

# ON THE CONVERGENCE OF LoRA-BASED FEDERATED LEARNING: A UNIFIED ANALYSIS OF AGGREGATION-BROADCAST OPERATORS

Anonymous authors  
Paper under double-blind review

## ABSTRACT

Federated Learning (FL) enables collaborative model training across decentralized data sources while preserving data privacy. However, the increasing scale of Machine Learning (ML) models poses significant communication and computation challenges in FL. Low-Rank Adaptation (LoRA) has recently been integrated into FL as a Parameter-Efficient Fine-Tuning (PEFT) strategy, substantially lowering communication costs by transmitting only a small set of trainable parameters. Nevertheless, how to aggregate LoRA-updated local models on the server remains a critical and understudied problem. This paper presents a comprehensive theoretical analysis of LoRA-based FL frameworks. We first classify existing aggregation schemes into two main categories: Sum-Product (SP) and Product-Sum (PS). We then introduce the Aggregation-Broadcast Operator (ABO) as a general class encompassing all aggregation-broadcast methods. Any method in this class ensures local or global convergence as long as the corresponding Weak or Strong Convergence Condition is satisfied. In particular, we prove that the SP and PS aggregation methods satisfy the weak and strong convergence conditions, respectively, but differ in their ability to achieve the optimal convergence rate. Moreover, we conducted extensive experiments on standard open datasets to verify our theoretical findings.

**AI Acknowledgment:** We acknowledge that AI tools were employed to assist in paper writing and polishing the text to improve readability.

## 1 INTRODUCTION

Federated Learning (FL) has emerged as a promising framework for training machine learning models across decentralized data sources while preserving data privacy McMahan et al. (2017); Kairouz et al. (2021). However, the growing complexity of modern deep neural networks poses significant challenges for communication efficiency and resource-constrained devices, which are central challenges in FL Deng et al. (2020). Low-Rank Adaptation (LoRA), a parameter-efficient fine-tuning technique initially proposed for large-scale language models Hu et al. (2022a), has recently attracted increasing interest in FL due to its ability to reduce communication overhead by updating only a small subset of trainable parameters. The core idea of LoRA is to constrain the weight update on the model by a low-rank decomposition:

$$W' = W_0 + \Delta W, \quad \Delta W = BA, \quad (1)$$

Where  $B \in \mathbb{R}^{d \times r}$  and  $A \in \mathbb{R}^{r \times n}$ , with  $r \ll d$ . By training only  $B$  and  $A$ , LoRA significantly reduces the number of trainable parameters while maintaining model performance. Using LoRA in an FL setting is an effective and resource-efficient strategy. The global objective of LoRA in FL can be expressed as:

$$\arg \min_{A, B} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{(x, y) \sim P_{XY}^{(i)}} [\mathcal{L}(W_0 + AB; (x, y))]. \quad (2)$$

where  $m$  denotes the number of clients, and  $(x, y) \sim P_{XY}^{(i)}$  indicates that the local data of client  $i$  follows the distribution  $P_{XY}^{(i)}$ . By leveraging LoRA adapters, clients can fine-tune large foundation

models with minimal computational overhead. Because only the low-rank adapter matrices need to be communicated with the central server, this approach greatly reduces communication overhead. Compared to full-parameter fine-tuning, LoRA offers a more scalable and efficient solution for improving model performance in collaborative learning environments.

Despite its advantages, integrating LoRA into FL introduces new challenges, particularly in how the locally updated low-rank parameters are aggregated on the server side Yang et al. (2025). Unlike traditional FL, where full model weights or gradients are averaged directly, LoRA-based training requires the design of specialized aggregation strategies that respect the low-rank structure. In recent studies, multiple LoRA aggregation methods have been proposed, which we broadly classify into the following categories in this paper:

**Sum-Product-Type (SP) Aggregation Method.** This method is referred to as the SP method throughout the paper. Such a method is also referred to as the ideal aggregation method, as it shares the same form as FedAvg McMahan et al. (2017). Its aggregation form is as follows:

$$\Delta W = \frac{1}{m} \sum_{i=1}^m B_i A_i \quad (3)$$

This aggregation form can unify several recent methods. For example, FlexLoRA Bai et al. (2024) was the first to aggregate local models to the server by Eq. (3), and then broadcast by Singular Value Decomposition(SVD). FedIT Zhang et al. (2024a) uploads locally fine-tuned LoRA parameters from each client, which are then aggregated on the server using FedAvg to update the global model. FLoRA Wang et al. (2024), a stacking-based LoRA aggregation method, further improves this process by reducing the impact of noise during aggregation.

**Product-Sum-Type (PS) Aggregation Method.** This method is referred to as the PS aggregation method throughout the paper. Its aggregation form is as follows:

$$\Delta W = \left( \frac{1}{m} \sum_{i=1}^m B_i \right) \left( \frac{1}{m} \sum_{i=1}^m A_i \right) \quad (4)$$

This form encompasses several existing methods, such as Zero-Padding Cho et al. (2023) and RBLA Chen et al. (2024a) for Heterogeneous LoRA aggregation. FFA-LoRA Sun et al. (2024), which freezes the LoRA matrix  $A_i = A_0$  and only aggregates the LoRA Matrix  $B_i$ . Moreover, RoLoRA Chen et al. (2024b) employed an alternating form, aggregating only  $B_i$  in odd rounds and  $A_i$  in even rounds. Different from these, FedSA-LoRA Guo et al. (2025) updates and learn both  $B_i$  and  $A_i$ , but only the  $A_i$  matrices are shared for aggregation to learn general knowledge, and saves  $B_i$  locally for capturing client-specific knowledge.

**Other Aggregation Method.** Some aggregation methods cannot be easily categorized as either SP or PS type aggregation methods, such as FedInc Qin & Li (2024), which proposed a clustering-based aggregation method, enabling more fine-grained and adaptive aggregation. FedEx-LoRA Singhal et al. (2024) and LoRA-fair Bian et al. (2024) introduce correction terms during aggregation to make the results closer to the SP aggregation(ideal aggregation). CoLRN Nguyen et al. (2024) adopts a hybrid aggregation strategy by locally learning matrix A while globally sharing matrix B through server-side decomposition. LoRA-A<sup>2</sup> Koo et al. (2024) employs alternating minimization with adaptive rank selection to reduce communication costs by focusing on the most important LoRA ranks.

**Motivation** Although the two methods mentioned above are now widely used, their underlying mechanisms remain unclear, as we still face the following questions:

- The SP aggregation method is referred to as the ideal aggregation method Yang et al. (2025); Guo et al. (2025), but is it truly the fastest in terms of convergence speed?
- The PS aggregation method is widely adopted Cho et al. (2023); Chen et al. (2024a;b), but can it really guarantee convergence for the global model? What is the difference between SP and PS aggregation methods in terms of convergence speed?

- For more general aggregation algorithms, under what conditions can they guarantee the convergence of the global model?

Addressing these questions is of both theoretical and practical significance. From a theoretical standpoint, a unified convergence analysis helps us understand the fundamental principles that govern the success or failure of various LoRA aggregation strategies. It also provides a rigorous framework for comparing different methods on an equal footing. From a practical perspective, identifying the conditions that guarantee convergence can guide the design of more effective aggregation algorithms, enabling faster training and better performance in real-world FL scenarios.

**Contribution** Our research is primarily motivated by addressing the above questions, and on this basis, we have proposed some more general conclusions. We summarize our contributions as follows:

- We formally define the Aggregation-Broadcast Operator. Under mild assumptions, we establish both weak and strong convergence conditions. We prove that when weak convergence conditions are satisfied, the Aggregated Broadcasting Operator (ABO) ensures convergence of local models in the LoRA subspace at a rate of  $O(1/\sqrt{T})$ . When strong convergence conditions are met, it guarantees convergence of global models in the same subspace at the same rate.
- Especially, we prove that the SP Aggregation-Broadcast Operator satisfies the weak convergence condition but cannot achieve the optimal convergence rate due to broadcast errors, whereas the PS operator satisfies the strong convergence condition and achieves both global convergence and the optimal convergence rate.
- We perform comprehensive empirical studies to validate our theoretical findings. In particular, we investigate the effects of LoRA rank and the number of local training epochs on the convergence behavior of PS and SP aggregation methods, demonstrating strong consistency with our analytical predictions.

## 2 RELATED WORK

**Federated Learning** McMahan et al. introduced FedAvg in McMahan et al. (2017) as a decentralized and privacy-aware model training approach. Since then, numerous works have addressed the challenges of non-IID data Zhao et al. (2018), communication efficiency Konečný et al. (2016), and personalization Smith et al. (2017). Several algorithms have been proposed to improve optimization in heterogeneous settings, including FedProx Li et al. (2020a), SCAFFOLD Karimireddy et al. (2020), and MOON Li et al. (2021). Recent efforts also explore fairness Li et al. (2019) and adaptive aggregation Wang et al. (2020) to balance performance across clients.

**LoRA** Low-Rank Adaptation (LoRA) has emerged as an efficient Parameter-Efficient-Fine-Tuning (PEFT) method for Large Language Models (LLMs) Hu et al. (2022b). Early work, such as Universal Language Model Fine-tuning (ULMFiT), also explored efficient adaptation methods to reduce overhead Howard & Ruder (2018). Based on this, QA-LoRA Xu et al. (2024) introduces quantization-aware adaptation by integrating LoRA with low-bit quantization to further reduce memory usage. Similarly, QLoRA Dettmers et al. (2024) extends this idea by using 4-bit quantized LoRA adapters, achieving competitive performance with significantly lower memory footprint. In vanilla LoRA, the model rank requires manual configuration. To solve this issue, AdaLoRA Zhang et al. (2023) and AutoLoRA Zhang et al. (2024c) automate the rank selection process, allowing the model to adaptively allocate parameters where most needed. In theoretical analysis, Sathika et al. Malladi et al. (2023) investigate LoRA fine-tuning through the lens of kernel theory and show that, in the lazy training regime, its behavior closely mirrors that of full fine-tuning. Moreover, Zhu et al. (2024) shows that tuning LoRA matrix B is more impactful than tuning LoRA matrix A. Zeng et al. Zeng & Lee (2024) provide a theoretical analysis of LoRA’s expressive power for both Fully Connected Neural Networks (FNNs) and Transformer Networks (TNs).

**LoRA-based Federated Learning** Based on the above advancements, numerous LoRA-based methods have been proposed for FL Wu et al. (2024); Cho et al. (2024); Yi et al. (2023); Bai et al.

(2024); Chen et al. (2024a); Zhang et al. (2024b); Guo et al. (2025). These approaches leverage low-rank LoRA adapters in place of full-rank local models to significantly reduce communication and computational overhead, while maintaining strong adaptability in heterogeneous environments.

### 3 AGGREGATION-BROADCAST OPERATOR

**Notation** Let  $W_i^{(t)}$ ,  $B_i^{(t)}$ , and  $A_i^{(t)}$  denote the local model parameters and the LoRA adapter matrices for client  $i$  at step  $t$  (where  $1 \leq i \leq m, 1 \leq t \leq T$ ). Correspondingly, let  $W^{(t)}$ ,  $B^{(t)}$ , and  $A^{(t)}$  represent the global model parameters at step  $t$ . The initial model  $W_0$  refers to the pretrained global model weight. We define the set of global synchronization steps  $\mathcal{I}_E$  as:

$$\mathcal{I}_E = \{nE \mid n \in \mathbb{N}^+\}, \quad (5)$$

where  $E$  denotes the communication interval. If  $t + 1 \in \mathcal{I}_E$ , then step  $t + 1$  corresponds to a communication round. Let  $\mathcal{L}_i(W_i^{(t)}; \xi_{i,t})$  be the local loss function for client  $i$  at step  $t$  with  $\mathbb{E}[\mathcal{L}_i(W_i^{(t)}; \xi_{i,t})] = \mathcal{L}_i(W_i^{(t)})$ , where  $\xi_{i,t}$  is sampled from the  $i$ -th client's local data uniformly at random at the training step  $t$ . Define the global objective as the weighted sum of local losses:

$$\mathcal{L}(W^{(t)}; \xi) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}_i(W^{(t)}; \xi_{i,t}). \quad (6)$$

To describe the aggregation and broadcast process more generally, we define the Aggregation-Broadcast Operator:

**Definition 1 (Aggregation-Broadcast Operator, ABO).** *The Aggregation-Broadcast Operator of the LoRA matrices  $B_i^{(t+1)}$  and  $A_i^{(t+1)}$  is defined by:*

$$\begin{aligned} \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) &:= \mathcal{P}(A_1^{(t+1)}, \dots, A_m^{(t+1)}, B_1^{(t+1)}, \dots, B_m^{(t+1)}) \\ \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) &:= \mathcal{Q}(A_1^{(t+1)}, \dots, A_m^{(t+1)}, B_1^{(t+1)}, \dots, B_m^{(t+1)}) \end{aligned}$$

If  $t + 1 \in \mathcal{I}_E$ , then a communication round is triggered, and each client enters the aggregation-broadcast phase. During this phase, the local LoRA matrices  $B_i^{(t+1)}$  and  $A_i^{(t+1)}$  is updated via Aggregate-Broadcast Operators  $\mathcal{P}$  and  $\mathcal{Q}$ , such that:

$$\begin{aligned} B_i^{(t+1)} &\leftarrow \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \\ A_i^{(t+1)} &\leftarrow \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \end{aligned}$$

for  $1 \leq i \leq m$ .

According to this definition, we establish both SP and PS-Type Aggregation-Broadcast here:

**SP-Type Aggregation-Broadcast.** When FL adopts the SP aggregation, the server first aggregates the global model based on local LoRA adapters  $A_i$  and  $B_i$  by using  $\Delta W = \frac{1}{m} \sum_{i=1}^m B_i A_i$ . After aggregation, the server applies SVD decomposition of  $\frac{1}{m} \sum_{i=1}^m B_i A_i$  to  $\tilde{U} \Sigma \tilde{V}^\top$ , and broadcasts the result to all clients. In this process, the updated local LoRA matrices after the aggregation-broadcast step can be expressed as:

$$B_i^{(t+1)} \leftarrow \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) = \tilde{U}[:, :r] \Sigma[:, :r] \quad (7)$$

$$A_i^{(t+1)} \leftarrow \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) = \tilde{V}^\top[:, :r] \quad (8)$$

**PS-Type Aggregation-Broadcast.** The PS method separately averages  $B_i^{(t+1)}$  and  $A_i^{(t+1)}$ , which means that the global model aggregates by using  $\Delta W = (\frac{1}{m} \sum_{i=1}^m B_i)(\frac{1}{m} \sum_{i=1}^m A_i)$ , then broad-

casting the mean of each:

$$B_i^{(t+1)} \leftarrow \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) = \frac{1}{m} \sum_{i=1}^m B_i^{(t+1)} \quad (9)$$

$$A_i^{(t+1)} \leftarrow \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) = \frac{1}{m} \sum_{i=1}^m A_i^{(t+1)} \quad (10)$$

The algorithm of FL with both SP and PS-Type Aggregation-Broadcast can be seen in A.2. Furthermore, the one-step update for the local LoRA matrices at round  $t + 1$  by such ABO  $\mathcal{P}$  and  $\mathcal{Q}$  can be described as follows:

$$\begin{pmatrix} B_i^{(t)} \\ A_i^{(t)} \end{pmatrix} \xrightarrow{\text{local update}} \begin{pmatrix} B_i^{(t+1)} = B_i^{(t)} - \eta \nabla_B \mathcal{L}_i(W_i^{(t)}; \xi_{i,t}) \\ A_i^{(t+1)} = A_i^{(t)} - \eta \nabla_A \mathcal{L}_i(W_i^{(t)}; \xi_{i,t}) \end{pmatrix} \xrightarrow{\text{if } t+1 \in \mathcal{I}_E} \begin{pmatrix} \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \\ \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \end{pmatrix}. \quad (11)$$

for  $1 \leq i \leq m$ .

## 4 ANALYSIS

In this section, we establish convergence conditions and corresponding theorems for arbitrary Aggregation-Broadcast Operators under certain assumptions. Our analysis is primarily inspired by the work in Zhou & Cong (2017); Guo et al. (2025); Li et al. (2020b). However, unlike the study in Guo et al. (2025), this paper analyzes arbitrary Aggregation-Broadcast Operators under milder conditions.

### 4.1 ASSUMPTION

To facilitate the theoretical analysis of LoRA-based aggregation in FL, we begin by introducing several standard assumptions that are widely used in the FL literature. These assumptions ensure that the local objective functions and model updates behave in a stable and analyzable manner. In particular, we assume the smoothness of the loss functions, uniform boundedness of their gradients, and uniform boundedness of the LoRA matrices during the model update process, which have been widely adopted in many theoretical analyses of FL Li et al. (2020b); Cho et al. (2021); Yu et al. (2019); Guo et al. (2025).

**Assumption 1.**  $\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_m$  are all  $L$ -smooth. For all  $V$  and  $W$ ,

$$\|\nabla \mathcal{L}(V) - \nabla \mathcal{L}(W)\|_F \leq L\|V - W\|_F$$

It is equivalent to

$$\mathcal{L}_i(V) \leq \mathcal{L}_i(W) + \langle V - W, \nabla \mathcal{L}_i(W) \rangle + \frac{L}{2}\|V - W\|_F^2$$

**Assumption 2.** The expected squared norm of the stochastic gradient is uniformly bounded, i.e.,  $\mathbb{E}[\|\nabla \mathcal{L}_i(W_i^{(t)}; \xi_{i,t})\|^2] = \|\nabla \mathcal{L}_i(W_i^{(t)})\|^2 \leq G^2$ , for all  $i = 1, 2, \dots, m$  and  $t = 0, \dots, T - 1$ . Here  $T$  denotes the total number of training steps for each client.

**Assumption 3.** Let  $W_i^{(t)} = W_0 + B_i^{(t)} A_i^{(t)}$ . There exist constants  $C_B > 0$  and  $C_A > 0$  such that  $\|B_i^{(t)}\|_F \leq C_B$ ,  $\|A_i^{(t)}\|_F \leq C_A$  for all  $i = 1, 2, \dots, m$  and  $t = 0, 1, \dots, T - 1$ .

### 4.2 CONVERGENCE RESULTS

We now discuss the convergence result for general Aggregate-Broadcast Operator(ABO). It's worth noting that: (1) we use the same learning rate  $\eta$  during the whole training process and (2) all devices are active. Under these assumptions, we give the following convergence condition and convergence theorem respectively for arbitrary Aggregate-Broadcast Operator(ABO).

#### 4.2.1 WEAK CONVERGENCE CONDITION

In what follows, we begin by introducing a convergence condition under which the local model is guaranteed to converge. We refer the condition as the **weak convergence condition**, which serves as the foundation for establishing convergence guarantee.

**Definition 2 (Weak Convergence Condition).** *The Aggregation-Broadcast Operators(ABO) of LoRA matrices  $B_i^{(t+1)}$  and  $A_i^{(t+1)}$  are said to satisfy the Weak Convergence Condition if there exists a constant  $R > 0$  such that:*

$$\mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} A_i^{(t+1)}\|_F^2 \right] \leq R^2 \eta^2 \quad (12)$$

for  $1 \leq t \leq T$ , where  $\eta$  is learning the rate.

Moreover, we show a sufficient condition of this definition as follows:

**Theorem 1 (Sufficient Condition 1).** *The Aggregation-Broadcast Operators(ABO) of LoRA matrices  $B_i^{(t+1)}$  and  $A_i^{(t+1)}$  satisfy the Weak Convergence Condition if there exists a constant  $R > 0$  such that:*

$$\mathbb{E} \left[ \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} A_i^{(t+1)}\|_F^2 \right] \leq R^2 \eta^2 \quad (13)$$

for  $1 \leq i \leq m$ ,  $1 \leq t \leq T$ , where  $\eta$  is the learning rate.

The details of this sufficient condition can be seen in A.3.1. Based on the Weak Convergence Condition introduced above, we are now in a position to analyze the convergence behavior of the model under this setting. This is a general convergence theorem for arbitrary ABO that satisfy the Weak Convergence Condition.

**Theorem 2 (Weak Convergence Theorem).** *Let Assumption 1, 2 and 3 hold. If the ABO satisfies the Weak Convergence Condition in Definition 2, the update for the local LoRA matrices follows Eq. (11). Then, for a learning rate  $\eta > \xi > 0$  for some  $\xi > 0$ , the gradient of the local loss in expectation satisfies:*

$$\frac{1}{mT} \sum_{i=1}^m \sum_{t=1}^T (\mathbb{E} [\|\nabla_B L_i(W_i^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A L_i(W_i^{(t)})\|_F^2]) \leq \frac{D}{\eta T} + M\eta \quad (14)$$

where  $T$  is the total number of training steps for each client,  $\mathcal{L}_i(W_i^0) - \mathcal{L}_i(W_i^*) \leq D$  for  $\forall i$ ,  $\frac{3}{2}L(\eta^2 C_A^2 C_B^2 G^4 + C_A^4 G^2 + C_B^2 G^2)\eta^2 + C_A C_B G^3 \eta^2 + \frac{L}{2} R^2 \eta^2 + \frac{1}{2}(R^2 + G^2)\eta \leq M\eta^2$ . Specifically, by choosing  $\eta = \sqrt{\frac{D}{MT}}$ , we obtain a convergence rate of  $2\sqrt{\frac{DM}{T}}$ .

The proof of Theorem 2 is provided in the Appendix. A.5. Theorem 2 establishes that if the Aggregation-Broadcast Operator (ABO) satisfies the Weak Convergence Condition, the local model converges to a stationary point within the subspace spanned by  $B$  and  $A$ , achieving a rate of  $\mathcal{O}(1/\sqrt{T})$ . Moreover, the convergence can be accelerated by reducing the values of  $M$ , which increase as  $R$  increases. As a result, the upper bound  $R$  in Definition 2, Eq. (12) directly affects the convergence speed: a larger  $R$  leads to slower convergence, whereas a smaller  $R$  results in faster convergence.

To further simplify the analysis, we ignore the dependency on the round  $t$ . It is known from the Weak Convergence Condition in Definition 2 that minimizing  $R$  is equivalent to solving the following optimization problem:

$$\min_{\mathcal{P}, \mathcal{Q}} \mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} A_i^{(t+1)}\|_F^2 \right] \quad (15)$$

By solving the optimal problem 15, we obtain:

**Corollary 1.** *The Aggregation-Broadcast Operators(ABO)  $\mathcal{P}$  and  $\mathcal{Q}$ , can satisfy the Weak Convergence Condition and achieve the optimal convergence rate of the local model shown in Theorem 2 if the following equation holds:*

$$\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) = \frac{1}{m} \sum_{i=1}^m B_i^{(t+1)} A_i^{(t+1)} \quad (16)$$

for  $1 \leq t \leq T$ , moreover, we can get  $R^2 = 8E^2 G^2 (C_A^4 + C_B^4)$ .

The proof of this corollary can be seen in Appendix A.7. We refer to Eq. (16) as the **optimality condition** under the Weak Convergence Condition. It's obvious that the SP Aggregation Method (As shown in Eq. (3)) satisfies the Eq. (16) in Corollary 1 during the aggregation phase. However, issues arise during the broadcast phase. As shown in Eq. (7) and Eq. (8), which means that:

$$\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) = \tilde{U}[:, : r] \Sigma[:, : r] \tilde{V}^\top[:, : r] \neq \frac{1}{m} \sum_{i=1}^m B_i^{(t+1)} A_i^{(t+1)} \quad (17)$$

These results indicate that the SP aggregate-broadcast strategy cannot achieve the optimal convergence rate due to the broadcast error, also referred to as broadcast loss. Notably, this error becomes more significant as the LoRA rank  $r$  decreases, thereby further slowing down the convergence of the global model. The convergence reaches its optimal rate only when the LoRA rank  $r$  equals the rank of the optimal model. A more comprehensive analysis of this phenomenon will be presented in Section 5.2.

#### 4.2.2 STRONG CONVERGENCE CONDITION

While the Weak Convergence Condition ensures the convergence of local models, it does not necessarily guarantee the convergence of the global model. In practice, the ultimate objective of FL is to achieve stable and efficient convergence of the global model. Therefore, we introduce a stronger convergence condition together with its corresponding theorem. This strong convergence condition imposes stricter requirements on the Aggregation-Broadcast Operator but provides the stronger guarantee that the global model will converge in the subspace spanned by the LoRA matrices  $B$  and  $A$ .

**Definition 3 (Strong Convergence Condition).** *The Aggregation-Broadcast Operators(ABO) of the LoRA matrices  $B_i^{(t+1)}$  and  $A_i^{(t+1)}$  are said to satisfy the Strong Convergence Condition if there exist constants  $P > 0$  and  $Q > 0$  such that:*

$$\mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)}\|_F^2 \right] \leq P^2 \eta^2 \quad (18)$$

$$\mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - A_i^{(t+1)}\|_F^2 \right] \leq Q^2 \eta^2 \quad (19)$$

for  $1 \leq i \leq m, 1 \leq t \leq T$ .

Similarly, we establish a sufficient condition for the Strong Convergence Condition:

**Theorem 3 (Sufficient Condition 2).** *The Aggregation-Broadcast Operators(ABO) of the LoRA matrices  $B_i^{(t+1)}$  and  $A_i^{(t+1)}$  satisfy the Strong Convergence Condition if there exist some constant  $P > 0$  and  $Q > 0$  such that:*

$$\mathbb{E} \left[ \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)}\|_F^2 \right] \leq P^2 \eta^2 \quad (20)$$

$$\mathbb{E} \left[ \|\mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - A_i^{(t+1)}\|_F^2 \right] \leq Q^2 \eta^2 \quad (21)$$

for  $1 \leq i \leq m, 1 \leq t \leq T$ . Where  $\eta$  is the learning rate.

The details of this sufficient condition can be seen in A.3.2. We define the global weight in step  $t$  as  $W^{(t)} = W_0 + \mathcal{P}(A_{1 \leq j \leq m}^{(t)}, B_{1 \leq j \leq m}^{(t)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t)}, B_{1 \leq j \leq m}^{(t)})$ . Then we obtain the following Strong Convergence Theorem under the Strong Convergence Condition.

**Theorem 4 (Strong Convergence Theorem).** *Let Assumption 1, 2 and 3 hold. If the ABO satisfies the Strong Convergence Condition in Definition 3, the update for the local LoRA matrices is followed by Eq. (11). Then, for a learning rate  $\eta > \xi > 0$  for some  $\xi > 0$ , the gradient of the global loss in expectation satisfy:*

$$\frac{1}{T} \sum_{t=1}^T (\mathbb{E} [\|\nabla_B \mathcal{L}(W^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A \mathcal{L}(W^{(t)})\|_F^2]) \leq 2(\frac{D}{\eta T} + (M + N)\eta) \quad (22)$$

for  $1 \leq i \leq m$ ,  $1 \leq t \leq T$ , where  $\mathcal{L}_i(W_0) - \mathcal{L}_i(W_i^*) \leq D$ ,  $R^2 = 4P^2Q^2\eta^2 + 3C_B^2Q^2 + 3C_A^2P^2$ ,  $4G^2(Q^2 + P^2)\eta^2 + 4(C_A^2 + C_B^2)L^2R^2\eta^2 \leq 2N\eta$  and  $\frac{3}{2}L(\eta^2C_A^2C_B^2G^4 + C_A^4G^2 + C_B^2G^2)\eta^2 + C_AC_BG^3\eta^2 + \frac{L}{2}R^2\eta^2 + \frac{1}{2}(R^2 + G^2)\eta \leq M\eta^2$ . Specifically, by choosing  $\eta = \sqrt{\frac{D}{(M+N)T}}$ , we obtain a convergence rate of  $4\sqrt{\frac{D(M+N)}{T}}$ .

The proof of this theorem is provided in Appendix A.8. Similarly, if we ignore the dependency on  $t$ , minimize  $R$  under Strong Convergence Condition in Definition 3 is equivalent to solve the following convex-optimization problem Boyd & Vandenberghe (2004):

$$\min_{\mathcal{P}} \mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)}\|_F^2 \right] \quad (23)$$

$$\min_{\mathcal{Q}} \mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - A_i^{(t+1)}\|_F^2 \right] \quad (24)$$

By solving the optimization problem, we can get the corollary:

**Corollary 2.** The Aggregation-Broadcast Operators (ABO)  $\mathcal{P}$  and  $\mathcal{Q}$  can satisfy the Strong Convergence Condition in Definition 3 and achieve the optimal convergence rate of the global model shown in Theorem 4 if the following equation holds:

$$\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) = \frac{1}{m} \sum_{i=1}^m B_i^{(t+1)} \quad (25)$$

$$\mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) = \frac{1}{m} \sum_{i=1}^m A_i^{(t+1)} \quad (26)$$

for  $1 \leq t \leq T$ , where  $P^2 = 4E^2G^2C_A^4$ ,  $Q^2 = 4E^2G^2C_B^2$ , which lead to  $R^2 = 64E^4G^4C_A^2C_B^2\eta^2 + 12E^2G^2(Q^2C_B^4 + P^2C_A^4)$ .

The proof of this corollary can be seen in A.9. We refer to Eq. (25) and (26) as the **optimality condition** under the Strong Convergence Condition. It is worth noting that the PS Aggregation Method satisfies this optimality condition. As a result, the PS Aggregation Method is capable of achieving the optimal convergence rate of the global model if all clients share the same LoRA rank, which suggests that the PS aggregation method is relatively robust to the choice of LoRA rank. We will provide a more detailed analysis of this issue in Section 5.2.

## 5 EXPERIMENT AND EVALUATION

In this section, we introduce the experiments we conduct to verify our proposed theory.

### 5.1 EXPERIMENT SETUP

We verify the proposed convergence theory using a Multi-Layer Perceptron (MLP) on the MNIST, FMNIST, QMNIST and KMNIST datasets with 10 clients under a highly non-IID distribution, where each client holds all samples from only one class. We pick representative method FlexLoRA Bai et al. (2024) and RBLA Chen et al. (2024a) to demonstrate SP and PS, respectively. All experiments are conducted with a fixed random seed 42. Next, we introduce some important parameters.

- **Rank scale ratio  $\delta$ :** The LoRA rank of each layer scaled by  $\delta$  as  $\max(\lfloor c \cdot \delta \rfloor, 1)$  Chen et al. (2024a), where  $c$  is a layer-specific constant. In our model, the based rank  $c$  for each layer is set to 160, 160, and 100, respectively.
- **Total number steps  $T$ :** Throughout the training process, we specify the total number of steps trained by each client as  $T$ , where  $T = \text{epoch} \cdot \text{rounds}$ .

The following experiments are mainly designed to examine the sensitivity of PS-ABO and SP-ABO to different rank ratios and epochs. Based on our theoretical analysis, we can make the following observations:



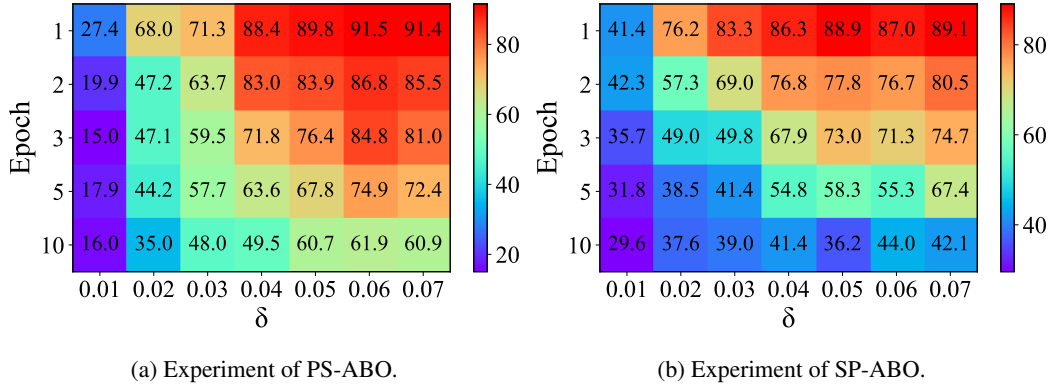


Figure 1: The highest accuracies (%) of PS-ABO and SP-ABO at different rank ratios and epochs when the total number of steps  $T = 300$  (Epoch \* Round) on the MNIST dataset. Each cell represents the highest test accuracy the global model can reach in experiments with different  $\delta$  and epoch.

**Different Rank.** We know that the smaller the rank of a model, the weaker its expressive power tends to be. Consequently, as the rank decreases, the model’s accuracy also drops. However, by Eq. 17, we know that due to the presence of broadcasting errors (caused by SVD), SP-ABO is more sensitive to the rank.

**Different Epoch.** From Corollary 1 and Corollary 2 we know that  $R_{sp}^2 = 8E^2G^2(C_A^4 + C_B^4)$ ,  $R_{ps}^2 = 64E^4G^4C_A^2C_B^2\eta^2 + 12E^2G^2(Q^2C_B^4 + P^2C_A^4)$ . These expressions indicate that, under a fixed total number of total steps  $T$  (i.e. total number of epochs) and fixed random seed (which means the fixed  $G$ ,  $C_A$  and  $C_B$ ), increasing the number of local epochs  $E$  per communication round leads to a larger value of  $R$ , and thus slower convergence for both SP-ABO and PS-ABO. In other words, when  $T$  is fixed, performing more local updates between communication rounds degrades overall convergence performance. Moreover, we observe that  $R_{sp}^2 = \mathcal{O}(E^2)$  and  $R_{ps}^2 = \mathcal{O}(E^4)$ , which clearly shows that PS-ABO is more sensitive to the choice of epoch number.

## 5.2 EXPERIMENT RESULT

As shown in Fig. 1, for example, at epoch 1, as the rank decreases from 0.07 to 0.01, SP-ABO drops sharply from 91.4% to 27.4%, whereas PS-ABO decreases more moderately from 89.1% to 41.4%. A similar phenomenon can be observed across different epochs. This contrast demonstrates that SP-ABO is considerably more sensitive to rank reduction than PS-ABO. Moreover, when  $\delta = 0.01$ , as the epoch increases from 1 to 10, the accuracy of SP-ABO decreases by approximately 30 percentage points (from 91.4% to 60.9%), whereas that of PS-ABO decreases by up to 47 percentage points (from 89.1% to 42.1%). A similar trend is observed across different rank ratios.

Overall, these results illustrate that PS-ABO is more sensitive to the number of local epochs, while SP-ABO is more sensitive to the LoRA rank. These results are consistent with our theoretical analysis. Detailed experiment setup and additional experiments can be seen in Appendix. A.10.

## 6 CONCLUSION

In this paper, we presented a unified theoretical framework for analyzing the convergence behavior of LoRA-based FL. By introducing the concept of Aggregation-Broadcast Operators, we established a general convergence condition along with several sufficient conditions. Our framework not only provides convergence guarantees for the widely used SP-ABO and PS-ABO, but also offers insights into designing new aggregation methods with provable performance. Extensive experiments on standard benchmarks corroborate our theoretical findings. The findings contribute to a clearer understanding of LoRA in federated scenarios and may assist in developing more efficient and reliable model adaptation strategies.

## REFERENCES

- Jiamu Bai, Daoyuan Chen, Bingchen Qian, Liuyi Yao, and Yaliang Li. Federated fine-tuning of large language models under heterogeneous tasks and client resources. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=gkOzoHBXUw>.
- Jieming Bian, Lei Wang, Letian Zhang, and Jie Xu. Lora-fair: Federated lora fine-tuning with aggregation and initialization refinement. *arXiv preprint arXiv:2411.14961*, 2024.
- Stephen P Boyd and Lieven Vandenbergh. *Convex optimization*. Cambridge university press, 2004.
- Shuaijun Chen, Omid Tavallaie, Niousha Nazemi, and Albert Y Zomaya. Rbla: Rank-based-lora-aggregation for fine-tuning heterogeneous models in flaaS. In *International Conference on Web Services*, pp. 47–62. Springer, 2024a.
- Shuangyi Chen, Yue Ju, Hardik Dalal, Zhongwen Zhu, and Ashish Khisti. Robust federated finetuning of foundation models via alternating minimization of lora. *arXiv preprint arXiv:2409.02346*, 2024b.
- Yae Jee Cho, Jianyu Wang, and Gauri Joshi. Client selection in federated learning: Convergence analysis and power-of-choice selection strategies, 2021. URL <https://openreview.net/forum?id=PYAFKBc8GL4>.
- Yae Jee Cho, Luyang Liu, Zheng Xu, Aldi Fahrezi, Matt Barnes, and Gauri Joshi. Heterogeneous loRA for federated fine-tuning of on-device foundation models. In *International Workshop on Federated Learning in the Age of Foundation Models in Conjunction with NeurIPS 2023*, 2023. URL <https://openreview.net/forum?id=EmV9sGpZ7q>.
- Yae Jee Cho, Luyang Liu, Zheng Xu, Aldi Fahrezi, and Gauri Joshi. Heterogeneous LoRA for federated fine-tuning of on-device foundation models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 12903–12913, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.717. URL <https://aclanthology.org/2024.emnlp-main.717/>.
- Shuiguang Deng, Hailiang Zhao, Weijia Fang, Jianwei Yin, Schahram Dustdar, and Albert Y. Zomaya. Edge intelligence: The confluence of edge computing and artificial intelligence. *IEEE Internet of Things Journal*, 7(8):7457–7469, 2020. doi: 10.1109/JIOT.2020.2984887.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- Pengxin Guo, Shuang Zeng, Yanran Wang, Huijie Fan, Feifei Wang, and Liangqiong Qu. Selective aggregation for low-rank adaptation in federated learning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=iX3uESGdsO>.
- Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 328–339, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1031. URL <https://aclanthology.org/P18-1031/>.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022a.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022b. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.

- Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and trends® in machine learning*, 14(1–2):1–210, 2021.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *ICML*, 2020.
- Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- Jabin Koo, Minwoo Jang, and Jungseul Ok. Towards robust and efficient federated low-rank adaptation with heterogeneous clients. *arXiv preprint arXiv:2410.22815*, 2024.
- Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In *CVPR*, 2021.
- Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. Fair resource allocation in federated learning. *arXiv preprint arXiv:1905.10497*, 2019.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, 2:429–450, 2020a.
- Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. In *International Conference on Learning Representations*, 2020b. URL <https://openreview.net/forum?id=HJxNANVtDS>.
- Sadhika Malladi, Alexander Wettig, Dingli Yu, Danqi Chen, and Sanjeev Arora. A kernel-based view of language model fine-tuning. In *International Conference on Machine Learning*, pp. 23610–23641. PMLR, 2023.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.
- Ngoc-Hieu Nguyen, Tuan-Anh Nguyen, Tuan Minh Nguyen, Vu Tien Hoang, Dung D. Le, and Kok-Seng Wong. Towards efficient communication and secure federated recommendation system via low-rank training. In *The Web Conference 2024*, 2024. URL <https://openreview.net/forum?id=GBIyyK22Cf>.
- Huangsiyuan Qin and Ying Li. Fedinc: One-shot federated tuning for collaborative incident recognition. In *International Conference on Artificial Neural Networks*, pp. 174–185. Springer, 2024.
- Raghav Singhal, Kaustubh Ponkshe, and Praneeth Vepakomma. Fedex-loRA: Exact aggregation for federated and efficient fine-tuning of foundation models, 2024. URL <https://openreview.net/forum?id=Yg998afEbH>.
- Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet Talwalkar. Federated multi-task learning. In *NeurIPS*, 2017.
- Youbang Sun, Zitao Li, Yaliang Li, and Bolin Ding. Improving loRA in privacy-preserving federated learning. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=NLPzL6HWNl>.
- Jihun Wang, Qing Liu, Haoran Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. In *NeurIPS*, 2020.
- Ziyao Wang, Zheyu Shen, Yexiao He, Guoheng Sun, Hongyi Wang, Lingjuan Lyu, and Ang Li. Flora: Federated fine-tuning large language models with heterogeneous low-rank adaptations. *Advances in Neural Information Processing Systems*, 37:22513–22533, 2024.

- Xinghao Wu, Xuefeng Liu, Jianwei Niu, Haolin Wang, Shaojie Tang, and Guogang Zhu. FedLoRA: When personalized federated learning meets low-rank adaptation, 2024. URL <https://openreview.net/forum?id=bZh06ptG9r>.
- Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhengsu Chen, XIAOPENG ZHANG, and Qi Tian. QA-LoRA: Quantization-aware low-rank adaptation of large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=WvFoJccpo8>.
- Yiyuan Yang, Guodong Long, Qinghua Lu, Liming Zhu, Jing Jiang, and Chengqi Zhang. Federated low-rank adaptation for foundation models: A survey. *arXiv preprint arXiv:2505.13502*, 2025.
- Liping Yi, Han Yu, Gang Wang, and Xiaoguang Liu. Fedlora: Model-heterogeneous personalized federated learning with lora tuning. *arXiv preprint arXiv:2310.13283*, 2023.
- Hao Yu, Sen Yang, and Shenghuo Zhu. Parallel restarted sgd with faster convergence and less communication: demystifying why model averaging works for deep learning. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence*, AAAI’19/IAAI’19/EAAI’19. AAAI Press, 2019. ISBN 978-1-57735-809-1. doi: 10.1609/aaai.v33i01.33015693. URL <https://doi.org/10.1609/aaai.v33i01.33015693>.
- Yuchen Zeng and Kangwook Lee. The expressive power of low-rank adaptation. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=likXVjmh3E>.
- Jianyi Zhang, Saeed Vahidian, Martin Kuo, Chunyuan Li, Ruiyi Zhang, Tong Yu, Guoyin Wang, and Yiran Chen. Towards building the federatedgpt: Federated instruction tuning. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6915–6919, 2024a. doi: 10.1109/ICASSP48485.2024.10447454.
- Jianyi Zhang, Saeed Vahidian, Martin Kuo, Chunyuan Li, Ruiyi Zhang, Tong Yu, Guoyin Wang, and Yiran Chen. Towards building the federatedgpt: Federated instruction tuning. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6915–6919. IEEE, 2024b.
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=lq62uWRJjiY>.
- Ruiyi Zhang, Rushi Qiang, Sai Ashish Somayajula, and Pengtao Xie. AutoLoRA: Automatically tuning matrix ranks in low-rank adaptation based on meta learning. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 5048–5060, Mexico City, Mexico, June 2024c. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.282. URL <https://aclanthology.org/2024.naacl-long.282/>.
- Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.
- Fan Zhou and Guojing Cong. On the convergence properties of a  $k$ -step averaging stochastic gradient descent algorithm for nonconvex optimization. *arXiv preprint arXiv:1708.01012*, 2017.
- Jiacheng Zhu, Kristjan H. Greenewald, Kimia Nadjahi, Haitz Sáez de Ocáriz Borde, Rickard Brühl Gabrielsson, Leshem Choshen, Marzyeh Ghassemi, Mikhail Yurochkin, and Justin Solomon. Asymmetry in low-rank adapters of foundation models. In *ICML*, 2024. URL <https://openreview.net/forum?id=txRZBD8tBV>.

## A APPENDIX

### A.1 ADDITIONAL NOTATION

We will adopt the notation used in Guo et al. (2025). Let  $W_i^{(t)}, B_i^{(t)}$  and  $A_i^{(t)}$  denote the local model and local LoRA matrices correspondingly for client  $i$  at step  $t$ . Similarly, let  $W^{(t)}, B^{(t)}$  and  $A^{(t)}$  represent the global model parameters at step  $t$ . Let  $W_0$  be the pre-trained model. Define  $\mathcal{I}_E = \{nE \mid n \in \mathbb{N}^+\}$ , which means that if  $t+1 \in \mathcal{I}_E$ , it indicates that the current step is a communication round. Moreover, we denote the Aggregation Broadcast Operator as follows :

$$\begin{aligned} P(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) &= P(A_1^{(t+1)}, \dots, A_m^{(t+1)}, B_1^{(t+1)}, \dots, B_m^{(t+1)}) \\ Q(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) &= Q(A_1^{(t+1)}, \dots, A_m^{(t+1)}, B_1^{(t+1)}, \dots, B_m^{(t+1)}) \end{aligned}$$

We define  $U_i^{(t)} = W_0 + B_i^{(t+1)} A_i^{(t+1)}$  as the immediate result of the local update, which captures the information during the local training phase. Define  $V_i^{(t)} = W_0 + P(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) Q(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)})$  as the aggregated update, which captures the information from the aggregation phase. Let  $W_i^{(t)} = W_0 + B_i^{(t)} A_i^{(t)}$  be the final update. Then we can get:

$$W_i^{(t+1)} = \begin{cases} U_i^{(t)} & \text{if } t+1 \notin \mathcal{I}_E, \\ V_i^{(t)} & \text{if } t+1 \in \mathcal{I}_E. \end{cases}$$

### A.2 SP AND PS AGGREGATION-BROADCAST

---

#### Algorithm 1 Federated Learning with SP Aggregation-Broadcast Operator

---

**Require:** Number of clients  $m$ , total steps  $T$ , local epochs  $E$ , learning rate  $\eta$

1: Initialize LoRA matrices  $(A_i^{(0)}, B_i^{(0)}) = (A_{\text{initial}}, B_{\text{initial}})$  for  $1 \leq i \leq m$

2: **for**  $t = 0, 1, \dots, T-1$  **do**

3:   **for** each client  $i \in \{1, \dots, m\}$  in parallel **do**

4:     Sample mini-batch  $\xi_{i,t}$

5:     Local Update:

$$B_i^{(t+1)} \leftarrow B_i^{(t)} - \eta \nabla_B \mathcal{L}_i(W_i^{(t)}, \xi_{i,t})$$

$$A_i^{(t+1)} \leftarrow A_i^{(t)} - \eta \nabla_A \mathcal{L}_i(W_i^{(t)}, \xi_{i,t})$$

6:   **end for**

7:   **if**  $t+1 \in \mathcal{I}_E$  **then**

8:     **Server aggregates:**

$$W^{(t)} = W_0 + \sum_{i=1}^m B_i^{(t+1)} A_i^{(t+1)}$$

9:      $[\tilde{U}, \Sigma, \tilde{V}^T] = \text{SVD}(\sum_{i=1}^m B_i^{(t+1)} A_i^{(t+1)})$

10:     Broadcast  $(B_i^{(t+1)}, A_i^{(t+1)}) = (\tilde{U}[:, :r], \Sigma[:, :r], \tilde{V}^T[:, :r])$  for  $1 \leq i \leq m$

11:   **end if**

12: **end for**

---

**Algorithm 2** Federated Learning with PS Aggregation-Broadcast Operator**Require:** Number of clients  $m$ , total steps  $T$ , local epochs  $E$ , learning rate  $\eta$ 

```

1: Initialize LoRA matrices  $(A_i^{(0)}, B_i^{(0)}) = (A_{\text{initial}}, B_{\text{initial}})$  for  $1 \leq i \leq m$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   for each client  $i \in \{1, \dots, m\}$  in parallel do
4:     Sample mini-batch  $\xi_{i,t}$ 
5:     Local Update:
6:        $B_i^{(t+1)} \leftarrow B_i^{(t)} - \eta \nabla_B \mathcal{L}_i(W_i^{(t)}, \xi_{i,t})$ 
7:        $A_i^{(t+1)} \leftarrow A_i^{(t)} - \eta \nabla_A \mathcal{L}_i(W_i^{(t)}, \xi_{i,t})$ 
8:   end for
9:   if  $t + 1 \in \mathcal{I}_E$  then
10:    Server aggregates:
11:     $W^{(t)} = W_0 + (\sum_{i=1}^m B_i^{(t+1)}) (\sum_{i=1}^m A_i^{(t+1)})$ 
12:    Broadcast  $(B_i^{(t+1)}, A_i^{(t+1)}) = (\sum_{j=1}^m B_j^{(t+1)}, \sum_{j=1}^m A_j^{(t+1)})$  for  $1 \leq i \leq m$ 
13:   end if
14: end for

```

## A.3 SUFFICIENT CONDITIONS

## A.3.1 SUFFICIENT CONDITION OF WEAK CONVERGENCE CONDITION

We recall the sufficient condition of the Weak Convergence Condition as follows:

$$\mathbb{E} \left[ \left\| \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} A_i^{(t+1)} \right\|_F^2 \right] \leq R^2 \eta^2 \quad (27)$$

Compared to the Weak Convergence Condition in Definition 2, Eq. 27 provides a more practical and easily verifiable sufficient condition. Theorem 1 can be directly proven by:

$$\frac{1}{m} \sum_{i=1}^m \left\| \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} A_i^{(t+1)} \right\|_F^2 \leq \frac{1}{m} \sum_{i=1}^m R^2 \eta^2 \leq R^2 \eta^2 \quad (28)$$

Based on Theorem 1, we can immediately derive that the local model also converges at the rate of  $\mathcal{O}(1/\sqrt{T})$ , as shown in Eq. (15). Moreover, minimizing the constant  $R$  in Theorem 1 leads to improved convergence. Actually, if we ignore the dependency on the communication round  $t$ , this minimization is equivalent to solving the following optimization problem:

$$\min_{\mathcal{P}, \mathcal{Q}} \max_{1 \leq i \leq m} \mathbb{E} \left[ \left\| \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} A_i^{(t+1)} \right\|_F^2 \right] \quad (29)$$

This formulation provides a clear objective for designing effective Aggregation-Broadcast Operators by minimizing the bound  $R$ .

## A.3.2 SUFFICIENT CONDITION OF STRONG CONVERGENCE CONDITION

We now turn to another sufficient condition. In contrast to the Strong Convergence Condition, it offers a more concise formulation and is more amenable to verification:

$$\mathbb{E} \left[ \left\| \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} \right\|_F^2 \right] \leq P^2 \eta^2 \quad (30)$$

$$\mathbb{E} \left[ \left\| \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - A_i^{(t+1)} \right\|_F^2 \right] \leq Q^2 \eta^2 \quad (31)$$

Then Theorem 3 can be easily proved by:

$$\mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)}\|_F^2 \right] \leq \frac{1}{m} \sum_{i=1}^m P^2 \eta^2 = P^2 \eta^2 \quad (32)$$

$$\mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - A_i^{(t+1)}\|_F^2 \right] \leq \frac{1}{m} \sum_{i=1}^m Q^2 \eta^2 = Q^2 \eta^2 \quad (33)$$

If we ignore the dependency on the round  $t$ , minimize  $R$  under Sufficient Condition 3 in Theorem 3 is equivalent to solve the following convex-optimal problem:

$$\min_{\mathcal{P}} \max_{1 \leq i \leq m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)}\|_F^2 \quad (34)$$

$$\min_{\mathcal{Q}} \max_{1 \leq i \leq m} \sum_{i=1}^m \|\mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - A_i^{(t+1)}\|_F^2 \quad (35)$$

#### A.4 KEY LEMMAS

To ensure a concise and clear proof of the theorem, we first provide the following lemmas. Detailed proofs of these lemmas will be presented after completing the proof of Theorem 2.

**Lemma 1.** Assume Assumption 2 and 3 hold. We have:

$$\mathbb{E} \left[ \langle U_i^{(t)} - W_i^{(t)}, \nabla_W \mathcal{L}_i(W_i^{(t)}) \rangle_F \right] \leq \eta_{i,t}^2 C_A C_B G^3 - \eta_{i,t} (\mathbb{E} [\|\nabla_B L_i(W_i^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A L_i(W_i^{(t)})\|_F^2])$$

**Lemma 2.** Assume Assumption 2 holds. We have:

$$\mathbb{E} [\|U_i^{(t)} - W_i^{(t)}\|_F^2] \leq 3\eta_{i,t}^4 C_A^2 C_B^2 G^4 + 3\eta_{i,t}^2 C_A^4 G^2 + 3\eta_{i,t}^2 C_B^4 G^2$$

#### A.5 PROOF OF THEOREM 2

*Proof.* First, by Assumption 1, we have:

$$\mathcal{L}_i(U_i^{(t)}) \leq \mathcal{L}_i(W_i^{(t)}) + \langle U_i^{(t)} - W_i^{(t)}, \nabla_W \mathcal{L}_i(W_i^{(t)}) \rangle_F + \frac{L}{2} \|U_i^{(t)} - W_i^{(t)}\|_F^2 \quad (36)$$

$$\mathcal{L}_i(V_i^{(t)}) \leq \mathcal{L}_i(U_i^{(t)}) + \langle V_i^{(t)} - U_i^{(t)}, \nabla_W \mathcal{L}_i(U_i^{(t)}) \rangle_F + \frac{L}{2} \|V_i^{(t)} - U_i^{(t)}\|_F^2 \quad (37)$$

If  $t+1 \notin \mathcal{I}_E$ , by lemma 1 and Eq. (36) we have:

$$\begin{aligned} \mathbb{E} [\mathcal{L}_i(W_i^{(t+1)})] &= \mathbb{E} [\mathcal{L}_i(U_i^{(t)})] \\ &\leq \mathbb{E} [\mathcal{L}_i(W_i^{(t)})] + \mathbb{E} [\langle U_i^{(t)} - W_i^{(t)}, \nabla_W \mathcal{L}_i(W_i^{(t)}) \rangle_F] + \frac{L}{2} \mathbb{E} [\|U_i^{(t)} - W_i^{(t)}\|_F^2] \\ &\leq \mathbb{E} [\mathcal{L}_i(W_i^{(t)})] + \frac{L}{2} \mathbb{E} [\|U_i^{(t)} - W_i^{(t)}\|_F^2] + \eta_{i,t}^2 C_A C_B G^3 \\ &\quad - \eta_{i,t} (\mathbb{E} [\|\nabla_B L_i(W_i^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A L_i(W_i^{(t)})\|_F^2]) \end{aligned} \quad (38)$$

If  $t + 1 \in \mathcal{I}_E$ , by lemma 1 and Eq. (36) and (37), we have:

$$\begin{aligned}
\mathbb{E} [\mathcal{L}_i(W_i^{(t+1)})] &= \mathbb{E} [\mathcal{L}_i(V_i^{(t)})] \\
&\leq \mathbb{E} [\mathcal{L}_i(U_i^{(t)})] + \mathbb{E} [\langle V_i^{(t)} - U_i^{(t)}, \nabla_W \mathcal{L}_i(U_i^{(t)}) \rangle_F] + \frac{L}{2} \mathbb{E} [\|V_i^{(t)} - U_i^{(t)}\|_F^2] \\
&\leq \mathbb{E} [\mathcal{L}_i(W_i^{(t)})] + \mathbb{E} [\langle U_i^{(t)} - W_i^{(t)}, \nabla_W \mathcal{L}_i(W_i^{(t)}) \rangle_F] + \frac{L}{2} \mathbb{E} [\|U_i^{(t)} - W_i^{(t)}\|_F^2] \\
&\quad + \mathbb{E} [\langle V_i^{(t)} - U_i^{(t)}, \nabla_W \mathcal{L}_i(U_i^{(t)}) \rangle_F] + \frac{L}{2} \mathbb{E} [\|V_i^{(t)} - U_i^{(t)}\|_F^2] \\
&\leq \mathbb{E} [\mathcal{L}_i(W_i^{(t)})] + \frac{L}{2} \mathbb{E} [\|U_i^{(t)} - W_i^{(t)}\|_F^2] + \eta_{i,t}^2 C_A C_B G^3 \\
&\quad - \eta_{i,t} (\mathbb{E} [\|\nabla_B L_i(W_i^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A L_i(W_i^{(t)})\|_F^2]) \\
&\quad + \mathbb{E} [\langle V_i^{(t)} - U_i^{(t)}, \nabla_W \mathcal{L}_i(U_i^{(t)}) \rangle_F] + \frac{L}{2} \mathbb{E} [\|V_i^{(t)} - U_i^{(t)}\|_F^2] \tag{39}
\end{aligned}$$

By comparing Eq. (38) and Eq. (39), we observe that Equation (6) holds universally for arbitrary values of  $t$ . So we get, for  $\forall t$  and  $\eta_{i,j} = \eta$ ,

$$\begin{aligned}
\mathbb{E} [\mathcal{L}_i(W_i^{(t+1)})] &\leq \mathbb{E} [\mathcal{L}_i(W_i^{(t)})] + \frac{L}{2} \mathbb{E} [\|U_i^{(t)} - W_i^{(t)}\|_F^2] + \eta^2 C_A C_B G^3 \\
&\quad - \eta (\mathbb{E} [\|\nabla_B L_i(W_i^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A L_i(W_i^{(t)})\|_F^2]) \\
&\quad + \mathbb{E} [\langle V_i^{(t)} - U_i^{(t)}, \nabla_W \mathcal{L}_i(U_i^{(t)}) \rangle_F] + \frac{L}{2} \mathbb{E} [\|V_i^{(t)} - U_i^{(t)}\|_F^2] \tag{40}
\end{aligned}$$

Moreover, we have:

$$\begin{aligned}
\eta (\mathbb{E} [\|\nabla_B L_i(W_i^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A L_i(W_i^{(t)})\|_F^2]) &\leq \mathbb{E} [\mathcal{L}_i(W_i^{(t)}) - \mathcal{L}_i(W_i^{(t+1)})] + \frac{L}{2} \mathbb{E} [\|U_i^{(t)} - W_i^{(t)}\|_F^2] + \eta^2 C_A C_B G^3 \\
&\quad + \mathbb{E} [\langle V_i^{(t)} - U_i^{(t)}, \nabla_W \mathcal{L}_i(U_i^{(t)}) \rangle_F] + \frac{L}{2} \mathbb{E} [\|V_i^{(t)} - U_i^{(t)}\|_F^2] \tag{41}
\end{aligned}$$

Since:

$$\begin{aligned}
\langle V_i^{(t)} - U_i^{(t)}, \nabla_W \mathcal{L}_i(U_i^{(t)}) \rangle_F &\leq \frac{1}{2} \left[ \frac{1}{\eta} \|V_i^{(t)} - U_i^{(t)}\|_F^2 + \eta \|\nabla_W \mathcal{L}_i(U_i^{(t)})\|_F^2 \right] \\
&\leq \frac{1}{2\eta} \|V_i^{(t)} - U_i^{(t)}\|_F^2 + \frac{\eta}{2} G^2 \tag{42}
\end{aligned}$$

The last inequality in 42 holds by Assumption 2. By Eq. (41) and (42), we have:

$$\begin{aligned}
\mathbb{E} [\|\nabla_B L_i(W_i^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A L_i(W_i^{(t)})\|_F^2] &\leq \frac{\mathbb{E} [\mathcal{L}_i(W_i^{(t)}) - \mathcal{L}_i(W_i^{(t+1)})]}{\eta} + \frac{1}{\eta} \left( \frac{L}{2} \mathbb{E} [\|U_i^{(t)} - W_i^{(t)}\|_F^2] \right. \\
&\quad \left. + \eta^2 C_A C_B G^3 + \left( \frac{1}{2\eta} + \frac{L}{2} \right) \mathbb{E} [\|V_i^{(t)} - U_i^{(t)}\|_F^2] + \frac{\eta}{2} G^2 \right) \tag{43}
\end{aligned}$$



Summing Eq. (43) over  $i = 1$  to  $m$ , by lemma. 2 and the convergence condition in Definition. 2, Eq. (43) can be transformed to:

$$\begin{aligned}
& \frac{1}{m} \sum_{i=1}^m (\mathbb{E} [\|\nabla_B L_i(W_i^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A L_i(W_i^{(t)})\|_F^2]) \\
& \leq \frac{1}{m} \sum_{i=1}^m \frac{\mathbb{E} [\mathcal{L}_i(W_i^{(t)}) - \mathcal{L}_i(W_i^{(t+1)})]}{\eta} + \frac{1}{\eta} \left( \frac{L}{2} \cdot \frac{1}{m} \sum_{i=1}^m \mathbb{E} [\|U_i^{(t)} - W_i^{(t)}\|_F^2] \right. \\
& \quad \left. + \eta^2 C_A C_B G^3 + \left( \frac{1}{2\eta} + \frac{L}{2} \right) \cdot \frac{1}{m} \sum_{i=1}^m \mathbb{E} [\|V_i^{(t)} - U_i^{(t)}\|_F^2] + \frac{\eta}{2} G^2 \right) \\
& \leq \frac{1}{m} \sum_{i=1}^m \frac{\mathbb{E} [\mathcal{L}_i(W_i^{(t)}) - \mathcal{L}_i(W_i^{(t+1)})]}{\eta} + \frac{1}{\eta} \left( \frac{L}{2} (3\eta^2 C_A^2 C_B^2 G^4 + 3C_A^4 G^2 \right. \\
& \quad \left. + 3C_B^4 G^2) \eta^2 + \eta^2 C_A C_B G^3 + \left( \frac{1}{2\eta} + \frac{L}{2} \right) R^2 \eta^2 + \frac{\eta}{2} G^2 \right) \\
& \leq \frac{1}{m} \sum_{i=1}^m \frac{\mathbb{E} [\mathcal{L}_i(W_i^{(t)}) - \mathcal{L}_i(W_i^{(t+1)})]}{\eta} + \frac{1}{\eta} \cdot M \eta^2 \\
& \leq \frac{1}{m} \sum_{i=1}^m \frac{\mathbb{E} [\mathcal{L}_i(W_i^{(t)}) - \mathcal{L}_i(W_i^{(t+1)})]}{\eta} + M \eta
\end{aligned} \tag{44}$$

Where  $\frac{3}{2} L (\eta^2 C_A^2 C_B^2 G^4 + C_A^4 G^2 + C_B^4 G^2) \eta^2 + C_A C_B G^3 \eta^2 + \frac{L}{2} R^2 \eta^2 + \frac{1}{2} (R^2 + G^2) \eta \leq M \eta^2$ , it holds if there exist some constant  $\epsilon$  such that  $\eta > \epsilon > 0$ . The second inequality holds by lemma. 2 and the convergence condition in Definition. 2. Summing Eq. (44) over  $t = 1$  to  $T$  yields:

$$\begin{aligned}
\frac{1}{mT} \sum_{t=1}^T \sum_{i=1}^m (\mathbb{E} [\|\nabla_B L_i(W_i^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A L_i(W_i^{(t)})\|_F^2]) & \leq \frac{1}{m} \sum_{i=1}^m \frac{\mathbb{E} [\mathcal{L}_i(W_i^{(0)}) - \mathcal{L}_i(W_i^{(T+1)})]}{\eta T} + M \eta \\
& \leq \frac{1}{m} \sum_{i=1}^m \frac{\mathbb{E} [\mathcal{L}_i(W_i^{(0)}) - \mathcal{L}_i(W_i^*)]}{\eta T} + M \eta \\
& \leq \frac{D}{\eta T} + M \eta
\end{aligned} \tag{45}$$

Where  $W_i^* = \operatorname{argmin}_W \mathcal{L}_i(W)$ ,  $W_i^{(0)} = W_0$  and we assume that  $\mathcal{L}_i(W_0) - \mathcal{L}_i(W_i^*) \leq D$  for  $\forall i$ . Based on the preceding reasoning and let  $\eta = \sqrt{\frac{D}{MT}}$ , we arrive at the following inequality:

$$\frac{1}{mT} \sum_{i=1}^m \sum_{t=1}^T (\mathbb{E} [\|\nabla_B L_i(W_i^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A L_i(W_i^{(t)})\|_F^2]) \leq 2\sqrt{\frac{DM}{T}} \tag{46}$$

□

## A.6 PROOF OF LEMMAS

Though Lemmas 1 and 2 replicate the proof methodology of Guo et al. (2025), we retain their proofs here to ensure a self-contained theoretical presentation, particularly given their critical role in deriving Theorem 2.

## A.6.1 PROOF OF LEMMA 1

*Proof.* We know that:

$$\begin{aligned}
U_i^{(t)} - W_i^{(t)} &= B_i^{(t+1)} A_i^{(t+1)} - B_i^{(t)} A_i^{(t)} \\
&= \left( B_i^{(t)} - \eta_{i,t} \nabla_B \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) \right) \left( A_i^{(t)} - \eta_{i,t} \nabla_A \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) \right) - B_i^{(t)} A_i^{(t)} \\
&= \eta_{i,t}^2 \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) A_i^{(t)\top} B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) - \eta_{i,t} \nabla_B \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) A_i^{(t)} \\
&\quad - \eta_{i,t} B_i^{(t)} \nabla_A \mathcal{L}_i(W_i^{(t)}, \xi_{i,t})
\end{aligned} \tag{47}$$

Where the third equation is from  $\nabla_B \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) = \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) A_i^{(t)\top}$  and  $\nabla_A \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) = B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t})$  since  $W_i^{(t)} = W_0 + B_i^{(t)} A_i^{(t)}$ . Then we have:

$$\begin{aligned}
\mathbb{E} \left\langle U_i^{(t)} - W_i^{(t)}, \nabla_W \mathcal{L}_i(W_i^{(t)}) \right\rangle_F &= \eta_{i,t}^2 \mathbb{E} \left\langle \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) A_i^{(t)\top} B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}), \nabla_W \mathcal{L}_i(W_i^{(t)}) \right\rangle_F \\
&\quad - \eta_{i,t} \mathbb{E} \left\langle \nabla_B \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) A_i^{(t)}, \nabla_W \mathcal{L}_i(W_i^{(t)}) \right\rangle_F \\
&\quad - \eta_{i,t} \mathbb{E} \left\langle B_i^{(t)} \nabla_A \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}), \nabla_W \mathcal{L}_i(W_i^{(t)}) \right\rangle_F \\
&= \eta_{i,t}^2 \left\langle \nabla_W \mathcal{L}_i(W_i^{(t)}) A_i^{(t)\top} B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}), \nabla_W \mathcal{L}_i(W_i^{(t)}) \right\rangle_F \\
&\quad - \eta_{i,t} \left\langle \nabla_B \mathcal{L}_i(W_i^{(t)}), \nabla_W \mathcal{L}_i(W_i^{(t)}) A_i^{(t)\top} \right\rangle_F \\
&\quad - \eta_{i,t} \left\langle \nabla_A \mathcal{L}_i(W_i^{(t)}), B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}) \right\rangle_F
\end{aligned} \tag{48}$$

By Assumption 2 and 3, we have:

$$\begin{aligned}
&\left\langle \nabla_W \mathcal{L}_i(W_i^{(t)}) A_i^{(t)\top} B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}), \nabla_W \mathcal{L}_i(W_i^{(t)}) \right\rangle_F \\
&\leq \|\nabla_W \mathcal{L}_i(W_i^{(t)}) A_i^{(t)\top} B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)})\|_F \cdot \|\nabla_W \mathcal{L}_i(W_i^{(t)})\|_F \\
&\leq C_A C_B G^3
\end{aligned} \tag{49}$$

By  $\nabla_B \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) = \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) A_i^{(t)\top}$ ,  $\nabla_A \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) = B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t})$ , we have:

$$\left\langle \nabla_B \mathcal{L}_i(W_i^{(t)}), \nabla_W \mathcal{L}_i(W_i^{(t)}) A_i^{(t)\top} \right\rangle_F = \|\nabla_B \mathcal{L}_i(W_i^{(t)})\|_F^2 \tag{50}$$

and

$$\left\langle \nabla_A \mathcal{L}_i(W_i^{(t)}), B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}) \right\rangle_F = \|\nabla_A \mathcal{L}_i(W_i^{(t)})\|_F^2 \tag{51}$$

Finally, we can get the result by Equation (48), (49), (50), (51):

$$\mathbb{E} \left[ \left\langle U_i^{(t)} - W_i^{(t)}, \nabla_W \mathcal{L}_i(W_i^{(t)}) \right\rangle_F \right] \leq \eta_{i,t}^2 C_A C_B G^3 - \eta_{i,t} (\mathbb{E} [\|\nabla_B \mathcal{L}_i(W_i^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A \mathcal{L}_i(W_i^{(t)})\|_F^2]) \tag{52}$$

□

## A.6.2 PROOF OF LEMMA 2

*Proof.* By Equation (47), Assumption 2 and 3, we can get:

$$\begin{aligned}
\mathbb{E} [\|U_i^{(t)} - W_i^{(t)}\|_F^2] &= \mathbb{E} [\eta_{i,t}^2 \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) A_i^{(t)\top} B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) A_i^{(t)\top} \\
&\quad - \eta_{i,t} \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) A_i^{(t)\top} A_i^{(t)} - \eta_{i,t} B_i^{(t)} B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) \|_F^2] \\
&\leq 3 \mathbb{E} [\|\eta_{i,t}^2 \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) A_i^{(t)\top} B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t})\|_F^2] \\
&\quad + 3 \mathbb{E} [\|\eta_{i,t} \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) A_i^{(t)\top} A_i^{(t)}\|_F^2] + 3 \mathbb{E} [\|\eta_{i,t} B_i^{(t)} B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t})\|_F^2] \\
&\leq 3 \eta_{i,t}^4 C_A^2 C_B^2 G^4 + 3 \eta_{i,t}^2 C_A^4 G^2 + 3 \eta_{i,t}^2 C_B^4 G^2
\end{aligned} \tag{53}$$

□

## A.7 PROOF OF COROLLARY 1

*Proof.* First, we show that if the Aggregation-Broadcast Operators (ABO)  $\mathcal{P}$  and  $\mathcal{Q}$ , can satisfy the convergence condition in Definition 2 if:

$$\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) = \frac{1}{m} \sum_{j=1}^m B_j^{(t+1)} A_j^{(t+1)} \quad (54)$$

It is know by Eq. (54) that:

$$\begin{aligned} & \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} A_i^{(t+1)}\|_F^2 \\ &= \frac{1}{m} \sum_{i=1}^m \left\| \frac{1}{m} \sum_{j=1}^m (B_j^{(t+1)} A_j^{(t+1)} - B_i^{(t+1)} A_i^{(t+1)}) \right\|_F^2 \\ &\leq \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m \|B_j^{(t+1)} A_j^{(t+1)} - B_i^{(t+1)} A_i^{(t+1)}\|_F^2 \\ &= \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m \|B_j^{(t+1)} (A_j^{(t+1)} - A_i^{(t+1)}) + (B_j^{(t+1)} - B_i^{(t+1)}) A_i^{(t+1)}\|_F^2 \\ &\leq \frac{2}{m^2} \sum_{i=1}^m \sum_{j=1}^m (\|B_j^{(t+1)} (A_j^{(t+1)} - A_i^{(t+1)})\|_F^2 + \|B_j^{(t+1)} - B_i^{(t+1)}\|_F^2 \|A_i^{(t+1)}\|_F^2) \\ &\leq \frac{2}{m^2} \sum_{i=1}^m \sum_{j=1}^m (C_B^2 \cdot \|A_j^{(t+1)} - A_i^{(t+1)}\|_F^2 + C_A^2 \cdot \|B_j^{(t+1)} - B_i^{(t+1)}\|_F^2) \end{aligned} \quad (55)$$

Next, we consider about  $\|A_j^{(t+1)} - A_i^{(t+1)}\|_F^2$  and  $\|B_j^{(t+1)} - B_i^{(t+1)}\|_F^2$ . If  $t+1 \in \mathcal{I}_E$ , then there comes to a communication round, therefore  $A_j^{(t+1)} = A_i^{(t+1)}$  and  $B_j^{(t+1)} = B_i^{(t+1)}$  for  $1 \leq i \leq m$ . On the other hand, if  $t+1 \notin \mathcal{I}_E$ , we suppose that  $nE < t+1 < (n+1)E$  for an non-negative integer  $n$ . Then we have:

$$\begin{aligned} A_j^{(t+1)} &= A_j^{(nE)} - \sum_{t_0=nE}^t \eta \nabla_A L_j(W_j^{(t_0)}, \xi_{j,t_0}) \\ &= A_j^{(nE)} - \sum_{t_0=nE}^t \eta B_j^{(t_0)\top} \nabla_W L_j(W_j^{(t_0)}, \xi_{j,t_0}) \end{aligned} \quad (56)$$

Where  $\xi_{i,t}$  is the  $i$ -th client's local data uniformly at random at the training step  $t$ . The second equation is from  $\nabla_A \mathcal{L}_i(W_i^{(t)}, \xi_{i,t}) = B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)}, \xi_{i,t})$ . Then by Eq. (56) and Assumption 3

we have:

$$\begin{aligned}
\|A_j^{(t+1)} - A_i^{(t+1)}\|_F^2 &= \|(A_j^{(nE)} - \sum_{t_0=nE}^t \eta B_j^{(t_0)\top} \nabla_W L_j(W_j^{(t_0)}, \xi_{j,t_0})) \\
&\quad - (A_i^{(nE)} - \sum_{t_0=nE}^t \eta B_i^{(t_0)\top} \nabla_W L_i(W_i^{(t_0)}, \xi_{i,t_0}))\|_F^2 \\
&= \|\eta \sum_{t_0=nE}^t (B_j^{(t_0)\top} \nabla_W L_j(W_j^{(t_0)}, \xi_{j,t_0}) - B_i^{(t_0)\top} \nabla_W L_i(W_i^{(t_0)}, \xi_{i,t_0}))\|_F^2 \\
&\leq \eta^2(t - nE + 1) \sum_{t_0=nE}^t \|B_j^{(t_0)\top} \nabla_W L_j(W_j^{(t_0)}, \xi_{j,t_0}) - B_i^{(t_0)\top} \nabla_W L_i(W_i^{(t_0)}, \xi_{i,t_0})\|_F^2 \\
&\leq \eta^2(t - nE + 1) \sum_{t_0=nE}^t (2\|B_j^{(t_0)\top} \nabla_W L_j(W_j^{(t_0)}, \xi_{j,t_0})\|_F^2 + 2\|B_i^{(t_0)\top} \nabla_W L_i(W_i^{(t_0)}, \xi_{i,t_0})\|_F^2) \\
&\leq \eta^2(t - nE + 1) \sum_{t_0=nE}^t (2C_B^2 \|L_j(W_j^{(t_0)}, \xi_{j,t_0})\|_F^2 + 2C_B^2 \|\nabla_W L_i(W_i^{(t_0)}, \xi_{i,t_0})\|_F^2)
\end{aligned} \tag{57}$$

Take the expectation on Eq. (57) and by Assumption 2 we can get:

$$\begin{aligned}
\mathbb{E} [\|A_j^{(t+1)} - A_i^{(t+1)}\|_F^2] &\leq \eta^2(t - nE + 1) \sum_{t_0=nE}^t (2C_B^2 G^2 + 2C_B^2 G^2) \\
&= 4\eta^2(t - nE + 1)^2 C_B^2 G^2 \\
&\leq 4E^2 C_B^2 G^2 \eta^2
\end{aligned} \tag{58}$$

By the same reason, we have:

$$\mathbb{E} [\|B_j^{(t+1)} - B_i^{(t+1)}\|_F^2] \leq 4E^2 C_A^2 G^2 \eta^2 \tag{59}$$

Plugging Eq. (58) and Eq. (59) into Eq. (55) we can obtain:

$$\begin{aligned}
&\mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} A_i^{(t+1)}\|_F^2 \right] \\
&\leq \frac{2}{m^2} \sum_{i=1}^m \sum_{j=1}^m (C_B^2 \cdot 4E^2 C_B^2 G^2 \eta^2 + C_A^2 \cdot 4E^2 C_A^2 G^2 \eta^2) \\
&= 8E^2 G^2 (C_A^4 + C_B^4) \eta^2
\end{aligned} \tag{60}$$

Let  $R^2 = 8E^2 G^2 (C_A^4 + C_B^4)$  we can proof that  $\mathcal{P}$  and  $\mathcal{Q}$  satisfy the convergence Condition in Definition 2 if Eq 54 holds. Next, we proof that  $\mathcal{P}$  and  $\mathcal{Q}$  can achieve the optimal convergence rate if Eq 54 holds. Let  $f(X) = \frac{1}{m} \sum_{i=1}^m \|X - X_i\|_F^2$ , then we have:

$$\begin{aligned}
f(X) &= \frac{1}{m} \sum_{i=1}^m \text{tr}((X - X_i)^\top (X - X_i)) \\
&= \text{tr}(X^\top X) - 2 \cdot \frac{1}{m} \sum_{i=1}^m \text{tr}(X_i^\top X) + \frac{1}{m} \text{tr}(X_i^\top X_i)
\end{aligned} \tag{61}$$

Therefore:

$$\nabla f(X) = 2X - \frac{2}{m} \sum_{i=1}^m X_i \tag{62}$$

It means that:

$$\operatorname{argmin}_X f(X) = \operatorname{argmin}_X \sum_{i=1}^m \|X - X_i\|_F^2 = \frac{1}{m} \sum_{i=1}^m X_i \quad (63)$$

Finally, let  $X = \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)})$  and  $X_i = B_i^{(t+1)} A_i^{(t+1)}$  we can get the proof.  $\square$

#### A.8 PROOF OF THEOREM 4

*Proof.* First, we show that the ABO  $\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)})$  and  $\mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)})$  satisfy the Weak Convergence Condition defined in 2.

$$\begin{aligned} & \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} A_i^{(t+1)}\|_F^2 \\ & \leq \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) (\mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - A_i^{(t+1)}) + (\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)}) A_i^{(t+1)}\|_F^2 \\ & \leq 2 \cdot \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)})\|_F^2 \cdot \|\mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - A_i^{(t+1)}\|_F^2 \\ & \quad + 2 \cdot \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)}\|_F^2 \cdot \|A_i^{(t+1)}\|_F^2 \end{aligned} \quad (64)$$

By the Eq. (18) and Assumption 3, we can obtain the estimate of  $\|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)})\|_F^2$ :

$$\begin{aligned} \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)})\|_F^2 &= \left\| \frac{1}{m} \sum_{i=1}^m (\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)}) + \frac{1}{m} \sum_{i=1}^m B_i^{(t+1)} \right\|_F^2 \\ &\leq 2 \left\| \frac{1}{m} \sum_{i=1}^m (\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)}) \right\|_F^2 + 2 \left\| \frac{1}{m} \sum_{i=1}^m B_i^{(t+1)} \right\|_F^2 \\ &\leq \frac{2}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)}\|_F^2 + \frac{2}{m} \sum_{i=1}^m \|B_i^{(t+1)}\|_F^2 \\ &\leq 2P^2\eta^2 + 2C_B^2 \end{aligned} \quad (65)$$

By combining Eq. (64), Eq. (65) and Assumption. 3 we can get:

$$\begin{aligned} & \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} A_i^{(t+1)}\|_F^2 \\ & \leq 2 \cdot (2P^2\eta^2 + 2C_B^2) \cdot \frac{1}{m} \sum_{i=1}^m \|\mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - A_i^{(t+1)}\|_F^2 + 2C_A^2 P^2 \eta^2 \\ & \leq 4(P^2\eta^2 + C_B^2) \cdot Q^2\eta^2 + 2C_A^2 P^2 \eta^2 \\ & = 4P^2 Q^2 \eta^4 + 4C_B^2 Q^2 \eta^2 + 2C_A^2 P^2 \eta^2 \end{aligned} \quad (66)$$

Similarly, by the symmetry of  $\mathcal{P}$  and  $\mathcal{Q}$ , we obtain:

$$\begin{aligned} & \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} A_i^{(t+1)}\|_F^2 \\ & \leq 4P^2 Q^2 \eta^4 + 2C_B^2 Q^2 \eta^2 + 4C_A^2 P^2 \eta^2 \end{aligned} \quad (67)$$

Take the average of Eq. (66) and Eq. (67) then yields:

$$\begin{aligned} & \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} A_i^{(t+1)}\|_F^2 \\ & \leq 4P^2 Q^2 \eta^4 + 3C_B^2 Q^2 \eta^2 + 3C_A^2 P^2 \eta^2 \\ & = R^2 \eta^2 \end{aligned} \quad (68)$$

Where  $R^2 = 4P^2Q^2\eta^2 + 3C_B^2Q^2 + 3C_A^2P^2$ . By Definition 2 we know that  $\mathcal{P}$  and  $\mathcal{Q}$  satisfy the Weak Convergence Condition, then by Eq. (45) we have:

$$\frac{1}{mT} \sum_{t=1}^T \sum_{i=1}^m (\mathbb{E} [\|\nabla_B \mathcal{L}_i(W_i^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A \mathcal{L}_i(W_i^{(t)})\|_F^2]) \leq \frac{D}{\eta T} + M\eta \quad (69)$$

It is worth noting that Eq. (69) provides an estimate of the convergence rate of the local models. However, the convergence of local models does not necessarily imply the convergence of the global model. Therefore, we next turn our attention to analyzing the convergence of the global loss  $\mathcal{L}(W^{(t)}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}_i(W^{(t)})$  with global model parameter  $W^{(t)} = W_0 + \mathcal{P}(A_{1 \leq j \leq m}^{(t)}, B_{1 \leq j \leq m}^{(t)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t)}, B_{1 \leq j \leq m}^{(t)})$  and  $\nabla_B \mathcal{L}_i(W_i^{(t)}) = \nabla_W \mathcal{L}_i(W_i^{(t)}) A_i^{(t)\top}$ ,  $\nabla_A \mathcal{L}_i(W_i^{(t)}) = B_i^{(t)\top} \nabla_W \mathcal{L}_i(W_i^{(t)})$ . We know that:

$$\begin{aligned} \|\nabla_B \mathcal{L}(W^{(t)})\|_F^2 &\leq \|\nabla_B \mathcal{L}(W^{(t)}) - \frac{1}{m} \sum_{i=1}^m \nabla_B \mathcal{L}_i(W_i^{(t)}) + \frac{1}{m} \sum_{i=1}^m \nabla_B \mathcal{L}_i(W_i^{(t)})\|_F^2 \\ &\leq 2\left\| \frac{1}{m} \sum_{i=1}^m \nabla_B \mathcal{L}_i(W^{(t)}) - \frac{1}{m} \sum_{i=1}^m \nabla_B \mathcal{L}_i(W_i^{(t)}) \right\|_F^2 + 2\left\| \frac{1}{m} \sum_{i=1}^m \nabla_B \mathcal{L}_i(W_i^{(t)}) \right\|_F^2 \\ &\leq \frac{2}{m} \sum_{i=1}^m \|\nabla_B \mathcal{L}_i(W^{(t)}) - \nabla_B \mathcal{L}_i(W_i^{(t)})\|_F^2 + \frac{2}{m} \sum_{i=1}^m \|\nabla_B \mathcal{L}_i(W_i^{(t)})\|_F^2 \end{aligned} \quad (70)$$

moreover, we can get:

$$\begin{aligned} &\frac{1}{m} \sum_{i=1}^m \|\nabla_B \mathcal{L}_i(W^{(t)}) - \nabla_B \mathcal{L}_i(W_i^{(t)})\|_F^2 \\ &= \frac{1}{m} \sum_{i=1}^m \|\nabla_W \mathcal{L}_i(W^{(t)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t)}, B_{1 \leq j \leq m}^{(t)})^\top - \nabla_W \mathcal{L}_i(W_i^{(t)}) A_i^{(t)\top}\|_F^2 \\ &= \frac{1}{m} \sum_{i=1}^m \|\nabla_W \mathcal{L}_i(W^{(t)}) (\mathcal{Q}(A_{1 \leq j \leq m}^{(t)}, B_{1 \leq j \leq m}^{(t)})^\top - A_i^{(t)\top}) + (\nabla_W \mathcal{L}_i(W^{(t)}) - \nabla_W \mathcal{L}_i(W_i^{(t)})) A_i^{(t)\top}\|_F^2 \\ &\leq \frac{2}{m} \sum_{i=1}^m \|\nabla_W \mathcal{L}_i(W^{(t)}) (\mathcal{Q}(A_{1 \leq j \leq m}^{(t)}, B_{1 \leq j \leq m}^{(t)})^\top - A_i^{(t)\top})\|_F^2 + \frac{2}{m} \sum_{i=1}^m \|(\nabla_W \mathcal{L}_i(W^{(t)}) - \nabla_W \mathcal{L}_i(W_i^{(t)})) A_i^{(t)\top}\|_F^2 \end{aligned} \quad (71)$$

by Assumption 2 and the Strong Convergence Condition, we can get:

$$\begin{aligned} &\frac{1}{m} \sum_{i=1}^m \|\nabla_W \mathcal{L}_i(W^{(t)}) (\mathcal{Q}(A_{1 \leq j \leq m}^{(t)}, B_{1 \leq j \leq m}^{(t)})^\top - A_i^{(t)\top})\|_F^2 \\ &\leq \frac{1}{m} \sum_{i=1}^m \|\nabla_W \mathcal{L}_i(W^{(t)})\|_F^2 \cdot \|\mathcal{Q}(A_{1 \leq j \leq m}^{(t)}, B_{1 \leq j \leq m}^{(t)})^\top - A_i^{(t)\top}\|_F^2 \\ &\leq G^2 \cdot \frac{1}{m} \sum_{i=1}^m \|\mathcal{Q}(A_{1 \leq j \leq m}^{(t)}, B_{1 \leq j \leq m}^{(t)})^\top - A_i^{(t)\top}\|_F^2 \\ &\leq G^2 Q^2 \eta^2 \end{aligned} \quad (72)$$

by Assumption 1, Assumption 3 and Eq. (69), we have:

$$\begin{aligned}
& \frac{1}{m} \sum_{i=1}^m \|(\nabla_W \mathcal{L}_i(W^{(t)}) - \nabla_W \mathcal{L}_i(W_i^{(t)})) A_i^{(t)\top}\|_F^2 \\
& \leq \frac{1}{m} \sum_{i=1}^m \|\nabla_W \mathcal{L}_i(W^{(t)}) - \nabla_W \mathcal{L}_i(W_i^{(t)})\|_F^2 \cdot \|A_i^{(t)\top}\|_F^2 \\
& \leq \frac{1}{m} \sum_{i=1}^m C_A^2 L^2 \|W^{(t)} - W_i^{(t)}\|_F^2 \\
& \leq C_A^2 L^2 \cdot \frac{1}{m} \sum_{i=1}^m \|\mathcal{P}(A_{1 \leq j \leq m}^{(t)}, B_{1 \leq j \leq m}^{(t)}) \mathcal{Q}(A_{1 \leq j \leq m}^{(t)}, B_{1 \leq j \leq m}^{(t)}) - B_i^{(t)} A_i^{(t)}\|_F^2 \\
& \leq C_A^2 L^2 R^2 \eta^2
\end{aligned} \tag{73}$$

where  $R^2 = 4P^2 Q^2 \eta^2 + 3C_B^2 Q^2 + 3C_A^2 P^2$ . Then by Eq. (70), Eq. (71), Eq. (72) and Eq. (73) we can get:

$$\|\nabla_B \mathcal{L}(W^{(t)})\|_F^2 \leq 4G^2 Q^2 \eta^2 + 4C_A^2 L^2 R^2 \eta^2 + \frac{2}{m} \sum_{i=1}^m \|\nabla_B \mathcal{L}_i(W_i^{(t)})\|_F^2 \tag{74}$$

for the same reason, we have:

$$\|\nabla_A \mathcal{L}(W^{(t)})\|_F^2 \leq 4G^2 P^2 \eta^2 + 4C_B^2 L^2 R^2 \eta^2 + \frac{2}{m} \sum_{i=1}^m \|\nabla_A \mathcal{L}_i(W_i^{(t)})\|_F^2 \tag{75}$$

summing Eq. (74) and Eq. (75) over  $t = 1$  to  $T$ , then we can get:

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T (\|\nabla_B \mathcal{L}(W^{(t)})\|_F^2 + \|\nabla_A \mathcal{L}(W^{(t)})\|_F^2) & \leq 4G^2(Q^2 + P^2)\eta^2 + 4(C_A^2 + C_B^2)L^2 R^2 \eta^2 \\
& + 2 \cdot \frac{1}{mT} \sum_{t=1}^T \sum_{i=1}^m (\|\nabla_B \mathcal{L}_i(W_i^{(t)})\|_F^2 + \|\nabla_A \mathcal{L}_i(W_i^{(t)})\|_F^2)
\end{aligned} \tag{76}$$

take the expectation on both side and by Eq. (69), we can get:

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T (\mathbb{E} [\|\nabla_B \mathcal{L}(W^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A \mathcal{L}(W^{(t)})\|_F^2]) \\
& \leq 4G^2(Q^2 + P^2)\eta^2 + 4(C_A^2 + C_B^2)L^2 R^2 \eta^2 + 2\left(\frac{D}{\eta T} + M\eta\right) \\
& \leq 2\left(\frac{D}{\eta T} + (M + N)\eta\right)
\end{aligned} \tag{77}$$

where we assume that  $4G^2(Q^2 + P^2)\eta^2 + 4(C_A^2 + C_B^2)L^2 R^2 \eta^2 \leq 2N\eta$ ,  $\frac{3}{2}L(\eta^2 C_A^2 C_B^2 G^4 + C_A^4 G^2 + C_B^2 G^2)\eta^2 + C_A C_B G^3 \eta^2 + \frac{L}{2}R^2 \eta^2 + \frac{1}{2}(R^2 + G^2)\eta \leq M\eta^2$  and  $\mathcal{L}_i(W_0) - \mathcal{L}_i(W_i^*) \leq D$ . let  $\eta = \sqrt{\frac{D}{(M+N)T}}$ , we can get the global model convergence rate as:

$$\frac{1}{T} \sum_{t=1}^T (\mathbb{E} [\|\nabla_B \mathcal{L}(W^{(t)})\|_F^2] + \mathbb{E} [\|\nabla_A \mathcal{L}(W^{(t)})\|_F^2]) \leq 4\sqrt{\frac{D(M+N)}{T}} \tag{78}$$

□

## A.9 PROOF OF COROLLARY 2

*Proof.* The proof of this corollary is similar to the proof of corollary 1 which can be seen in Appendix A.7. First, we show that the Aggregation-Broadcast Operators (ABO)  $\mathcal{P}$  and  $\mathcal{Q}$  can satisfy

the Sufficient Condition in Theorem 3 if

$$\mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) = \frac{1}{m} \sum_{i=1}^m B_i^{(t+1)} \quad (79)$$

$$\mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) = \frac{1}{m} \sum_{i=1}^m A_i^{(t+1)} \quad (80)$$

By Eq. (79) we know that:

$$\begin{aligned} \mathbb{E} \left[ \left\| \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} \right\|_F^2 \right] &\leq \mathbb{E} \left[ \left\| \frac{1}{m} \sum_{j=1}^m (B_j^{(t+1)} - B_i^{(t+1)}) \right\|_F^2 \right] \\ &\leq \frac{1}{m} \sum_{j=1}^m \mathbb{E} \left[ \left\| B_j^{(t+1)} - B_i^{(t+1)} \right\|_F^2 \right] \\ &\leq 4E^2 C_A^2 G^2 \eta^2 \end{aligned} \quad (81)$$

The last equation holds by Eq. (59). Then by Eq. (81) we can get:

$$\begin{aligned} \mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \left\| \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} \right\|_F^2 \right] &\leq \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[ \left\| \mathcal{P}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - B_i^{(t+1)} \right\|_F^2 \right] \\ &\leq 4E^2 C_A^2 G^2 \eta^2 \end{aligned} \quad (82)$$

Similarly, by Eq 58 we have

$$\mathbb{E} \left[ \left\| \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - A_i^{(t+1)} \right\|_F^2 \right] \leq 4E^2 C_B^2 G^2 \eta^2 \quad (83)$$

and

$$\begin{aligned} \mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \left\| \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - A_i^{(t+1)} \right\|_F^2 \right] &\leq \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[ \left\| \mathcal{Q}(A_{1 \leq j \leq m}^{(t+1)}, B_{1 \leq j \leq m}^{(t+1)}) - A_i^{(t+1)} \right\|_F^2 \right] \\ &\leq 4E^2 C_B^2 G^2 \eta^2 \end{aligned} \quad (84)$$

Let  $P^2 = 4E^2 C_A^2 G^2$  and  $Q^2 = 4E^2 C_B^2 G^2$  we can proof that  $\mathcal{P}$  and  $\mathcal{Q}$  can satisfy the Sufficient Condition in Theorem 3 under Eq. (79) and Eq. (80). Meanwhile, by the same step form Eq. (61) to Eq. (63), we can also proof that  $\mathcal{P}$  and  $\mathcal{Q}$  can achieve the optimal convergence rate of the global model.  $\square$

## A.10 EXPERIMENTS CONFIGURATIONS AND ADDITIONAL RESULTS

In this section, we provide the remaining experimental results on FMNIST, KMNIST, and QMNIST. For completeness, the model configuration is detailed as follows:

**Model configurations:** The base model is a classic Multi-Layer Perceptron (MLP) with three fully-connected layers. To enable LoRA, we replace each dense layer with a LoRA layer. Specifically, the architecture consists of two hidden LoRA layers with output size 200 and a final output LoRA layer with size 10. The full model architecture is shown in Table. 1.

**Optimizer and loss function configurations:** All experiments use the SGD optimizer with 0.01 learning rate and cross-entropy loss function.

**Rank settings:** In our experiments, the global model’s rank ratio  $\delta_g$  is set to  $\delta_g = 1$ , and rest clients hold the same  $\delta$ , which is different under different scenarios.

Table 1: LoRA-augmented model architecture (rank-scaled).

Layer	Output Shape	Activation	LoRA Rank (A,B)
Input	(784,)	None	–
LoRALayer	(200,)	ReLU	$(r_1 = \max(160 \cdot \delta, 1))$
LoRALayer	(200,)	ReLU	$(r_2 = \max(160 \cdot \delta, 1))$
LoRALayer	(10,)	Softmax	$(r_3 = \max(100 \cdot \delta, 1))$



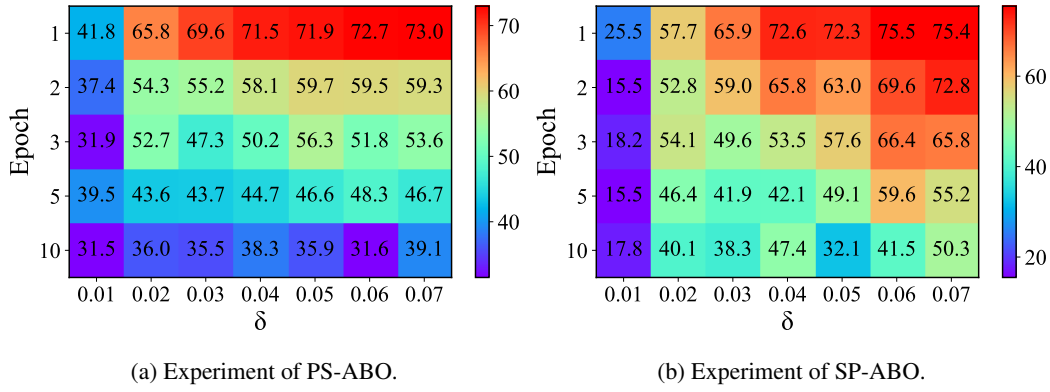


Figure 2: The highest accuracies (%) of PS-ABO and SP-ABO at different rank ratios and epochs when the total number of steps  $T = 300$  on FMNIST dataset.

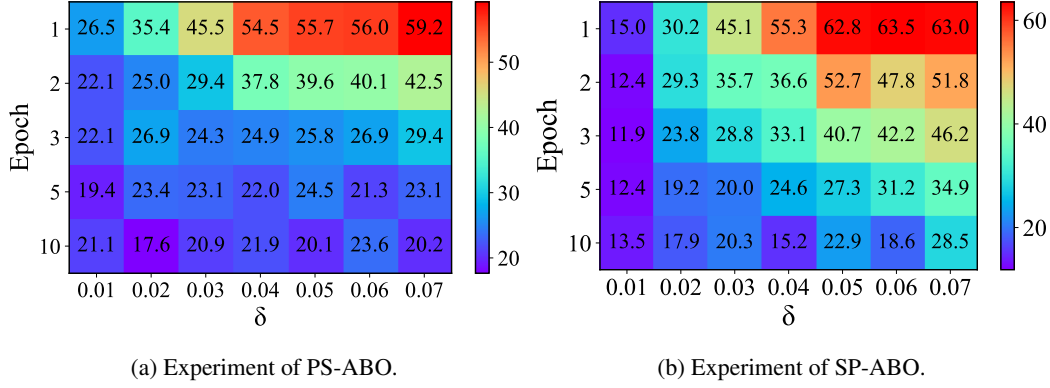


Figure 3: The highest accuracies (%) of PS-ABO and SP-ABO at different rank ratios and epochs when the total number of steps  $T = 300$  on KMNIST dataset.

Figure 2, 3, and 4 show the accuracy trend of PS-ABO and SP-ABO under different datasets (FMNIST, KMNIST, QMNIST), with varying rank ratios  $\delta$  and local epochs. From these results, we can clearly observe distinct sensitivities of the two operators to the experimental settings.

First, **SP-ABO is more sensitive to the rank ratio  $\delta$** . For example, on FMNIST with epoch = 1, the accuracy of SP-ABO decreases from 75.4% at  $\delta = 0.07$  to 25.5% at  $\delta = 0.01$ , a drop of about

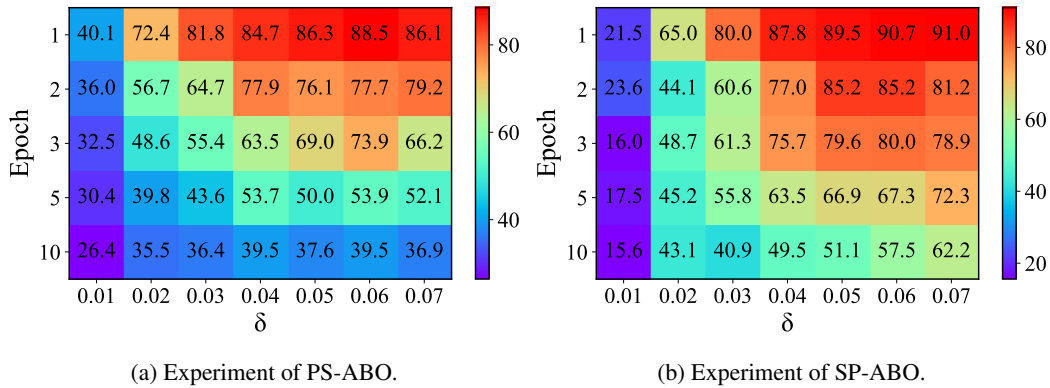


Figure 4: The highest accuracies (%) of PS-ABO and SP-ABO at different rank ratios and epochs when the total number of steps  $T = 300$  on QMNIST dataset.

50 percentage points. In contrast, PS-ABO under the same condition only decreases from 73.0% to 41.80% (about 31 points). Similar trends are observed on KMNIST (SP: 63.0% to 15.0%, PS: 59.2% to 26.5%) and QMNIST (SP: 91.0% to 21.5%, PS: 86.1% to 40.1%). Across all datasets, the performance gap between low and high rank settings is much larger for SP-ABO, indicating that its performance critically depends on the chosen LoRA rank.

Second, **PS-ABO is more sensitive to the number of local epochs**. On FMNIST with  $\delta = 0.07$ , PS-ABO decreases from 73% at epoch = 10 to 39.1% at epoch = 1 (a drop of 34 points), whereas SP-ABO decreases from 75% to 50% (a drop of 25 points). On QMNIST, the effect is clearer: PS-ABO falls from 86.1% to 36.9% when the epoch increases from 1 to 10, while SP-ABO drops only from 91% to 62.2%. Similar epoch-driven degradations are also observed on KMNIST. Overall, PS-ABO is less robust to larger epoch counts, showing a stronger dependence on synchronization frequency.

Overall, these empirical results demonstrate a clear contrast between the two operators: SP-ABO exhibits strong sensitivity to the choice of rank ratio, while PS-ABO exhibits strong sensitivity to the number of local epochs. This observation not only highlights the different convergence behaviors of the two methods but also aligns well with our theoretical analysis.