply Lemma C.6 instead of Lemma C.5 during the proof. \square

Lemma C.9. (This lemma directly follows from Lemma 3 in FedAdam (Reddi et al., 2020). For local learning rate which satisfying $\eta \leq \frac{1}{8KL}$, the local model difference after k ($\forall k \in \{0; 1; \dots; K - 1\}$) steps local updates satisfies

$$\frac{1}{m} \sum_{i=1}^m \mathbb{E}[\|\mathbf{x}_{t,k}^i - \mathbf{x}_t\|^2] \leq 5K \sum_{i=1}^2 \left(\frac{1}{i} + 6K \frac{1}{i} \right) + 30K^2 \sum_{i=1}^2 \mathbb{E}[\|\nabla f(\mathbf{x}_t)\|^2]. \quad (\text{C.18})$$

Proof. The proof of Lemma C.9 is exactly same as the proof of Lemma 3 in Reddi et al. (2020). \square

D. Additional Discussions

The additional server-to-worker communication: Our current analysis only focus on one-way compression from worker to server while the server-to-worker broadcasting is still uncompressed since of cost of broadcasting is in general cheaper than worker to server uploading. Note that it is also straightforward to compress \mathbf{x} for server-to-worker communication in the full participation scheme (with guarantees). So we can indeed achieve high communication efficiency even for two-way compression. Table 1 shows the communication bits comparison for scaled sign and top- k compressors, where T is the total iteration of training and d denotes the dimension of \mathbf{x} . Specifically, Table 2 shows the communication bits corresponding to the experiments showing by Figure 4. However, for the partial participating setting, it will encounter a synchronization issue, which is highly non-trivial to solve. Thus we leave the two-way compression strategy for future work.

Method	Uncompressed	One-way Compression	Two-way Compression
Scaled sign	$32d \times 2T$	$(32 + d) \times T + 32d \times T$	$(32 + d) \times 2T$
Top- k	$32d \times 2T$	$\approx 32(2k + d) \times T$	$\approx 32 \times 2k \times 2T$

Table 1. Communication bits comparisons for scaled sign and top- k compressors.

Method	Uncompressed	One-way Compression	Two-way Compression
Scaled sign	$3:58 \times 10^{11}$	$1:84 \times 10^{11}$	$1:12 \times 10^{10}$
Top- k with $r = 1=64$	$3:58 \times 10^{11}$	$1:84 \times 10^{11}$	$1:12 \times 10^{10}$
Top- k with $r = 1=128$	$3:58 \times 10^{11}$	$1:82 \times 10^{11}$	$5:59 \times 10^9$
Top- k with $r = 1=256$	$3:58 \times 10^{11}$	$1:80 \times 10^{11}$	$2:79 \times 10^9$

Table 2. Approximate communication bits comparisons for scaled sign and top- k with $r = 1/64$, $r = 1/128$ and $r = 1/256$ compressors when training CIFAR-10 on ResNet-18 model for 500 rounds.

E. Additional Experimental Results

E.1. Hyperparameter Settings

We conduct detailed hyperparameter searches to find the best hyperparameters for each baseline methods including ours. In details, we grid search over the local learning rate $\eta \in \{0.0001; 0.001; 0.01; 0.1; 1.0\}$, the global learning rate $\alpha \in \{0.001; 0.01; 0.1; 1.0\}$ for all methods. For adaptive federated optimization methods, we set $\alpha_1 = 0.9$, $\alpha_2 = 0.99$. For FedAdam, FedYogi, and FedAMSGrad, we search the best η from $\{10^{-8}; 10^{-4}; 10^{-3}; 10^{-2}; 10^{-1}; 10^0\}$. For FedAMS and FedCAMS, we search the max stabilization η from $\{10^{-8}; 10^{-4}; 10^{-3}; 10^{-2}; 10^{-1}; 10^0\}$.

Specifically, for our ResNet-18 experiments, we set the local learning rate $\eta_l = 0.01$ and the global learning rate $\eta_g = 1.0$ for FedAvg, set $\eta_l = 0.01$, $\beta = 0.1$ and $\epsilon = 0.1$ for FedAdam and FedAMSGrad, set $\eta_l = 0.01$, $\beta = 1.0$ and $\epsilon = 0.1$ for FedYogi, set $\eta_l = 0.01$, $\beta = 1.0$ and max stabilization $\epsilon = 0.001$ for FedAMS and FedCAMS. For our ConvMixer-256-8 experiments, we set the local learning rate $\eta_l = 0.01$ and the global learning rate $\eta_g = 1.0$ for FedAvg, set $\eta_l = 0.01$, $\beta = 1.0$ and $\epsilon = 0.1$ for FedAdam, FedYogi and FedAMSGrad, set $\eta_l = 0.01$, $\beta = 1.0$ and max stabilization $\epsilon = 0.001$ for FedAMS and FedCAMS.

E.2. Additional Experiments

Figure 6 shows the effect of parameter n on the convergence rate of FedCAMS with choosing groups of parameters: $n \in \{5; 10; 20\}$. For both ResNet-18 and ConvMixer-256-8 models, it is shown that a larger number of participating clients n achieves a faster convergence rate, this backs up our theory.

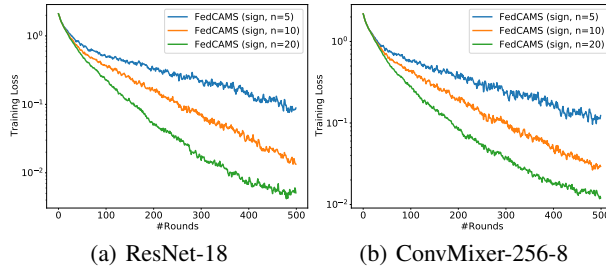


Figure 6. The learning curves for FedCAMS with different participating number of clients n in training CIFAR-10 data on ResNet-18 and ConvMixer-256-8 models.

Figure 7 shows the convergence result of FedAMS and other federated learning baselines on training CIFAR-100 dataset with the ResNet-18 model and the ConvMixer-256-8 model. We compare the training loss and test accuracy against the global rounds for each model. For the ResNet-18 model, FedAMS and FedYogi achieve significantly better performance comparing with other three baselines. In particular, FedYogi has a fast convergence rate at the beginning status, while FedAMS performs the best in terms of the final training loss and test accuracy. FedAvg achieves a slightly better training loss to FedAdam and FedAMSGrad but much higher test accuracy which is close to FedYogi and FedAMS.

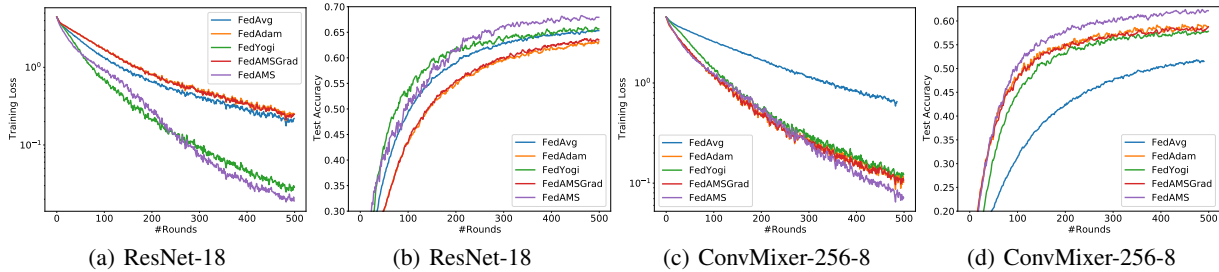


Figure 7. The learning curves for FedAMS and other federated learning baselines on training CIFAR-100 data (a)(b) show the results for ResNet-18 model and (c)(d) show the results for ConvMixer-256-8 model.

For the ConvMixer-256-8 model, which is typically trained via adaptive gradient method, we observe that all adaptive federated optimization methods (FedAdam, FedYogi, FedAMSGrad and FedAMS) achieve much better performance in terms of both training loss and test accuracy compared with FedAvg. In details, FedAMS again achieves a significantly better result than other baselines in terms of training loss and test accuracy. Other adaptive methods, including FedAdam, FedYogi, and FedAMSGrad, have similar convergence behaviour when training the ConvMixer-256-8 model. Such results empirically show the effectiveness of our proposed FedAMS method.

E.3. Additional Ablation Study

The ablation on ϵ : We conduct an ablation study with $\epsilon \in \{10^{-1}; 10^{-2}; 10^{-3}; 10^{-4}; 10^{-6}; 10^{-8}\}$ on CIFAR-10 in Table 3 and our ϵ value in experiments is chosen by its relatively higher test accuracy.

	10^{-1}	10^{-2}	10^{-3}	10^{-4}	10^{-6}	10^{-8}
Test acc (%)	90.45	90.51	90.94	90.72	90.49	90.30

Table 3. Ablation study on ϵ .