
Federated Sketching LoRA: On-Device Collaborative Fine-Tuning of Large Language Models

Wenzhi Fang¹ Dong-Jun Han² Liangqi Yuan¹ Seyyedali Hosseinalipour³ Christopher G. Brinton¹

Abstract

Fine-tuning large language models (LLMs) on devices is attracting increasing interest. Recent works have fused low-rank adaptation (LoRA) techniques with federated fine-tuning to mitigate challenges associated with device model sizes and data scarcity. Still, the heterogeneity of computational resources remains a critical bottleneck: while higher-rank modules generally enhance performance, varying device capabilities constrain LoRA’s feasible rank range. Existing approaches attempting to resolve this issue either lack analytical justification or impose additional computational overhead, leaving a wide gap for an efficient and theoretically-grounded solution. To address these challenges, we propose federated sketching LoRA (FSLoRA), which leverages a sketching mechanism to enable devices to selectively update submatrices of global LoRA modules maintained by the server. By adjusting the sketching ratios, which determine the ranks of the submatrices on the devices, FSLoRA flexibly adapts to device-specific communication and computational constraints. We provide a rigorous convergence analysis of FSLoRA that characterizes how the sketching ratios affect the convergence rate. Through comprehensive experiments on multiple datasets and LLM models, we demonstrate FSLoRA’s superior performance compared to various baselines.

1. Introduction

On-device LLMs have recently gained significant attention as a promising complement to cloud-based large language models (LLMs) (Fan et al., 2024). They align with the

¹Department of Electrical and Computer Engineering, Purdue University ²Department of Computer Science and Engineering, Yonsei University ³Department of Electrical Engineering, University at Buffalo-SUNY. Correspondence to: Wenzhi Fang <fang375@purdue.edu>.

typical paradigm of LLMs: they are initialized with a base model pre-trained on extensive datasets to capture linguistic patterns, semantics, and contextual nuances, and then fine-tuned on specific datasets to achieve better performance in specialized or domain-specific tasks. However, an LLM fine-tuned on a single device often achieves unsatisfactory performance due to the limited data available on each device. Fortunately, federated learning (McMahan et al., 2017; Chen et al., 2023) offers an effective solution here, enabling the model to be fine-tuned across a distributed group of clients within the same task domain, without any data sharing.

However, federated learning imposes significant computational and memory costs, as each device must fine-tune the LLM using its local dataset and send updates to the server for model aggregation. Recently, many parameter-efficient fine-tuning methods have been proposed (Lester et al., 2021; Li & Liang, 2021; Hu et al., 2021) to reduce the cost associated with model adaptation. Among them, low-rank adaptation (LoRA) (Hu et al., 2021) stands out as a particularly effective approach due to its flexibility. LoRA enables efficient fine-tuning by approximating weight updates $\Delta\mathbf{W}$ through a low-rank decomposition $\Delta\mathbf{W} = \mathbf{B}\mathbf{A}$, where matrices \mathbf{B} and \mathbf{A} contain significantly fewer trainable parameters than the original weight matrix. To support distributed on-device LLM, Zhang et al. (2024); Ye et al. (2024) incorporated LoRA with FedAvg (McMahan et al., 2017), significantly reducing the fine-tuning cost by cutting down the number of parameters that need to be synchronized across distributed devices.

Challenges. While integrating federated learning with LoRA reduces the number of trainable parameters through matrix decomposition, *communication costs still increase linearly with the rank of the decomposition*. This poses challenges when complex tasks demand higher-rank LoRA modules, particularly on resource-constrained mobile devices. Furthermore, the *heterogeneity in computational and communication capabilities across distributed devices makes a uniform rank inefficient*: a fixed rank r may be too large for some devices while being too small for more powerful ones, resulting in underutilized resources. Consequently, an approach that reduces communication overhead while adapting LoRA ranks to heterogeneous device capa-

bilities is highly desirable for collaborative fine-tuning of LLMs. Although some existing approaches have attempted to provide a solution here (Cho et al., 2024; Bai et al., 2024; Wang et al., 2024), they either lack theoretical justification or impose additional computational overhead, leaving a gap for an efficient and theoretically-grounded solution. As discussed in Section 2.2, a comprehensive approach that retains LoRA’s advantages while addressing heterogeneous on-device fine-tuning under tight resource constraints, with theoretical guarantees, has remained elusive.

1.1. Contributions

Motivated by these limitations, this work aims to develop a methodology for collaborative on-device LLM fine-tuning that (i) retains the flexibility of LoRA, (ii) provides theoretical convergence guarantees, and (iii) addresses the challenges posed by system heterogeneity and resource constraints across distributed devices. As depicted in Figure 1, our key idea is to introduce a sketching-based LoRA update to the fine-tuning process, which allows devices to selectively update a subset of columns and rows of the LoRA modules during each round, reducing the computation and communication consumption through the system. Additionally, our method customizes the fine-tuning process by adjusting the sparsity level of the sketching matrix, i.e., the size of the updated submatrices for each device in each iteration, as illustrated in Figure 1. As we will see, the impact of the introduced sketching mechanism on the overall optimization landscape requires careful consideration, posing additional challenges for the theoretical analysis which we investigate.

Overall, we make the following contributions:

- We propose federated sketching LoRA (FSLoRA), which leverages a sketching mechanism to enable devices to selectively update submatrices of global LoRA modules maintained by the server. By adjusting the sketching ratios, which determine the ranks of the submatrices on devices, FSLoRA effectively adapts to device-specific communication and computational constraints.
- We present a rigorous convergence analysis of FSLoRA under non-uniform submatrix update scenarios (i.e., heterogeneous LoRA configurations) across devices, revealing how the sketching ratios affect the convergence rate. In particular, we show how increasing the sketching ratios improves convergence theoretically but raises communication and computation costs, highlighting a delicate trade-off in selecting the sketching ratios.
- We conduct extensive experiments across multiple datasets and LLM models with diverse parameter settings, demonstrating FSLoRA’s superior performance compared to various baselines in terms of testing accuracy, training time, and resource utilization, validating the effectiveness

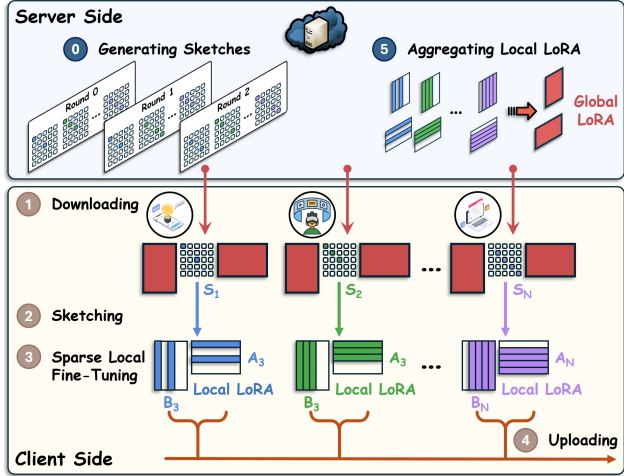


Figure 1: An illustration of our proposed methodology where the server maintains a pair of global LoRA modules while the devices adaptively update submatrices of the global LoRA modules through sketching during each round.

of the sketching mechanism.

1.2. Related Works

LoRA-based parameter-efficient fine-tuning. LoRA was first introduced in (Hu et al., 2021) as a parameter-efficient alternative to full model fine-tuning by utilizing low-rank matrix approximation. Subsequently, Kalajdzievski (2023) proposed rank-stabilized LoRA (rsLoRA), an approach that enhances LoRA’s performance in high-rank scenarios by modifying the scaling factor. Shuttleworth et al. (2024) demonstrated that with this novel scaling factor design, rsLoRA could approach the performance of full model fine-tuning as the rank of LoRA modules increases. Malinovskiy et al. (2024) investigated an alternating update of LoRA modules and provided a theoretical analysis. Han et al. (2024) introduced a sparse matrix in parallel with LoRA modules to improve the fitting capability of LoRA. Lialin et al. (2023); Xia et al. (2024) proposed ReLoRA and Chain of LoRA, which periodically merge learned LoRA modules into the full model to achieve a high-rank fine-tuning. However, the works mentioned above consider centralized scenarios, assuming that the data required for fine-tuning is available at the central server.

Collaborative Fine-tuning via Federated LoRA. Leveraging the efficiency of LoRA, federated LoRA is recently gaining recognition as a promising approach for collaborative fine-tuning of LLMs across distributed devices (Chen et al., 2023). Sun et al. (2024) examined the performance of federated LoRA with the incorporation of differential privacy. To address communication overhead, Kuo et al. (2024) proposed integrating communication compression with federated LoRA. Meanwhile, Bai et al. (2024); Cho

et al. (2024); Byun & Lee (2024); Wang et al. (2024); Koo et al. (2024) explored the challenges of resource heterogeneity among distributed devices and introduced heterogeneous LoRA as a solution. However, the approaches proposed in (Cho et al., 2024; Koo et al., 2024; Byun & Lee, 2024) lack a theoretical justification. FlexLoRA, introduced in (Bai et al., 2024), incurs additional computational overhead due to its reliance on singular value decomposition (SVD). Furthermore, the stack LoRA method proposed in (Bai et al., 2024; Wang et al., 2024) requires the devices to integrate the LoRA modules into the base model, thereby compromising the inherent flexibility of LoRA. Overall, there is still a lack of a systematic and theoretically grounded solution that can effectively tackle the challenges of heterogeneity in collaborative on-device LLM fine-tuning systems.

2. Problem Background

2.1. LoRA-based Federated LLM Fine-tuning

Following the decomposition of LoRA, the federated LoRA fine-tuning problem can be formulated as

$$\begin{aligned} \min_{\mathbf{B}, \mathbf{A}} f(\mathbf{B}, \mathbf{A}) &:= \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{B}, \mathbf{A}) \\ f_i(\mathbf{B}, \mathbf{A}) &:= \mathbb{E}_{\xi \sim \mathcal{D}_i} [\ell(\mathbf{W}_0 + \mathbf{B}\mathbf{A}, \xi)], \end{aligned} \quad (1)$$

where \mathbf{W}_0 denotes the frozen base model, $\mathbf{B} \in \mathbb{R}^{m \times r}$, $\mathbf{A} \in \mathbb{R}^{r \times n}$ are LoRA modules, N denotes the number of devices, ξ denotes a data sample, and \mathcal{D}_i is the local dataset on device i . ℓ , f_i , and f are the sample loss function, the local loss for device i , and the global loss, respectively. Problem (1) aligns with the conventional federated optimization formulation, which thus can be solved using the FedAvg algorithm. Based on the FedAvg framework, Zhang et al. (2024) developed the federated LoRA algorithm, which applies a uniform rank r across all devices, overlooking system heterogeneity. This one-size-fits-all approach leads to resource mismatches, where computationally constrained devices may struggle, while more powerful devices remain underutilized with a fixed rank.

2.2. Aren't the Existing Solutions Good Enough?

To address this issue, researchers have proposed heterogeneous federated LoRA approaches, where devices maintain non-uniform LoRA modules with varying ranks. They also introduce mechanisms to overcome the challenges of directly aggregating matrices with different dimensions. However, these methods often lack theoretical justification or introduce significant computational overhead, as outlined below.

HeteroLoRA (Cho et al., 2024) lets the server pad the updates from the devices with smaller ranks to match the size

of the largest rank during aggregation. During model dissemination, devices receive a truncated version of the global LoRA modules from the server. HeteroLoRA, while easy to implement, lacks a solid theoretical foundation. Additionally, its dependence on zero-padding diminishes optimization efficiency, potentially limiting overall performance.

FlexLoRA (Bai et al., 2024) requires the server to collect the individual LoRA matrices \mathbf{B}_i and \mathbf{A}_i from the devices and then computes their product $\mathbf{B}_i \mathbf{A}_i$. To support the initialization of non-uniform LoRA modules, the server applies truncated SVD to the averaged product $\frac{1}{N} \sum_{i=1}^N \mathbf{B}_i \mathbf{A}_i$. However, this approach introduces extra computational and memory overhead due to truncated SVD, and the associated error may limit the performance.

FedStackLoRA (Wang et al., 2024) introduces a stacking mechanism where the server concatenates LoRA modules from the devices. The concatenated matrices are then sent back to the devices, which compute their product and merge it into the base model before initializing new LoRA modules for the next fine-tuning round. However, this approach increases communication complexity linearly with the number of devices, imposes higher computation and memory demands on the devices, and compromises LoRA's flexibility to support multiple adapters for different tasks.

In summary, a theoretically-grounded solution that preserves LoRA's benefits while effectively addressing system heterogeneity and the constraints of resource-limited devices remains lacking.

3. Methodology

Driven by the limitations of existing methods, we propose a new federated LoRA reformulation. Building on this foundation, we develop FSLoRA, a resource-adaptive algorithm that preserves LoRA's flexibility while accounting for system heterogeneity and resource constraints.

3.1. Our Formulation

We propose a sketching-based LoRA formulation for collaborative LLM fine-tuning as follows:

$$\begin{aligned} \min_{\mathbf{B}, \mathbf{A}} f^{\mathcal{S}}(\mathbf{B}, \mathbf{A}) &:= \frac{1}{N} \sum_{i=1}^N f_i^{\mathcal{S}}(\mathbf{B}, \mathbf{A}) \\ f_i^{\mathcal{S}}(\mathbf{B}, \mathbf{A}) &:= \mathbb{E}_{\mathbf{S} \sim \mathcal{S}_i; \xi \sim \mathcal{D}_i} [\ell(\mathbf{W}_0 + \mathbf{B}\mathbf{S}\mathbf{A}, \xi)], \end{aligned} \quad (2)$$

where $\mathbf{B} \in \mathbb{R}^{m \times r}$, $\mathbf{A} \in \mathbb{R}^{r \times n}$ are LoRA modules, $f_i^{\mathcal{S}}$ is the local loss function at device i with sketching, and \mathbf{S} denotes a sketching matrix randomly sampled from the diagonal matrix set $\mathcal{S}_i = \mathcal{S}(r, k_i)$. The set $\mathcal{S}(r, k_i)$ comprises diagonal matrices of size $r \times r$ with exactly k_i non-zero entries. The formal definition of $\mathcal{S}(r, k)$ is provided below:

Definition 3.1 (Random- k sketching). A random- k diago-

nal matrix set is defined as:

$$\mathcal{S}(r, k) = \left\{ \mathbf{S} \mid \mathbf{S} = \frac{r}{k} \sum_{j \in \mathcal{I}} \mathbf{e}_j \mathbf{e}_j^\top, \mathcal{I} \subseteq \{1, \dots, r\}, |\mathcal{I}| = k \right\},$$

where $\mathbf{e}_1, \dots, \mathbf{e}_r \in \mathbb{R}^r$ are standard unit basis vectors and index set \mathcal{I} is a random subset of $[r] := \{1, 2, \dots, r\}$ sampled uniformly from all subsets of $[r]$ with cardinality k .

With \mathbf{S} being a matrix sampled from \mathcal{S}_i , we have

$$\mathbf{BSA} = \frac{r}{k_i} \sum_{j \in \mathcal{I}_i} \mathbf{B} \mathbf{e}_j \mathbf{e}_j^\top \mathbf{A},$$

where \mathcal{I}_i corresponds to the index set of non-zero entries of \mathbf{S} . $\mathbf{B} \mathbf{e}_j$ extracts the j -th column of \mathbf{B} while $\mathbf{e}_j^\top \mathbf{A}$ extracts the j -th row of \mathbf{A} . In other words, only k_i columns and rows in the LoRA modules \mathbf{B} and \mathbf{A} are activated by the sketching matrix in the loss $\ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)$ at device i . On the other hand, the sketching matrix \mathbf{S} satisfies $\mathbb{E}_{\mathbf{S} \sim \mathcal{S}_i}[\mathbf{S}] = \mathbf{I}_r$ where \mathbf{I}_r is a r -dimensional identity matrix. Based upon this property, $\mathbf{W}_0 + \mathbf{BSA}$ can be treated as an unbiased estimate of $\mathbf{W}_0 + \mathbf{BA}$.

Intuition. A larger rank allows LoRA modules to be more expressive, leading to better performance. However, resource-constrained devices cannot afford the computational or communication demands of large-rank modules. Our formulation (2) leverages the sketching matrix to balance the expressiveness of high-rank LoRA modules with the resource constraints of different devices. With the sketching mechanism introduced, the local gradients with respect to the LoRA modules on the devices will exhibit structured sparsity. By adjusting the sketching ratio k_i/r , we can tailor the sparsity of the gradient to match the capabilities of each device, ensuring affordable training while maintaining performance across heterogeneous systems, as elaborated in the following subsection. Overall, compared to the original objective in (1), our formulation offers a more resource-adaptive and flexible framework, tailored to address the diverse capabilities of heterogeneous devices.

3.2. Sparsity in the Gradients

In this subsection, we analyze the gradient structure of LoRA modules and highlight the gradients' sparsity properties under a given sketching matrix. To begin, we present the gradient expressions for the LoRA modules \mathbf{B} and \mathbf{A} in the following lemma. The proof is provided in Appendix D.2.

Lemma 3.2 (Gradient Formulation). *For a given sketching matrix \mathbf{S} , the gradients of $\ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)$ with respect to \mathbf{B} and \mathbf{A} take the following form*

$$\begin{aligned} \nabla_{\mathbf{B}} \ell(\mathbf{W}_0 + \mathbf{BSA}, \xi) &= \nabla \ell(\mathbf{W}_0 + \mathbf{BSA}, \xi) \mathbf{A}^\top \mathbf{S}^\top \\ \nabla_{\mathbf{A}} \ell(\mathbf{W}_0 + \mathbf{BSA}, \xi) &= \mathbf{S}^\top \mathbf{B}^\top \nabla \ell(\mathbf{W}_0 + \mathbf{BSA}, \xi), \end{aligned} \quad (3)$$

where $\nabla_{\mathbf{B}} \ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)$, $\nabla_{\mathbf{A}} \ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)$, and $\nabla \ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)$ represent the gradients of $\ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)$ with respect to \mathbf{B} , \mathbf{A} , and $\mathbf{W}_0 + \mathbf{BSA}$, respectively.

In particular, a random- k diagonal sketching matrix selectively samples k rows or columns of a matrix through left product or right product, respectively. With \mathbf{S} being a random- k diagonal matrix, the gradients of $\ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)$ with respect to LoRA modules \mathbf{B} and \mathbf{A} , as shown in (3), naturally become structurally sparse matrices. This sparsity reduces the computational and memory overhead during training, allowing for faster gradient computation and parameter updates. Additionally, sparse training enables better scalability across distributed devices by reducing communication costs, as only the non-zero elements need to be transmitted. The sparsity level of these gradients at each device is determined by the corresponding sketching matrix set \mathcal{S}_i .

Remark 3.3 (Sparsity Level Control). A key advantage of our formulation is its flexible control over the sparsity level of local gradients, achieved by configuring the parameter k_i of the sketching matrix set $\mathcal{S}_i = \mathcal{S}(r, k_i)$. This mechanism allows each device to tailor its local updates according to its communication and computation resource constraints, ensuring efficient and scalable fine-tuning in heterogeneous federated systems. Lowering k_i helps resource-constrained devices reduce computation and communication overhead, while more capable devices can increase k_i to conduct more informative local updates.

3.3. Algorithm

Based on the formulation in (2), we propose a resource-adaptive algorithm termed FSLoRA for collaborative on-device fine-tuning. FSLoRA allows each device to update submatrices of the original modules \mathbf{B} and \mathbf{A} in each round. The server maintains a pair of global LoRA modules \mathbf{B} and \mathbf{A} and periodically updates them by aggregating sparse local updates received from distributed devices. Specifically, the procedure of FSLoRA at each round t is detailed below.

- The server begins by generating sketching matrices $\{\mathbf{S}_i^t \sim \mathcal{S}_i\}_{i=1}^N$ for all devices, where \mathcal{S}_i represents the set of possible sketching matrices for device i . These sketches are then sent to the corresponding devices. Additionally, the server broadcasts the current global LoRA modules $[\mathbf{B}^t; \mathbf{A}^t]$ to all devices.
- Devices perform local fine-tuning using sketch \mathbf{S}_i^t . Specifically, guided by sketching matrix \mathbf{S}_i^t , the update at device i during the h -th iteration of the t -th round is given by:

$$\begin{bmatrix} \mathbf{B}_i^{t,h+1} \\ \mathbf{A}_i^{t,h+1} \end{bmatrix} = \begin{bmatrix} \mathbf{B}_i^{t,h} \\ \mathbf{A}_i^{t,h} \end{bmatrix} - \gamma \begin{bmatrix} \Delta \mathbf{B}_i^{t,h} (\mathbf{S}_i^t)^\top \\ (\mathbf{S}_i^t)^\top \Delta \mathbf{A}_i^{t,h} \end{bmatrix}, \quad (4)$$

Algorithm 1 Federated Sketching LoRA (FSLoRA)

Require: Base model \mathbf{W}_0 , initial LoRA modules \mathbf{B}_0 and \mathbf{A}_0 , learning rate γ , and sketching matrix set $\{\mathcal{S}_i\}_{i=1}^N$

- 1: **for** $t = 0, 1, \dots, T - 1$ **do**
- 2: Server generates sketching matrices $\{\mathbf{S}_i^t \sim \mathcal{S}_i\}_{i=1}^N$ and sends them back to the devices
- 3: Server broadcasts the current global LoRA modules to the devices
- 4: **for** $h = 0, 1, \dots, H - 1$ **do**
- 5: Devices update local LoRA modules via (4)
- 6: **end for**
- 7: Devices upload non-zero columns of $(\mathbf{B}_i^{t,H} - \mathbf{B}_i^{t,0})$ and non-zero rows of $(\mathbf{A}_i^{t,H} - \mathbf{A}_i^{t,0})$ to the server
- 8: Server updates the global LoRA modules via (5)
- 9: **end for**

where γ denotes the learning rate and $[\Delta \mathbf{B}_i^{t,h}; \Delta \mathbf{A}_i^{t,h}]$ is a shorthand representation for:

$$\begin{bmatrix} \Delta \mathbf{B}_i^{t,h} \\ \Delta \mathbf{A}_i^{t,h} \end{bmatrix} = \begin{bmatrix} \nabla \ell(\mathbf{W}_0 + \mathbf{B}_i^{t,h} \mathbf{S}_i^t \mathbf{A}_i^{t,h}; \xi_i^{t,h}) (\mathbf{A}_i^{t,h})^\top \\ (\mathbf{B}_i^{t,h})^\top \nabla \ell(\mathbf{W}_0 + \mathbf{B}_i^{t,h} \mathbf{S}_i^t \mathbf{A}_i^{t,h}; \xi_i^{t,h}) \end{bmatrix}.$$

The update direction in (4) corresponds to the negative stochastic gradient of $\ell(\mathbf{W}_0 + \mathbf{B} \mathbf{S} \mathbf{A}, \xi)$ with respect to $[\mathbf{B}; \mathbf{A}]$ for a given sketch \mathbf{S}_i^t , as established in Lemma 3.2. The total update for device i during one round of training, consisting of H local steps, can be expressed as follows:

$$\begin{bmatrix} \mathbf{B}_i^{t,H} - \mathbf{B}_i^{t,0} \\ \mathbf{A}_i^{t,H} - \mathbf{A}_i^{t,0} \end{bmatrix} = \begin{bmatrix} \gamma \left(\sum_{h=0}^{H-1} \Delta \mathbf{B}_i^{t,h} \right) (\mathbf{S}_i^t)^\top \\ \gamma (\mathbf{S}_i^t)^\top \left(\sum_{h=0}^{H-1} \Delta \mathbf{A}_i^{t,h} \right) \end{bmatrix}.$$

From the above equation, we see that only the columns of \mathbf{B} and the rows of \mathbf{A} corresponding to the nonzero entries of \mathbf{S}_i^t are updated during the t -th round at device i . In essence, \mathbf{S}_i^t selectively activates specific columns of \mathbf{B} and rows of \mathbf{A} for each round. Afterward, devices transmit these nonzero columns and rows of the sparse model updates to the server.

- Using the sketch information, the server reconstructs the corresponding sparse matrices from the received updates and aggregates them to update the global model:

$$\begin{bmatrix} \mathbf{B}^{t+1} \\ \mathbf{A}^{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{B}^t \\ \mathbf{A}^t \end{bmatrix} + \frac{1}{N} \sum_{i=1}^N \begin{bmatrix} \mathbf{B}_i^{t,H} - \mathbf{B}_i^{t,0} \\ \mathbf{A}_i^{t,H} - \mathbf{A}_i^{t,0} \end{bmatrix}. \quad (5)$$

The above procedure is repeated for $t = 0, 1, \dots, T - 1$ across T global rounds. Algorithm 1 summarizes the overall process of FSLoRA.

3.4. Comparison with Communication Compression

Although both the sketching approach in FSLoRA and communication compression (Kuo et al., 2024) reduce communication overhead, the sketching approach fundamentally

differs from traditional compression techniques. Notably, these two methods are orthogonal and can be combined to achieve greater efficiency. Specifically, the compression can be applied to the transmission of non-zero columns of \mathbf{B} and the non-zero rows of \mathbf{A} in FSLoRA to further enhance communication efficiency. We demonstrate the effectiveness of this combination in Appendix C.2. Additionally, the compression focuses solely on reducing the transmission load, leaving the gradient computation and model updates unchanged from the vanilla federated LoRA, FSLoRA goes beyond communication savings by also reducing gradient computation and model update overhead through sparse training.

4. Analysis

In this section, we analyze the convergence of the proposed FSLoRA algorithm. We will show that the iterate sequence generated by FSLoRA algorithm converges to the stationary point of function (2). In our analysis, we will use the following notations.

Notations: We define $\tilde{\ell}(\mathbf{B}, \mathbf{A}, \xi; \mathbf{S}) = \ell(\mathbf{W}_0 + \mathbf{B} \mathbf{S} \mathbf{A}, \xi)$ and $\tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) = \mathbb{E}_{\xi \sim \mathcal{D}_i} [\ell(\mathbf{W}_0 + \mathbf{B} \mathbf{S} \mathbf{A}, \xi)]$ for a given \mathbf{S} and $f_i^{\mathcal{S}}(\mathbf{B}, \mathbf{A}) = \mathbb{E}_{\mathbf{S} \sim \mathcal{S}_i} [\tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S})]$. For simplicity, we denote $\mathbf{X} = [\mathbf{B}; \mathbf{A}]$ and rewrite $f(\mathbf{B}, \mathbf{A})$, $f_i(\mathbf{B}, \mathbf{A})$, $f^{\mathcal{S}}(\mathbf{B}, \mathbf{A})$, $f_i^{\mathcal{S}}(\mathbf{B}, \mathbf{A})$, $\tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S})$, and $\tilde{\ell}(\mathbf{B}, \mathbf{A}, \xi; \mathbf{S})$ as $f(\mathbf{X})$, $f_i(\mathbf{X})$, $f^{\mathcal{S}}(\mathbf{X})$, $f_i^{\mathcal{S}}(\mathbf{X})$, $\tilde{f}_i(\mathbf{X}; \mathbf{S})$, and $\tilde{\ell}(\mathbf{X}, \xi; \mathbf{S})$ respectively. Additionally, we use $\|\cdot\|$ to denote the Frobenius norm in our analysis.

We conduct analysis based on the following assumptions.

Assumption 4.1. $f_i(\mathbf{X})$ is differentiable and L -smooth, i.e., there exists a positive constant L such that $\forall \mathbf{X}, \mathbf{Y}$,

$$\|\nabla f_i(\mathbf{X}) - \nabla f_i(\mathbf{Y})\| \leq L \|\mathbf{X} - \mathbf{Y}\|, \forall i.$$

Assumption 4.2. The variance of the gradient $\nabla_{\mathbf{X}} \tilde{f}_i(\mathbf{X}; \mathbf{S})$ from the sketching matrix $\mathbf{S} \sim \mathcal{S}_i$ can be bounded as

$$\mathbb{E}_{\mathbf{S} \sim \mathcal{S}_i} \|\nabla_{\mathbf{X}} \tilde{f}_i(\mathbf{X}; \mathbf{S}) - \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X})\|^2 \leq \sigma_s^2, \forall i,$$

where $f_i^{\mathcal{S}}(\mathbf{X}) = \mathbb{E}_{\mathbf{S} \sim \mathcal{S}_i} [\tilde{f}_i(\mathbf{X}; \mathbf{S})]$. In addition, for a given \mathbf{S} , the variance of the stochastic gradient $\nabla_{\mathbf{X}} \tilde{\ell}(\mathbf{X}, \xi; \mathbf{S})$ due to data sampling $\xi \sim \mathcal{D}_i$ can be bounded as

$$\mathbb{E}_{\xi \sim \mathcal{D}_i} \|\nabla_{\mathbf{X}} \tilde{\ell}(\mathbf{X}, \xi; \mathbf{S}) - \nabla_{\mathbf{X}} \tilde{f}_i(\mathbf{X}; \mathbf{S})\|^2 \leq \sigma_g^2, \forall i,$$

where $\tilde{f}_i(\mathbf{X}; \mathbf{S}) = \mathbb{E}_{\xi \sim \mathcal{D}_i} [\tilde{\ell}(\mathbf{X}, \xi; \mathbf{S})]$.

Assumption 4.3. The gradient dissimilarity between the global loss $f^{\mathcal{S}}(\mathbf{X})$ and each local loss $f_i^{\mathcal{S}}(\mathbf{X})$ satisfies

$$\|\nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}) - \nabla_{\mathbf{X}} f^{\mathcal{S}}(\mathbf{X})\|^2 \leq c_h \|\nabla_{\mathbf{X}} f^{\mathcal{S}}(\mathbf{X})\|^2 + \delta_h^2, \forall i,$$

where $c_h \geq 0$ and $f^{\mathcal{S}}(\mathbf{X}) = \frac{1}{N} \sum_{i=1}^N f_i^{\mathcal{S}}(\mathbf{X})$.

Assumptions 4.1 and 4.2 are standard in stochastic optimization (Demidovich et al., 2023; Fang et al., 2024), while Assumption 4.3 is commonly used in distributed optimization (Fang et al., 2022; Yi et al., 2022) to characterize data heterogeneity. Building on these assumptions, we analyze the convergence behavior of FSLoRA. Our main result is summarized in the following theorem.

Theorem 4.4. *Suppose that Assumptions 4.1-4.3 hold and the learning rate satisfies $\gamma \leq \min\left\{\frac{1}{6\sqrt{(c_h+1)\bar{L}LH}}, \frac{1}{NH^2L}\right\}$.*

Then the iterate sequence $\{\mathbf{X}^t\}_{t=0}^{T-1}$ generated by FSLoRA satisfies

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2 &\leq 4 \frac{f^S(\mathbf{X}^0) - f^*}{\gamma TH} \\ &+ 6 \frac{\gamma \bar{L}}{N} (\sigma_g^2 + \sigma_s^2) + 12 \frac{\gamma \tilde{L}}{N} (\sigma_g^2 + \sigma_s^2 + 3\sigma_h^2), \end{aligned} \quad (6)$$

where $\bar{L} = \left(\frac{1}{N} \sum_{i=1}^N \frac{r}{k_i}\right) L$, $\tilde{L} = \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2}\right) L$, and f^* denotes the lower bound of $f^S(\mathbf{X})$.

A novel step in the proof of Theorem 4.4 is the characterization of how the introduced sketching mechanism impacts the optimization landscape. This analysis delves into the way the sketching operation modifies the smoothness properties of the objective function, particularly how it introduces scaled smoothness constants, $\frac{r}{k_i}L$ and $\frac{r^2}{k_i^2}L$, that impacts the overall optimization. Further details are presented in Appendix D.3.

Based on the results in Theorem 4.4, we have the following corollary by applying an appropriate learning rate γ to Algorithm 1.

Corollary 4.5. *Under the same assumptions of Theorem 4.4, let $\mathcal{F}_0 = f^S(\mathbf{X}^0) - f^*$ and the learning rate $\gamma = \sqrt{N}/\sqrt{TH\tilde{L}}$ which satisfies the condition outlined in Theorem 4.4 when T is large enough. Then the iterate sequence $\{\mathbf{X}^t\}_{t=0}^{T-1}$ generated by FSLoRA satisfies*

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2 &\leq \mathcal{O}\left(\frac{\mathcal{F}_0 \sqrt{\tilde{L}}}{\sqrt{THN}}\right) \\ &+ \mathcal{O}\left(\frac{(\sigma_g^2 + \sigma_s^2 + \sigma_h^2) \sqrt{\tilde{L}}}{\sqrt{THN}}\right). \end{aligned} \quad (7)$$

Discussion: The results obtained in Corollary 4.5 show the impacts of the variances associated with both the sketching matrix \mathbf{S} (i.e., σ_s^2) and the data sample ξ (i.e., σ_g^2), as well as the data heterogeneity (i.e., σ_h^2), on FSLoRA’s convergence. Furthermore, the parameter \tilde{L} provides insight into how the sketching operation influences the convergence rate. Increasing k_i would lead to a faster convergence rate

for FSLoRA. However, on the other hand, as k_i increases, the communication and computation costs at device i will increase. In other words, there is a trade-off in the selection of the sketching ratios. Additionally, Corollary 4.5 suggests that FSLoRA achieves a similar linear speedup in the number of local updates and the number of devices as FedAvg (Yu et al., 2019; Khaled et al., 2020).

5. Experiments

In this section, we present experiments to evaluate the performance of the proposed FSLoRA. All the experiments are conducted on a cluster equipped with 4 NVIDIA A100 GPUs, each with 40 GB of memory. The number of devices is set to 20 in our experiments. Other hyperparameter configurations can be found in Appendix A.

5.1. Models and Datasets

Our experiments are based on RoBERTa (125M) (Liu, 2019) and LLaMA-3.2-3B (3.21B) (Dubey et al., 2024). For RoBERTa, we fine-tune and evaluate it on the GLUE benchmark (Wang, 2018), which includes QNLI, MRPC, CoLA, MNLI, RTE, SST-2, and QQP tasks. For LLaMA-3.2-3B, we fine-tune and evaluate it on the Commonsense170K dataset (Hu et al., 2023), covering eight commonsense reasoning question-answering tasks: ARC-c, ARC-e, BoolQ, HellaSwag, OBQA, PIQA, SIQA, and WinoGrande. Further details on these datasets are provided in Appendix B.

5.2. Main Results

Baselines for Heterogeneous LoRA Settings: We consider the following state-of-the-art baselines that integrate LoRA with FL: HeteroLoRA (Cho et al., 2024), FlexLoRA (Bai et al., 2024), and FedStackLoRA (Wang et al., 2024). In our approach, the rank of the global LoRA modules is fixed as $r = 64$, while the sketching ratio for device i is sampled from the set $\{0.125, 0.25, 0.5, 0.75\}$. For a fair comparison, we apply the same rank configuration to all other baselines as in FSLoRA.

RoBERTa Model: Table 5.1 presents a performance comparison between the proposed FSLoRA and baseline methods in the heterogeneous LoRA scenario on the GLUE benchmark with the RoBERTa model. As shown in Table 5.1, compared with the baselines, our approach boosts the average performance across these seven tasks by a noticeable margin. Concretely, FSLoRA consistently outperforms HeteroLoRA across all the considered tasks. In addition, our approach significantly outperforms FlexLoRA in most tasks, with marginal performance differences only in QNLI and QQP. Similarly, it outperforms FedStackLoRA on all tasks with the exception of the QQP task. It is worth noting that FSLoRA and HeteroLoRA maintain similar simplicity

Table 5.1: Testing accuracy over 3 independent runs for fine-tuning the RoBERTa model on the GLUE benchmark. FSLoRA achieves a notable improvement in average performance compared to the baselines.

Method	GPU hours	QNLI	MRPC	CoLA	MNLI	RTE	SST-2	QQP	Avg.
HeteroLoRA	10.7h	87.5 \pm 0.5	84.4 \pm 0.9	75.3 \pm 1.2	66.3 \pm 0.8	69.0 \pm 1.7	89.5 \pm 0.0	85.3 \pm 0.1	79.6
FlexLoRA	12.6h	88.5 \pm 0.2	81.2 \pm 0.4	77.5 \pm 1.2	63.0 \pm 0.5	62.2 \pm 1.9	92.8 \pm 0.4	87.4 \pm 0.1	78.9
FedStackLoRA	12.3h	87.2 \pm 0.3	78.1 \pm 0.7	77.4 \pm 1.7	74.6 \pm 0.5	54.4 \pm 2.1	93.4 \pm 0.1	87.1 \pm 0.3	78.9
FSLoRA	10.9h	88.0 \pm 0.3	87.3 \pm 0.2	82.2 \pm 0.5	76.4 \pm 0.2	69.8 \pm 1.2	93.5 \pm 0.1	85.8 \pm 0.1	83.3

Table 5.2: Testing accuracy over 3 independent runs for fine-tuning the LLaMA-3.2-3B model on the commonsense reasoning benchmark. FSLoRA demonstrates consistent performance improvement across these tasks compared to baselines.

Method	GPU hours	ARC-c	ARC-e	BoolQ	HellaSwag	OBQA	PIQA	SIQA	WinoGrande	Avg.
HeteroLoRA	44.2h	69.2 \pm 0.2	84.6 \pm 0.2	68.4 \pm 0.5	80.0 \pm 0.7	69.9 \pm 0.0	77.3 \pm 0.0	68.7 \pm 0.3	72.0 \pm 0.3	73.8
FlexLoRA	60.8h	69.9 \pm 0.3	84.7 \pm 0.4	66.9 \pm 0.2	80.5 \pm 0.4	72.3 \pm 0.1	78.1 \pm 0.2	70.4 \pm 0.4	73.3 \pm 0.5	74.5
FedStackLoRA	56.2h	67.5 \pm 0.7	83.1 \pm 0.5	65.8 \pm 0.9	78.4 \pm 0.5	69.2 \pm 0.7	75.5 \pm 0.6	67.1 \pm 0.3	71.5 \pm 0.5	72.3
FSLoRA	44.5h	73.8 \pm 0.6	86.2 \pm 0.1	68.5 \pm 0.1	83.1 \pm 1.1	78.7 \pm 0.3	82.0 \pm 0.2	75.8 \pm 0.0	74.8 \pm 0.6	77.9

in computation and communication complexity, whereas FlexLoRA and FedStackLoRA incur additional overhead, as detailed in Section 2.2. This increased computational complexity is also reflected in GPU hours, which represent the total computational time accumulated across all devices. As reported in Table 5.1, FSLoRA and HeteroLoRA demonstrate comparable computational efficiency, requiring similar GPU hours, whereas FlexLoRA and FedStackLoRA incur higher computational costs.

In Figure 2, we present the averaged testing accuracy of the proposed algorithm alongside the baseline methods across seven tasks described in Section 5.1, plotted against the communication rounds, showcasing the convergence behavior of the proposed algorithm. FlexLoRA achieves fast initial convergence but falls behind in final accuracy due to the approximation errors introduced by its use of truncated SVD, which is effective early on but limits long-term performance. Similarly, HeteroLoRA’s reliance on zero-padding reduces optimization efficiency, preventing it from achieving higher accuracy. FedStackLoRA also underperforms due to the random reinitialization of LoRA modules at each communication round, disrupting tuning continuity. In contrast, FSLoRA overcomes these limitations, achieving the highest final accuracy among all methods.

LLaMA-3.2-3B Model: In Table 5.2, we scale up to LLaMA-3.2-3B, which contains 3.21 billion parameters, to evaluate the effectiveness of our proposed algorithm with increasing model complexity. The experimental results demonstrate that FSLoRA consistently achieves superior performance compared to baseline methods across all benchmark tasks. This advantage is particularly significant because fine-tuning LLaMA-3.2-3B, with its substantially larger parameter count, introduces greater challenges for deployment on heterogeneous distributed devices with limited

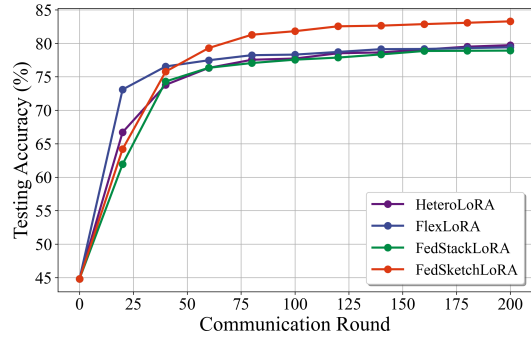


Figure 2: Convergence behavior of FSLoRA and baselines on the GLUE benchmark with the RoBERTa model. Testing accuracy is averaged over seven tasks.

computational and communication resources. This further demonstrates the superiority of the proposed FSLoRA.

5.3. Ablation Study

Impact of Sketching: In Figures 3 and 4(a), we compare the performance of FSLoRA with and without sketching on fine-tuning the RoBERTa model and the LLaMA-3.2-3B model, respectively. For FSLoRA with sketching, we apply a uniform sketching ratio of $k_i/r = 0.5$ across all distributed devices. Notably, FSLoRA without sketching is equivalent to the vanilla federated LoRA. The upload budget for each device is set to 100 and 400 times the size of the full global LoRA modules at the corresponding rank for the RoBERTa and the LLaMA-3.2-3B models, respectively. As shown in Figures 3 and 4(a), both FSLoRA with and without sketching achieve higher accuracy when the rank r increases, due to the availability of more tunable parameters. In addition, FSLoRA consistently outperforms its non-sketched counterpart across all the ranks and datasets. The use of sketching

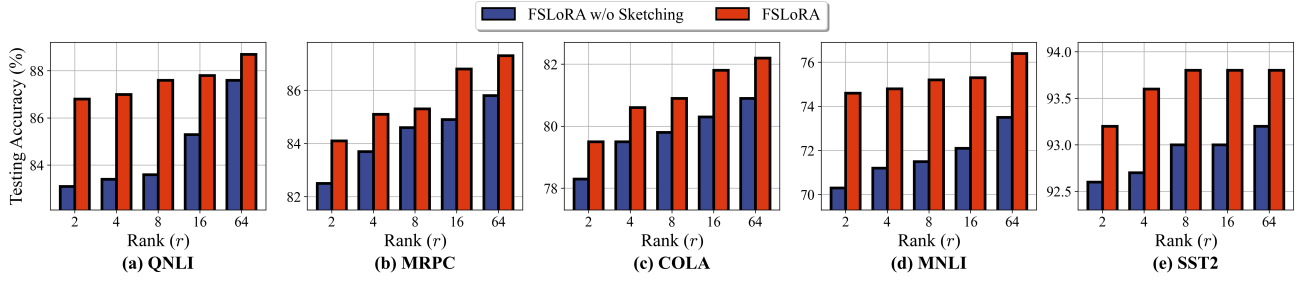
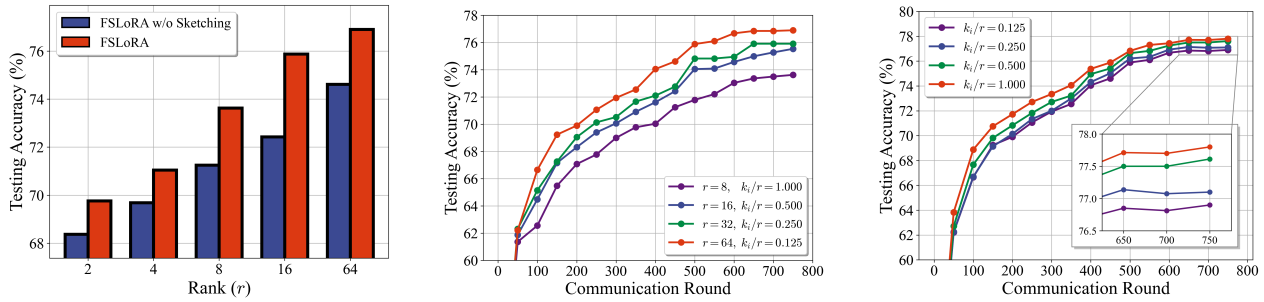


Figure 3: Comparison between FSLoRA with and without sketching, where the upload budget for devices is set to $100\times$ the full global LoRA modules at the corresponding rank. The experiment is performed on the GLUE benchmark and the RoBERTa model. FSLoRA with sketching obtains a better performance, validating the effectiveness of sketching.



(a) Comparison of FSLoRA with and without sketching, with an upload budget $400\times$ the global LoRA module size at each rank.

(b) Impact of the rank of global LoRA modules on FSLoRA, given a fixed rank for the updated submatrices at the devices.

(c) Impact of the sketching ratio on FSLoRA under a fixed rank $r = 64$ for the global LoRA modules.

Figure 4: Fine-tuning the LLaMA-3.2-3B model on the commonsense reasoning benchmark. The results are averaged over eight tasks as described in Section 5.1.

increases the communication frequency for devices within the same communication budget, thereby facilitating the optimization process and enhancing fine-tuning efficiency.

Impact of the Global Rank: In Figure 4(b), we investigate the impact of the rank of the global LoRA modules on FSLoRA’s performance. We vary the rank of the global LoRA modules while keeping the rank of submatrices updated by the devices to be consistent (i.e., $k_i = 8$). This ensures that the communication and computational resources on the client side remain unchanged. As illustrated in Figure 4(b), FSLoRA demonstrates improved performance as the global rank increases within the range considered. This observation validates that using a large rank for the global LoRA modules results in a more expressive model. Moreover, the proposed sketching mechanism enables resource-constrained systems to effectively benefit from a large rank.

Impact of Sketching Ratio: Finally, we investigate the impact of the sketching ratio on FSLoRA’s performance by maintaining a constant global LoRA rank $r = 64$ while varying the sketching ratio k_i/r in the range $\{0.125, 0.25, 0.5, 1\}$. From Figure 4(c), we see that there is a slight performance degradation as the sketching ratio decreases, which aligns with our theoretical analysis. This observation reflects the inherent tradeoff: while a larger

sketching ratio enables better convergence, a smaller sketching ratio reduces both computational and communication overhead. Notably, the slight performance degradation observed demonstrates FSLoRA’s ability to effectively balance efficiency and accuracy, underscoring its advantage in scenarios with limited resources.

Further Experiments: Additional results, including detailed comparisons on each task in the commonsense reasoning benchmark corresponding to Figures 4(a) and 4(b), the integration of communication compression and sketching, and the experiments with more devices, are provided in Appendix C.

6. Conclusion and Future Work

We have proposed FSLoRA, a novel on-device collaborative LLM fine-tuning framework that introduces a sketching mechanism to enhance both performance and efficiency in resource-constrained systems. By maintaining large-rank LoRA modules on the server and allowing devices to selectively update submatrices based on the sketching ratios, FSLoRA effectively adapts to heterogeneous communication and computational constraints. We provide a rigorous convergence analysis of FSLoRA that characterizes how

the sketching ratios affect the convergence rate. Finally, we confirmed the effectiveness of FSLoRA through extensive experiments across multiple datasets and models. A direction for future work is to extend FSLoRA beyond LLM fine-tuning and explore its performance in pretraining, which remains an open area for further investigation.

Impact Statement

This paper makes important contributions to on-device LLM fine-tuning by developing a resource-adaptive algorithm for collaborative fine-tuning in resource-constrained systems. The focus of this work is on the technical advancement of LLM fine-tuning algorithms. While this research has potential societal impacts, it primarily addresses technical challenges and does not necessitate a specific discussion on societal consequences.

References

- Bai, J., Chen, D., Qian, B., Yao, L., and Li, Y. Federated fine-tuning of large language models under heterogeneous tasks and client resources, 2024. URL <https://arxiv.org/abs/2402.11505>.
- Bisk, Y., Zellers, R., Gao, J., Choi, Y., et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.
- Byun, Y. and Lee, J. Towards federated low-rank adaptation of language models with rank heterogeneity. *arXiv preprint arXiv:2406.17477*, 2024.
- Chen, C., Feng, X., Zhou, J., Yin, J., and Zheng, X. Federated large language model: A position paper. *arXiv preprint arXiv:2307.08925*, 2023.
- Cho, Y. J., Liu, L., Xu, Z., Fahrezi, A., and Joshi, G. Heterogeneous lora for federated fine-tuning of on-device foundation models. *arXiv preprint arXiv:2401.06432*, 2024.
- Clark, C., Lee, K., Chang, M.-W., Kwiatkowski, T., Collins, M., and Toutanova, K. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*, 2019.
- Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- Demidovich, Y., Malinovsky, G., Shulgin, E., and Richtárik, P. Mast: Model-agnostic sparsified training. *arXiv preprint arXiv:2311.16086*, 2023.
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Fan, D., Messmer, B., and Jaggi, M. On-device collaborative language modeling via a mixture of generalists and specialists. *arXiv preprint arXiv:2409.13931*, 2024.
- Fang, W., Yu, Z., Jiang, Y., Shi, Y., Jones, C. N., and Zhou, Y. Communication-efficient stochastic zeroth-order optimization for federated learning. *IEEE Transactions on Signal Processing*, 70:5058–5073, 2022.
- Fang, W., Han, D.-J., Chen, E., Wang, S., and Brinton, C. G. Hierarchical federated learning with multi-timescale gradient correction. *arXiv preprint arXiv:2409.18448*, 2024.
- Han, A., Li, J., Huang, W., Hong, M., Takeda, A., Jawanpuria, P., and Mishra, B. Sltrain: a sparse plus low-rank approach for parameter and memory efficient pretraining. *arXiv preprint arXiv:2406.02214*, 2024.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Hu, Z., Wang, L., Lan, Y., Xu, W., Lim, E.-P., Bing, L., Xu, X., Poria, S., and Lee, R. K.-W. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models. *arXiv preprint arXiv:2304.01933*, 2023.
- Kalajdziewski, D. A rank stabilization scaling factor for fine-tuning with lora, 2023. URL <https://arxiv.org/abs/2312.03732>.
- Khaled, A., Mishchenko, K., and Richtárik, P. Tighter theory for local sgd on identical and heterogeneous data. In *International Conference on Artificial Intelligence and Statistics*, pp. 4519–4529. PMLR, 2020.
- Koo, J., Jang, M., and Ok, J. Towards robust and efficient federated low-rank adaptation with heterogeneous clients, 2024. URL <https://arxiv.org/abs/2410.22815>.
- Kuo, K., Raje, A., Rajesh, K., and Smith, V. Federated lora with sparse communication, 2024.
- Lester, B., Al-Rfou, R., and Constant, N. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- Li, X. L. and Liang, P. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021.

- Lialin, V., Muckatira, S., Shivagunde, N., and Rumshisky, A. Relora: High-rank training through low-rank updates. In *The Twelfth International Conference on Learning Representations*, 2023.
- Liu, Y. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 364, 2019.
- Malinovsky, G., Michieli, U., Hammoud, H. A. A. K., Ceritli, T., Elesedy, H., Ozay, M., and Richtárik, P. Randomized asymmetric chain of lora: The first meaningful theoretical framework for low-rank adaptation. *arXiv preprint arXiv:2410.08305*, 2024.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.
- Mihaylov, T., Clark, P., Khot, T., and Sabharwal, A. Can a suit of armor conduct electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*, 2018.
- Sakaguchi, K., Bras, R. L., Bhagavatula, C., and Choi, Y. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- Sap, M., Rashkin, H., Chen, D., LeBras, R., and Choi, Y. Socialliqa: Commonsense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- Shuttleworth, R., Andreas, J., Torralba, A., and Sharma, P. Lora vs full fine-tuning: An illusion of equivalence. *arXiv preprint arXiv:2410.21228*, 2024.
- Sun, Y., Li, Z., Li, Y., and Ding, B. Improving lora in privacy-preserving federated learning. *arXiv preprint arXiv:2403.12313*, 2024.
- Wang, A. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018.
- Wang, J., Liu, Q., Liang, H., Joshi, G., and Poor, H. V. A novel framework for the analysis and design of heterogeneous federated learning. *IEEE Transactions on Signal Processing*, 69:5234–5249, 2021.
- Wang, Z., Shen, Z., He, Y., Sun, G., Wang, H., Lyu, L., and Li, A. Flora: Federated fine-tuning large language models with heterogeneous low-rank adaptations. *arXiv preprint arXiv:2409.05976*, 2024.
- Xia, W., Qin, C., and Hazan, E. Chain of lora: Efficient fine-tuning of language models via residual learning. *arXiv preprint arXiv:2401.04151*, 2024.
- Ye, R., Wang, W., Chai, J., Li, D., Li, Z., Xu, Y., Du, Y., Wang, Y., and Chen, S. Openfedllm: Training large language models on decentralized private data via federated learning. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 6137–6147, 2024.
- Yi, X., Zhang, S., Yang, T., and Johansson, K. H. Zeroth-order algorithms for stochastic distributed nonconvex optimization. *Automatica*, 142:110353, 2022.
- Yu, H., Yang, S., and Zhu, S. Parallel restarted sgd with faster convergence and less communication: Demystifying why model averaging works for deep learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 5693–5700, 2019.
- Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
- Zhang, J., Vahidian, S., Kuo, M., Li, C., Zhang, R., Yu, T., Wang, G., and Chen, Y. 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, 2024.

A. Details of Hyperparameters

Table A.1: The hyperparameters for RoBERTa & GLUE and LLaMA-3.2-3B & Commensense Reasoning benchmarks.

Hyperparameter	RoBERTa & GLUE	LLaMA-3.2-3B & Commensense Reasoning
Batch size	16	16
LoRA dropout rate	0.1	0.1
Learning rate, γ	5e-4	3e-4
Communication round, T	200	750
Local iteration number, H	50	20
Number of edge devices, N	20	20
Target module	["query", "value", "classification head"]	["q_proj", "k_proj", "v_proj", "up_proj", "down_proj"]

B. Details of Datasets

B.1. GLUE Benchmark

GLUE is a widely recognized benchmark designed to assess the natural language understanding capabilities of language models (Wang, 2018).

- **CoLA** focuses on whether a given sentence is acceptable according to linguistic rules. It evaluates a model’s ability to recognize well-formed sentences.
 - ▷ Input: A single sentence.
 - ☆ Output: A label indicating whether the sentence is acceptable or unacceptable.
- **SST-2** is designed for sentiment classification on movie reviews or short texts. It tests whether a model can correctly identify positive or negative sentiment in a given sentence.
 - ▷ Input: A single sentence.
 - ☆ Output: A label indicating positive or negative sentiment.
- **MRPC** checks if two sentences are paraphrases of each other, i.e., if they mean the same thing.
 - ▷ Input: Two sentences (‘sentence1’ and ‘sentence2’).
 - ☆ Output: A label indicating either equivalent or not equivalent.
- **QQP** tests a model’s ability to determine if two questions ask the same thing.
 - ▷ Input: Two questions.
 - ☆ Output: A label indicating duplicate or not duplicate.
- **MNLI** tests whether a given hypothesis is entailed, contradicted, or neutral with respect to a premise.
 - ▷ Input: A premise (first sentence) and a hypothesis (second sentence).
 - ☆ Output: A label indicating entailment, contradiction, or neutral.
- **QNLI** aims to determine if a context sentence correctly answers a given question.
 - ▷ Input: A question and a sentence.
 - ☆ Output: A label indicating the sentence answers the question or it does not.
- **RTE** provides pairs of sentences to see if one implies the other.
 - ▷ Input: Two sentences (‘sentence1’ and ‘sentence2’)
 - ☆ Output: A label indicating whether the meaning of one sentence is entailed from the other one.

B.2. Commonsense Reasoning Benchmark

Table B.1: The prompt template of the Commonsense170K dataset (Hu et al., 2023).

Dataset	Input Template
ARC-c/e	<p>Please choose the correct answer to the question: [QUESTION] Answer1: [ANSWER_1] Answer2: [ANSWER_2] Answer3: [ANSWER_3] Answer4: [ANSWER_4] Answer format: answer1/answer2/answer3/answer4 the correct answer is [ANSWER]</p>
BoolQ	<p>Please answer the following question with true or false, question: [QUESTION] Answer format: true/false the correct answer is [ANSWER]</p>
HellaSwag	<p>Please choose the correct ending to complete the given sentence: [ACTIVITY_LABEL]: [CONTEXT] Ending1: [ENDING_1] Ending2: [ENDING_2] Ending3: [ENDING_3] Ending4: [ENDING_4] Answer format: ending1/ending2/ending3/ending4 the correct answer is [ANSWER]</p>
OBQA	<p>Please choose the correct answer to the question: [QUESTION] Answer1: [ANSWER_1] Answer2: [ANSWER_2] Answer3: [ANSWER_3] Answer4: [ANSWER_4] Answer format: answer1/answer2/answer3/answer4 the correct answer is [ANSWER]</p>
PIQA	<p>Please choose the correct solution to the question: [QUESTION] Solution1: [SOLUTION_1] Solution2: [SOLUTION_2] Answer format: solution1/solution2 the correct answer is [ANSWER]</p>
SIQA	<p>Please choose the correct answer to the question: [QUESTION] Answer1: [ANSWER_1] Answer2: [ANSWER_2] Answer3: [ANSWER_3] Answer format: answer1/answer2/answer3 the correct answer is [ANSWER]</p>
WinoGrande	<p>Please choose the correct answer to fill in the blank to complete the given sentence: [SENTENCE] Option1: [OPTION_1] Option2: [OPTION_2] the correct answer is [ANSWER]</p>

The Commonsense170K dataset is a mixture of multiple datasets including about 170K training samples from ARC-c/e (Clark et al., 2018), BoolQ (Clark et al., 2019), HellaSwag (Zellers et al., 2019), OBQA (Mihaylov et al., 2018), PIQA (Bisk et al., 2020), SIQA (Sap et al., 2019), and WinoGrande (Sakaguchi et al., 2021) datasets.

- **ARC-c/e** contains the challenge and easy question set from the ARC dataset of genuine grade-school level, multiple-choice science questions.
- **BoolQ** is a question-answering dataset with yes/no questions derived from natural, real-world scenarios.
- **HellaSwag** includes questions for commonsense natural language inference, where a context and multiple endings are given, requiring the most coherent ending to be selected.
- **OBQA** involves multi-step problem-solving that combines commonsense knowledge, reasoning, and comprehension of accompanying textual information.
- **PIQA** focuses on questions requiring physical commonsense to solve. Each question offers two answer choices.
- **SIQA** targets reasoning about human actions and their social implication.
- **WinoGrande** is designed as a binary-choice fill-in-the-blank task, this dataset evaluates the ability to resolve ambiguous sentences through commonsense reasoning.

The input template, i.e., prompt format for these datasets is detailed in Table B.1.

C. Further Experiments

In this section, we provide additional results, including detailed per-task comparisons from the commonsense reasoning benchmark corresponding to Figures 4(a) and 4(b), the investigation of the integration of communication compression and sketching, and the experiments with more devices.

C.1. Further Details of Ablation Study

Impact of Sketching: In Figure 5, we compare the performance of FSLoRA with and without sketching on eight tasks from the commonsense reasoning benchmark using the LLaMA-3.2-3B model. For FSLoRA with sketching, we apply a uniform sketching ratio of $k_i/r = 0.5$ across all distributed devices. The uploading budget for each device is set to 400 times the size of the full global LoRA modules at the corresponding rank. It is clear that FSLoRA with sketching consistently outperforms its non-sketched counterpart across these eight tasks, demonstrating the effectiveness of sketching in improving performance.

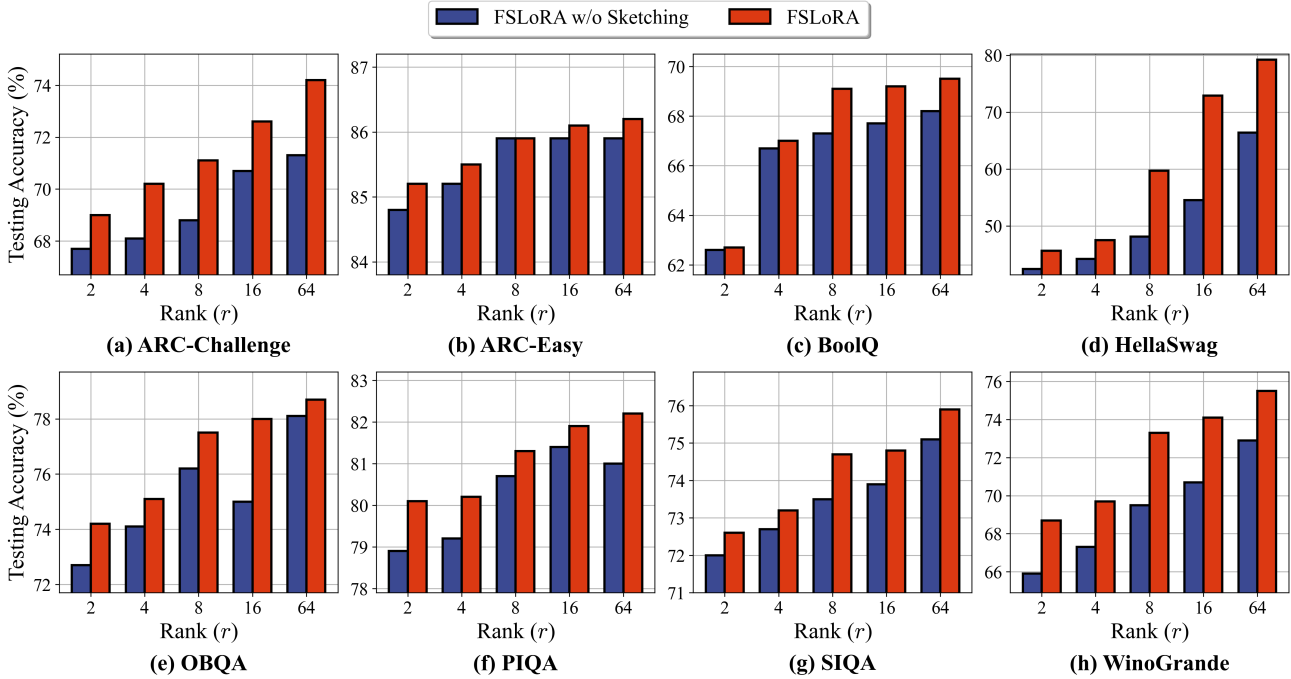


Figure 5: Comparison of FSLoRA with and without sketching, with an upload budget $400\times$ the global LoRA module size at each rank. This is based on the commonsense reasoning benchmark and the LLaMA-3.2-3B model. We observe that the sketching mechanism improves performance across all considered tasks. The average accuracy of the eight tasks is shown in Figure 4(a).

Impact of the Global Rank: In Figure 6, we present the impact of the rank of the global LoRA modules on FSLoRA’s performance across eight tasks from the commonsense reasoning benchmark. We consider four configurations: 1) $r = 8$, $k_i/r = 1$, 2) $r = 16$, $k_i/r = 0.5$, 3) $r = 32$, $k_i/r = 0.25$, and 4) $r = 64$, $k_i/r = 0.125$. The rank of submatrices updated by the devices at each iteration remains consistent across all configurations (i.e., $k_i = 8$), ensuring that the communication and computational resources on the client side are kept fixed for all cases. We see that the third subfigure exhibits oscillations as the sketching ratio increases. One potential explanation for this behavior is that the BoolQ task may be more sensitive to variations in the sketching ratio. Overall, FSLoRA demonstrates improved performance as the global rank increases.

C.2. Integration of Sketching and Top-k Compression

In Figure 7, we investigate the integration of sketching and communication compression, two orthogonal techniques. We employ top-k compression, a widely used compression method that Kuo et al. (2024) introduced to LoRA, at each device to

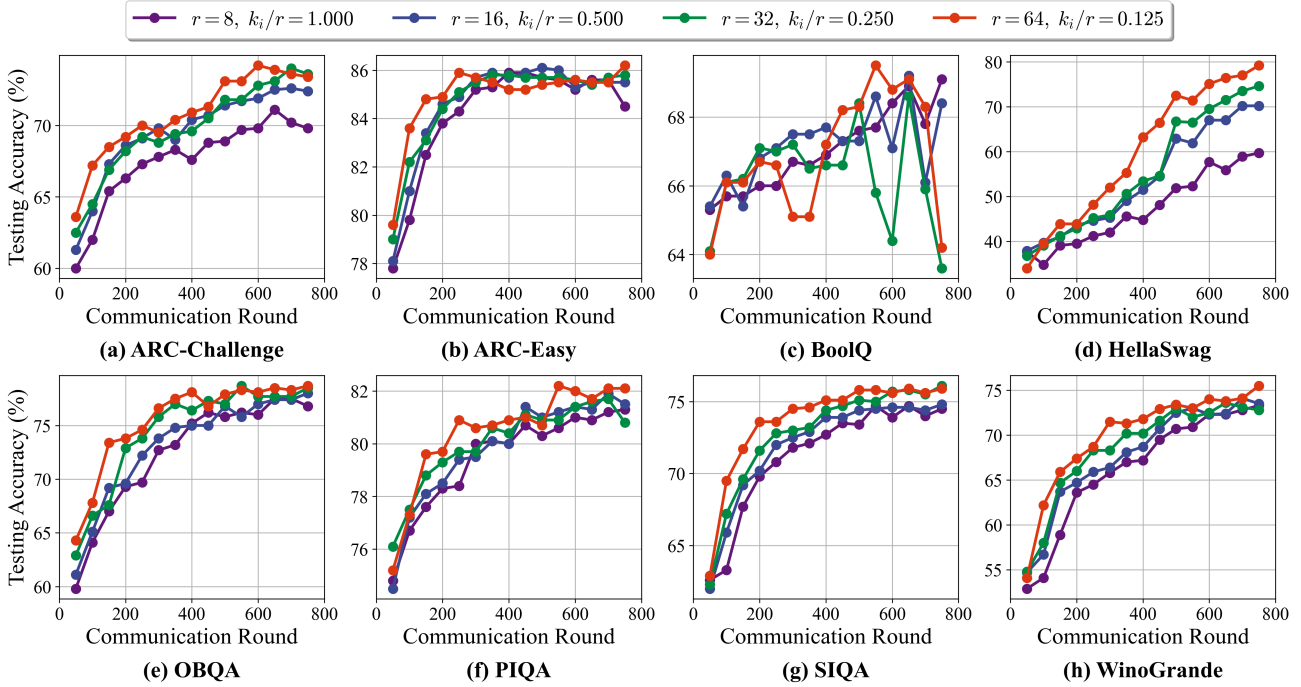


Figure 6: Impact of the rank of global LoRA modules on FSLoRA, given a fixed rank for the updated submatrices at the devices. This is based on the commonsense reasoning benchmark and the LLaMA-3.2-3B model. Overall, FSLoRA demonstrates improved performance as the global rank increases. The average accuracy of the eight tasks is shown in Figure 4(b).

reduce communication overhead before uploading model updates. The compression ratio is fixed at 0.5 for all considered methods, while the sketching ratio k_i/r varies in the range $\{0.125, 0.25, 0.5, 1\}$. Notably, FSLoRA with sketching ratio $k_i/r = 1$ is equivalent to the vanilla federated LoRA, meaning no sketching is applied. Figure 7 plots testing accuracy against communication overhead, where the x-axis represents the per-device upload communication load (MB). The results demonstrate that incorporating sketching further enhances efficiency, with lower sketching ratios leading to higher testing accuracy for the same communication cost, highlighting the benefits of combining these two orthogonal techniques.

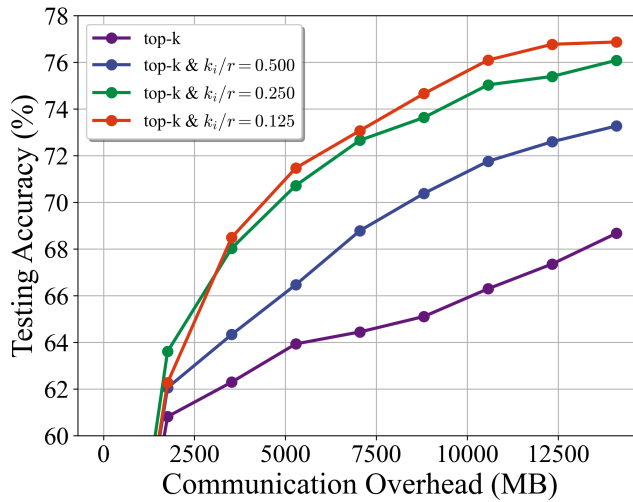


Figure 7: Comparison of top-k compression and its integration with sketching, evaluated on the commonsense reasoning benchmark using the LLaMA-3.2-3B model. The results show that combining these two orthogonal techniques significantly enhances performance, demonstrating the benefits of integrating sketching with top-k compression.

C.3. Experiments with More Devices

Finally, we increase the number of devices to 50 and repeat the comparison of FSLoRA with and without sketching. Figure 8 shows the detailed per-task comparison while Figure 9 shows the averaged performance comparison across these eight tasks from the commonsense reasoning benchmark. We vary the rank of the global LoRA modules in the range $\{4, 16, 64\}$. The sketching ratio is set to 0.5 for FSLoRA. This experiment is also based on the LLaMA-3.2-3B model. As shown in Figures 8 and 9, the advantages of FSLoRA are preserved as the number of devices increases.

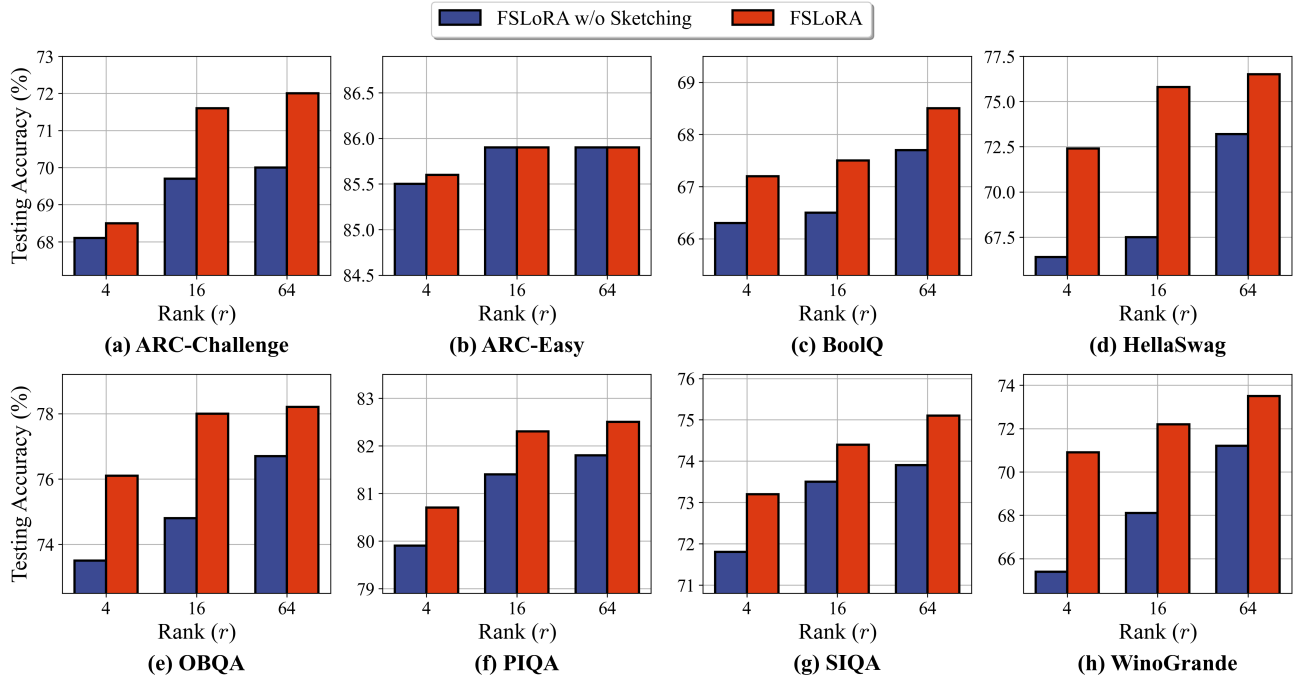


Figure 8: Comparison of FSLoRA with and without sketching, with an upload budget $400 \times$ the global LoRA module size at each rank, evaluated on the commonsense reasoning benchmark and the LLaMA-3.2-3B model. The number of devices is set to 50. We observe that the sketching mechanism improves performance across all considered tasks. The average accuracy of the eight tasks is shown in Figure 9.

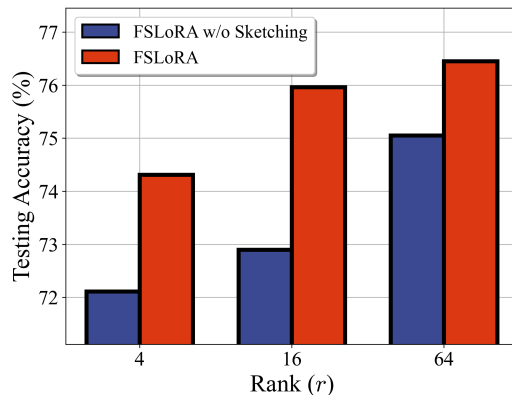


Figure 9: Comparison of FSLoRA with and without sketching, with an upload budget $400 \times$ the global LoRA module size at each rank, evaluated on the LLaMA-3.2-3B model. The number of devices is set to 50. The results are averaged over eight tasks from the commonsense reasoning benchmark.

D. Proof of the Theoretical Results

D.1. Preliminaries

Before presenting the proof of the main results, we first introduce some preliminary facts that will be used later.

Lemma D.1. *Suppose a sequence of independent random matrices $\{\mathbf{P}_i\}_{i=1}^N$ satisfy $\mathbb{E}[\mathbf{P}_i] = \mathbf{0}, \forall i$. Then,*

$$\mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \mathbf{P}_i \right\|^2 = \frac{1}{N^2} \sum_{i=1}^N \mathbb{E} \|\mathbf{P}_i\|^2.$$

Lemma D.2. (Wang et al., 2021) *Suppose a sequence of random matrices $\{\mathbf{P}_i\}_{i=1}^N$ satisfy $\mathbb{E}[\mathbf{P}_i | \mathbf{P}_{i-1}, \mathbf{P}_{i-2}, \dots, \mathbf{P}_1] = \mathbf{0}, \forall i$. Then,*

$$\mathbb{E} \left[\left\| \sum_{i=1}^N \mathbf{P}_i \right\|^2 \right] = \sum_{i=1}^N \mathbb{E} \left[\|\mathbf{P}_i\|^2 \right].$$

Lemma D.3 (Random sketching bounds). *Let \mathbf{S} be a random diagonal sketching matrix of the form*

$$\mathbf{S} = \frac{r}{k} \sum_{j \in \mathcal{I}} \mathbf{e}_j \mathbf{e}_j^\top,$$

where $\mathbf{e}_1, \dots, \mathbf{e}_r \in \mathbb{R}^r$ are standard unit basis vectors and $\mathcal{I} \subseteq \{1, \dots, r\}$ is chosen uniformly at random with $|\mathcal{I}| = k$. Then any matrix \mathbf{X} we have

$$\|\mathbf{X} \mathbf{S}\|^2 \leq \frac{r^2}{k^2} \|\mathbf{X}\|^2, \quad (8)$$

and in expectation we have

$$\mathbb{E}_{\mathbf{S}} \left[\|\mathbf{X} \mathbf{S}\|^2 \right] \leq \frac{r}{k} \|\mathbf{X}\|^2. \quad (9)$$

Proof. Since \mathbf{S} is diagonal with exactly k diagonal entries equal to $\frac{r}{k}$ and the rest zero, its largest eigenvalue is $\frac{r}{k}$. Squaring gives

$$\mathbf{S} \mathbf{S}^\top = \mathbf{S}^2 \preceq \frac{r^2}{k^2} \mathbf{I},$$

Equivalently,

$$\mathbf{x}^\top (\mathbf{S} \mathbf{S}^\top) \mathbf{x} \leq \frac{r^2}{k^2} \|\mathbf{x}\|^2, \forall \mathbf{x}.$$

Setting $\mathbf{x} = \mathbf{x}_j$ to be the j -th column of \mathbf{X} and summing over j implies

$$\|\mathbf{X} \mathbf{S}\|^2 = \sum_j \|\mathbf{S}^\top \mathbf{x}_j\|^2 = \sum_j \mathbf{x}_j^\top (\mathbf{S} \mathbf{S}^\top) \mathbf{x}_j \leq \frac{r^2}{k^2} \sum_j \|\mathbf{x}_j\|^2 = \frac{r^2}{k^2} \|\mathbf{X}\|^2,$$

which proves (8).

For the expected bound (9), note that each diagonal index $j \in \{1, \dots, r\}$ is included in \mathcal{I} with probability $\frac{k}{r}$. Hence the expectation of \mathbf{S}^2 satisfies

$$\mathbb{E}[\mathbf{S}^2] = \frac{r^2}{k^2} \mathbb{E} \left[\sum_{j \in \mathcal{I}} \mathbf{e}_j \mathbf{e}_j^\top \right] = \frac{r^2}{k^2} \frac{k}{r} \mathbf{I} = \frac{r}{k} \mathbf{I}.$$

Thus for any vector \mathbf{x} ,

$$\mathbb{E}_{\mathbf{S}} \left[\|\mathbf{S}^\top \mathbf{x}\|^2 \right] = \mathbb{E}_{\mathbf{S}} \left[\mathbf{x}^\top \mathbf{S} \mathbf{S}^\top \mathbf{x} \right] = \mathbf{x}^\top \left(\mathbb{E}[\mathbf{S}^2] \right) \mathbf{x} = \frac{r}{k} \|\mathbf{x}\|^2.$$

Summing over columns of \mathbf{X} again establishes

$$\mathbb{E}_{\mathbf{S}} \left[\|\mathbf{X} \mathbf{S}\|^2 \right] = \sum_j \mathbb{E}_{\mathbf{S}} \left[\|\mathbf{S}^\top \mathbf{x}_j\|^2 \right] = \sum_j \mathbf{x}_j^\top \left(\mathbb{E}[\mathbf{S}^2] \right) \mathbf{x}_j = \frac{r}{k} \|\mathbf{X}\|^2.$$

This completes the proof of Lemma D.3. □

D.2. Proof of Lemma 3.2

From the chain rule for matrix calculus, we know that:

$$\nabla_{\mathbf{Y}}g(\mathbf{XY}) = \mathbf{X}^\top \nabla g(\mathbf{XY}), \quad \nabla_{\mathbf{X}}g(\mathbf{XY}) = \nabla g(\mathbf{XY})\mathbf{Y}^\top,$$

where $\nabla g(\mathbf{XY})$ denotes the gradient of g to \mathbf{XY} . Applying this to $\ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)$, we proceed as follows: To compute the gradient with respect to \mathbf{B} , set $\mathbf{X} = \mathbf{B}$ and $\mathbf{Y} = \mathbf{SA}$:

$$\nabla_{\mathbf{B}}\ell(\mathbf{W}_0 + \mathbf{BSA}, \xi) = \nabla\ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)(\mathbf{SA})^\top.$$

Similarly, to compute the gradient with respect to \mathbf{A} , set $\mathbf{X} = \mathbf{BS}$ and $\mathbf{Y} = \mathbf{A}$:

$$\nabla_{\mathbf{A}}\ell(\mathbf{W}_0 + \mathbf{BSA}, \xi) = \mathbf{S}^\top \mathbf{B}^\top \nabla\ell(\mathbf{W}_0 + \mathbf{BSA}, \xi).$$

D.3. Proof of Theorem 4.4

The proof of Theorem 4.4 relies on the following proposition.

Proposition D.4. *Under Assumption 4.1, $\tilde{f}_i(\mathbf{X}; \mathbf{S}) = f_i(\mathbf{BS}, \mathbf{A})$, $\mathbf{S} \in \mathcal{S}_i$, $f_i^{\mathcal{S}}(\mathbf{X}) = \mathbb{E}_{\mathbf{S} \sim \mathcal{S}_i}[\tilde{f}_i(\mathbf{X}; \mathbf{S})]$, and $f^{\mathcal{S}}(\mathbf{X}) = \frac{1}{N} \sum_{i=1}^N f_i^{\mathcal{S}}(\mathbf{X})$ are smooth with parameters $L_{\frac{r^2}{k_i^2}}$, $L_{\frac{r}{k_i}}$, and $\left(\frac{1}{N} \sum_{i=1}^N \frac{r}{k_i}\right) L$, respectively.*

The proof of Proposition D.4 is deferred to Appendix D.4. With this proposition, we are ready to prove Theorem 4.4.

In FSLoRA, the update direction in (4) corresponds to the negative stochastic gradient of $\ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)$ with respect to $[\mathbf{B}; \mathbf{A}]$ for a given sketch \mathbf{S}_i^t . We have defined $\tilde{\ell}(\mathbf{X}, \xi; \mathbf{S}) = \ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)$. The iterative equation for the proposed FSLoRA algorithm thus can be written as

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \gamma \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} \tilde{\ell}(\mathbf{X}_i^{t,h}, \xi_i^{t,h}; \mathbf{S}_i^t), \quad (10)$$

where $\mathbf{g}_i^{t,h}$ denotes the stochastic gradient $\nabla_{\mathbf{X}} \tilde{\ell}(\mathbf{X}_i^{t,h}, \xi_i^{t,h}; \mathbf{S}_i^t)$. Based on the smoothness of $f^{\mathcal{S}}(\mathbf{X})$, i.e., Proposition D.4, we have

$$\mathbb{E}[f^{\mathcal{S}}(\mathbf{X}^{t+1})] \leq \mathbb{E}[f^{\mathcal{S}}(\mathbf{X}^t)] - \underbrace{\mathbb{E} \left\langle \nabla_{\mathbf{X}} f^{\mathcal{S}}(\mathbf{X}^t), \gamma \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \mathbf{g}_i^{t,h} \right\rangle}_{T_1} + \frac{\gamma^2 L_s}{2} \underbrace{\mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \mathbf{g}_i^{t,h} \right\|^2}_{T_2}, \quad (11)$$

where $L_s = \left(\frac{1}{N} \sum_{i=1}^N \frac{r}{k_i}\right) L$.

For T_1 , we have

$$\begin{aligned}
 T_1 &= -H\mathbb{E}\left\langle \nabla_{\mathbf{X}} f^{\mathcal{S}}(\mathbf{X}^t), \gamma \frac{1}{NH} \sum_{i=1}^N \sum_{h=0}^{H-1} \mathbf{g}_i^{t,h} \right\rangle \\
 &= -H\mathbb{E}\left\langle \nabla_{\mathbf{X}} f^{\mathcal{S}}(\mathbf{X}^t), \gamma \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} \tilde{f}_i(\mathbf{X}_i^{t,h}; \mathbf{S}_i^t) \right\rangle \\
 &= -H\mathbb{E}\left\langle \nabla_{\mathbf{X}} f^{\mathcal{S}}(\mathbf{X}^t), \gamma \frac{1}{NH} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}_i^{t,h}) \right\rangle \\
 &= -\frac{\gamma H}{2} \mathbb{E} \|\nabla_{\mathbf{X}} f^{\mathcal{S}}(\mathbf{X}^t)\|^2 - \frac{\gamma H}{2} \mathbb{E} \left\| \frac{1}{NH} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}_i^{t,h}) \right\|^2 \\
 &\quad + \frac{\gamma H}{2} \mathbb{E} \left\| \nabla_{\mathbf{X}} f^{\mathcal{S}}(\mathbf{X}^t) - \frac{1}{NH} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}_i^{t,h}) \right\|^2 \\
 &\leq -\frac{\gamma H}{2} \mathbb{E} \|\nabla_{\mathbf{X}} f^{\mathcal{S}}(\mathbf{X}^t)\|^2 - \frac{\gamma H}{2} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}_i^{t,h}) \right\|^2 \\
 &\quad + \frac{\gamma}{2} \sum_{h=0}^{H-1} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}^t) - \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}_i^{t,h}) \right\|^2 \\
 &\leq -\frac{\gamma H}{2} \mathbb{E} \|\nabla_{\mathbf{X}} f^{\mathcal{S}}(\mathbf{X}^t)\|^2 - \frac{\gamma}{2H} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}_i^{t,h}) \right\|^2 + \frac{\gamma HL^2}{2} \frac{1}{NH} \sum_{i=1}^N \frac{r^2}{k_i^2} \sum_{h=0}^{H-1} \mathbb{E} \|\mathbf{X}_i^{t,h} - \mathbf{X}^t\|^2, \quad (12)
 \end{aligned}$$

where the last inequalities follow Jensen's inequality and Proposition D.4.

For T_2 , we have

$$\begin{aligned}
 &\mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \mathbf{g}_i^{t,h} \right\|^2 \\
 &= \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \mathbf{g}_i^{t,h} \mp \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} \tilde{f}_i(\mathbf{X}_i^{t,h}; \mathbf{S}_i^t) \mp \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}_i^{t,h}) \right\|^2 \\
 &\leq 3\mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \mathbf{g}_i^{t,h} - \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} \tilde{f}_i(\mathbf{X}_i^{t,h}; \mathbf{S}_i^t) \right\|^2 \\
 &\quad + 3\mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} \tilde{f}_i(\mathbf{X}_i^{t,h}; \mathbf{S}_i^t) - \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}_i^{t,h}) \right\|^2 + 3\mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}_i^{t,h}) \right\|^2 \\
 &= \frac{3}{N^2} \sum_{i=1}^N \mathbb{E} \left\| \sum_{h=0}^{H-1} \mathbf{g}_i^{t,h} - \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} \tilde{f}_i(\mathbf{X}_i^{t,h}; \mathbf{S}_i^t) \right\|^2 \\
 &\quad + \frac{3}{N^2} \sum_{i=1}^N \mathbb{E} \left\| \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} \tilde{f}_i(\mathbf{X}_i^{t,h}; \mathbf{S}_i^t) - \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}_i^{t,h}) \right\|^2 + 3\mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}_i^{t,h}) \right\|^2 \\
 &\leq 3H \frac{\sigma_g^2 + \sigma_s^2}{N} + 3\mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^{\mathcal{S}}(\mathbf{X}_i^{t,h}) \right\|^2, \quad (13)
 \end{aligned}$$

where the last equality comes from Lemma D.1 while the last inequality follows Lemma D.2 and Assumption 4.2.

Combining (11), (12), and (13) gives rise to

$$\begin{aligned}
 \mathbb{E}[f^S(\mathbf{X}^{t+1})] &\leq \mathbb{E}[f^S(\mathbf{X}^t)] - \frac{\gamma H}{2} \mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2 - \frac{\gamma}{2H} \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^S(\mathbf{X}_i^{t,h}) \right\|^2 \\
 &\quad + \frac{\gamma H L_s^2}{2} \frac{1}{NH} \sum_{i=1}^N \sum_{h=0}^{H-1} \mathbb{E} \|\mathbf{X}_i^{t,h} - \mathbf{X}^t\|^2 \\
 &\quad + \frac{3\gamma^2 L_s}{2} \left\{ H \frac{\sigma_g^2 + \sigma_s^2}{N} + \mathbb{E} \left\| \frac{1}{N} \sum_{i=1}^N \sum_{h=0}^{H-1} \nabla_{\mathbf{X}} f_i^S(\mathbf{X}_i^{t,h}) \right\|^2 \right\} \\
 &\leq \mathbb{E}[f^S(\mathbf{X}^t)] - \frac{\gamma H}{2} \mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2 + \frac{3\gamma^2 L_s}{2} H \frac{\sigma_g^2 + \sigma_s^2}{N} \\
 &\quad + \underbrace{\frac{\gamma H L^2}{2} \frac{1}{NH} \sum_{i=1}^N \frac{r^2}{k_i^2} \sum_{h=0}^{H-1} \mathbb{E} \|\mathbf{X}_i^{t,h} - \mathbf{X}^t\|^2}_{T_3}, \tag{14}
 \end{aligned}$$

where the second inequality follows the condition $\gamma \leq \frac{1}{3HL_s}$.

For T_3 , we have

$$\begin{aligned}
 T_3 &= \frac{1}{NH} \sum_{i=1}^N \frac{r^2}{k_i^2} \sum_{h=0}^{H-1} \mathbb{E} \|\mathbf{X}_i^{t,h} - \mathbf{X}^t\|^2 \\
 &= \frac{1}{NH} \sum_{i=1}^N \frac{r^2}{k_i^2} \sum_{h=0}^{H-1} \mathbb{E} \left\| \gamma \sum_{\tau=0}^{h-1} \mathbf{g}_i^{t,h} \right\|^2 \\
 &= \gamma^2 \frac{1}{NH} \sum_{i=1}^N \frac{r^2}{k_i^2} \sum_{h=0}^{H-1} \mathbb{E} \left\| \sum_{\tau=0}^{h-1} \mathbf{g}_i^{t,h} \mp \sum_{\tau=0}^{h-1} \nabla_{\mathbf{X}} f_i(\mathbf{X}_i^{t,\tau}, \mathbf{S}_i^t) \mp \sum_{\tau=0}^{h-1} \nabla_{\mathbf{X}} f_i^S(\mathbf{X}_i^{t,\tau}) \right\|^2 \\
 &\leq 3\gamma^2 H (\sigma_g^2 + \sigma_s^2) \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) + 3\gamma^2 \frac{1}{NH} \sum_{i=1}^N \frac{r^2}{k_i^2} \sum_{h=0}^{H-1} h \sum_{\tau=0}^{h-1} \mathbb{E} \|\nabla_{\mathbf{X}} f_i^S(\mathbf{X}_i^{t,\tau})\|^2 \\
 &= 3\gamma^2 H (\sigma_g^2 + \sigma_s^2) \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) + \gamma^2 \frac{3}{NH} \sum_{i=1}^N \sum_{h=0}^{H-1} h \sum_{\tau=0}^{h-1} \mathbb{E} \|\nabla_{\mathbf{X}} f_i^S(\mathbf{X}_i^{t,\tau}) \mp \nabla_{\mathbf{X}} f_i^S(\mathbf{X}^t) \mp \nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2 \\
 &\leq 3\gamma^2 H (\sigma_g^2 + \sigma_s^2) \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) + 9\gamma^2 L^2 \frac{1}{NH} \sum_{i=1}^N \frac{r^2}{k_i^2} \sum_{h=0}^{H-1} h \sum_{\tau=0}^{h-1} \mathbb{E} \|\mathbf{X}_i^{t,\tau} - \mathbf{X}^t\|^2 \\
 &\quad + 9 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) c_h \gamma^2 H^2 \mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2 + 9 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) \gamma^2 H^2 \sigma_h^2 + 9 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) \gamma^2 H^2 \mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2 \\
 &\leq 3 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) \gamma^2 H (\sigma_g^2 + \sigma_s^2) + 9 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) \gamma^2 H^2 \sigma_h^2 \\
 &\quad + 9\gamma^2 L^2 H^2 \underbrace{\frac{1}{NH} \sum_{i=1}^N \frac{r^2}{k_i^2} \sum_{\tau=0}^{H-1} \mathbb{E} \|\mathbf{X}_i^{t,\tau} - \mathbf{X}^t\|^2}_{T_3} + 9 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) (c_h + 1) \gamma^2 H^2 \mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2, \tag{15}
 \end{aligned}$$

where the first inequality equality comes from Lemma D.2 and Assumption 4.2 and the second inequality is due to Assumption 4.3. Therefore, it follows that

$$(1 - 9\gamma^2 L^2 H^2) T_3 \leq 3 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) \gamma^2 H^2 (\sigma_g^2 + \sigma_s^2 + 3\sigma_h^2) + 9 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) (c_h + 1) \gamma^2 H^2 \mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2.$$

Based on the condition on γ outlined in Theorem 4.4, we have $\gamma \leq \frac{1}{\sqrt{18HL}}$. We thus obtain $1 - 9\gamma^2 L^2 H^2 \geq \frac{1}{2}$. It follows that

$$T_3 \leq 6 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) \gamma^2 H^2 (\sigma_g^2 + \sigma_s^2 + 3\sigma_h^2) + 18 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) (c_h + 1) \gamma^2 H^2 \mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2. \quad (16)$$

Plugging (16) into (14) gives rise to

$$\begin{aligned} \mathbb{E}[f^S(\mathbf{X}^{t+1})] &\leq \mathbb{E}[f^S(\mathbf{X}^t)] - \left(\frac{\gamma H}{2} - 9 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) (c_h + 1) \gamma^3 H^3 L^2 \right) \mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2 \\ &\quad + 3 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) \gamma^3 H^3 L^2 (\sigma_g^2 + \sigma_s^2 + 3\sigma_h^2) + \frac{3\gamma^2 L_s}{2} H \frac{\sigma_g^2 + \sigma_s^2}{N} \\ &\leq \mathbb{E}[f^S(\mathbf{X}^t)] - \frac{\gamma H}{4} \mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2 \\ &\quad + 3\gamma^3 H^3 L_s^2 (\sigma_g^2 + \sigma_s^2 + 3\sigma_h^2) + \frac{3\gamma^2 L_s}{2} \frac{H}{N} (\sigma_g^2 + \sigma_s^2). \end{aligned} \quad (17)$$

Due to $\gamma \leq \frac{1}{6\sqrt{\left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2}\right)(c_h+1)HL}}$, we have $\frac{\gamma H}{2} - 9 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) (c_h + 1) \gamma^3 H^3 L^2 \geq \frac{\gamma H}{4}$. Reorganizing the above inequality, we have

$$\mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2 \leq 4 \frac{\mathbb{E}[f^S(\mathbf{X}^t)] - \mathbb{E}[f^S(\mathbf{X}^{t+1})]}{\gamma H} + 12 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) \gamma^2 H^2 L^2 (\sigma_g^2 + \sigma_s^2 + 3\sigma_h^2) + 6\gamma L_s \frac{\sigma_g^2 + \sigma_s^2}{N}.$$

Telescoping the above inequality from $t = 0$ to $T - 1$ and utilizing $\gamma \leq \frac{1}{NH^2L}$, we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla_{\mathbf{X}} f^S(\mathbf{X}^t)\|^2 \leq 4 \frac{f^S(\mathbf{X}^0) - \mathbb{E}[f^S(\mathbf{X}^T)]}{\gamma TH} + 12 \left(\frac{1}{N} \sum_{i=1}^N \frac{r^2}{k_i^2} \right) \frac{\gamma L}{N} (\sigma_g^2 + \sigma_s^2 + 3\sigma_h^2) + 6 \left(\frac{1}{N} \sum_{i=1}^N \frac{r}{k_i} \right) \gamma L \frac{\sigma_g^2 + \sigma_s^2}{N}.$$

This completes the proof of Theorem 4.4.

D.4. Proof of Proposition D.4

i) For illustration, we need to recover \mathbf{X} to $[\mathbf{B}; \mathbf{A}]$ in this proof. According to the definition of $\tilde{f}_i(\mathbf{X}; \mathbf{S})$ and $f_i(\mathbf{B}, \mathbf{A})$, we have

$$\tilde{f}_i(\mathbf{X}; \mathbf{S}) = \tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) \quad (18)$$

$$\begin{aligned} &= \mathbb{E}_{\xi \sim \mathcal{D}_i} [\ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)] \\ &= f_i(\mathbf{BS}, \mathbf{A}). \end{aligned} \quad (19)$$

As $f_i(\mathbf{B}, \mathbf{A})$ is L -smooth, we have

$$f_i(\mathbf{BS} + \Delta\mathbf{BS}, \mathbf{A} + \Delta\mathbf{A}) \leq f_i(\mathbf{BS}, \mathbf{A}) + \left\langle \begin{bmatrix} \nabla_{\mathbf{BS}} f_i(\mathbf{BS}, \mathbf{A}) \\ \nabla_{\mathbf{A}} f_i(\mathbf{BS}, \mathbf{A}) \end{bmatrix}, \begin{bmatrix} \Delta\mathbf{BS} \\ \Delta\mathbf{A} \end{bmatrix} \right\rangle + \frac{L}{2} \left\| \begin{bmatrix} \Delta\mathbf{BS} \\ \Delta\mathbf{A} \end{bmatrix} \right\|^2. \quad (20)$$

According to (18) and (19), we have $\tilde{f}_i(\mathbf{B} + \Delta\mathbf{B}, \mathbf{A} + \Delta\mathbf{A}; \mathbf{S}) = f_i(\mathbf{BS} + \Delta\mathbf{BS}, \mathbf{A} + \Delta\mathbf{A})$ and $\tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) = f_i(\mathbf{BS}, \mathbf{A})$. Combining these with (20) gives rise to

$$\tilde{f}_i(\mathbf{B} + \Delta\mathbf{B}, \mathbf{A} + \Delta\mathbf{A}; \mathbf{S}) \leq \tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) + \left\langle \begin{bmatrix} \nabla_{\mathbf{BS}} f_i(\mathbf{BS}, \mathbf{A}) \\ \nabla_{\mathbf{A}} f_i(\mathbf{BS}, \mathbf{A}) \end{bmatrix}, \begin{bmatrix} \Delta\mathbf{BS} \\ \Delta\mathbf{A} \end{bmatrix} \right\rangle + \frac{L}{2} \left\| \begin{bmatrix} \Delta\mathbf{BS} \\ \Delta\mathbf{A} \end{bmatrix} \right\|^2. \quad (21)$$

We denote

$$L(\mathbf{W}_0 + \mathbf{BSA}) = \tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) = \mathbb{E}_{\xi \sim \mathcal{D}_i} [\ell(\mathbf{W}_0 + \mathbf{BSA}, \xi)]. \quad (22)$$

Note that $\nabla_{\mathbf{BS}} f_i(\mathbf{BS}, \mathbf{A}) = \nabla L(\mathbf{W}_0 + \mathbf{BSA})\mathbf{A}^\top$ and $\nabla_{\mathbf{A}} f_i(\mathbf{BS}, \mathbf{A}) = \mathbf{S}^\top \mathbf{B}^\top \nabla L(\mathbf{W}_0 + \mathbf{BSA})$. We thus have

$$\begin{aligned} \left\langle \begin{bmatrix} \nabla_{\mathbf{BS}} f_i(\mathbf{BS}; \mathbf{A}) \\ \nabla_{\mathbf{A}} f_i(\mathbf{BS}; \mathbf{A}) \end{bmatrix}, \begin{bmatrix} \Delta \mathbf{BS} \\ \Delta \mathbf{A} \end{bmatrix} \right\rangle &= \left\langle \begin{bmatrix} \nabla L(\mathbf{W}_0 + \mathbf{BSA})\mathbf{A}^\top \\ \mathbf{S}^\top \mathbf{B}^\top \nabla L(\mathbf{W}_0 + \mathbf{BSA}) \end{bmatrix}, \begin{bmatrix} \Delta \mathbf{BS} \\ \Delta \mathbf{A} \end{bmatrix} \right\rangle \\ &= \left\langle \begin{bmatrix} \nabla L(\mathbf{W}_0 + \mathbf{BSA})\mathbf{A}^\top \mathbf{S}^\top \\ \mathbf{S}^\top \mathbf{B}^\top \nabla L(\mathbf{W}_0 + \mathbf{BSA}) \end{bmatrix}, \begin{bmatrix} \Delta \mathbf{B} \\ \Delta \mathbf{A} \end{bmatrix} \right\rangle \\ &= \left\langle \begin{bmatrix} \nabla_{\mathbf{B}} \tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) \\ \nabla_{\mathbf{A}} \tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) \end{bmatrix}, \begin{bmatrix} \Delta \mathbf{B} \\ \Delta \mathbf{A} \end{bmatrix} \right\rangle, \end{aligned} \quad (23)$$

where the last equality follows the fact that $\tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) = L(\mathbf{W}_0 + \mathbf{BSA})$ defined in (22) and

$$\begin{bmatrix} \nabla_{\mathbf{B}} \tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) \\ \nabla_{\mathbf{A}} \tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) \end{bmatrix} = \begin{bmatrix} \nabla L(\mathbf{W}_0 + \mathbf{BSA})\mathbf{A}^\top \mathbf{S}^\top \\ \mathbf{S}^\top \mathbf{B}^\top \nabla L(\mathbf{W}_0 + \mathbf{BSA}) \end{bmatrix}.$$

Plugging (23) into (21) gives rise to

$$\tilde{f}_i(\mathbf{B} + \Delta \mathbf{B}, \mathbf{A} + \Delta \mathbf{A}; \mathbf{S}) \leq \tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) + \left\langle \begin{bmatrix} \nabla_{\mathbf{B}} \tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) \\ \nabla_{\mathbf{A}} \tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S}) \end{bmatrix}, \begin{bmatrix} \Delta \mathbf{B} \\ \Delta \mathbf{A} \end{bmatrix} \right\rangle + \frac{L}{2} \left\| \begin{bmatrix} \Delta \mathbf{BS} \\ \Delta \mathbf{A} \end{bmatrix} \right\|^2. \quad (24)$$

In particular, $\left\| \begin{bmatrix} \Delta \mathbf{BS} \\ \Delta \mathbf{A} \end{bmatrix} \right\|^2 = \|\Delta \mathbf{BS}\|^2 + \|\Delta \mathbf{A}\|^2$. From (8), we know $\|\Delta \mathbf{BS}\|^2 \leq \frac{r^2}{k_i^2} \|\Delta \mathbf{B}\|^2$. Therefore, we have $\left\| \begin{bmatrix} \Delta \mathbf{BS} \\ \Delta \mathbf{A} \end{bmatrix} \right\|^2 = \frac{r^2}{k_i^2} \left\| \begin{bmatrix} \Delta \mathbf{B} \\ \Delta \mathbf{A} \end{bmatrix} \right\|^2$. As a result, $\tilde{f}_i(\mathbf{B}, \mathbf{A}; \mathbf{S})$ (i.e., $\tilde{f}_i(\mathbf{X}, \mathbf{S})$) is $L \frac{r^2}{k_i^2}$ -smooth.

ii) Note that $f_i^{\mathcal{S}}(\mathbf{X}) = \mathbb{E}_{\mathbf{S} \sim \mathcal{S}_i} [\tilde{f}_i(\mathbf{X}, \mathbf{S})]$. Therefore, we further take expectation for (24) over $\mathbf{S} \sim \mathcal{S}_i$, leading to

$$f_i^{\mathcal{S}}(\mathbf{B} + \Delta \mathbf{B}, \mathbf{A} + \Delta \mathbf{A}) \leq f_i^{\mathcal{S}}(\mathbf{B}, \mathbf{A}) + \left\langle \begin{bmatrix} \nabla_{\mathbf{B}} f_i^{\mathcal{S}}(\mathbf{B}, \mathbf{A}) \\ \nabla_{\mathbf{A}} f_i^{\mathcal{S}}(\mathbf{B}, \mathbf{A}) \end{bmatrix}, \begin{bmatrix} \Delta \mathbf{B} \\ \Delta \mathbf{A} \end{bmatrix} \right\rangle + \frac{L}{2} \mathbb{E}_{\mathbf{S} \sim \mathcal{S}_i} \left\| \begin{bmatrix} \Delta \mathbf{BS} \\ \Delta \mathbf{A} \end{bmatrix} \right\|^2.$$

In particular, $\mathbb{E}_{\mathbf{S} \sim \mathcal{S}_i} \left\| \begin{bmatrix} \Delta \mathbf{BS} \\ \Delta \mathbf{A} \end{bmatrix} \right\|^2 = \mathbb{E}_{\mathbf{S} \sim \mathcal{S}_i} \|\Delta \mathbf{BS}\|^2 + \|\Delta \mathbf{A}\|^2$. From (9), we know $\mathbb{E}_{\mathbf{S} \sim \mathcal{S}_i} \|\Delta \mathbf{BS}\|^2 \leq \frac{r}{k_i} \|\Delta \mathbf{B}\|^2$. In other words, $\mathbb{E}_{\mathbf{S} \sim \mathcal{S}_i} \left\| \begin{bmatrix} \Delta \mathbf{BS} \\ \Delta \mathbf{A} \end{bmatrix} \right\|^2 = \frac{r}{k_i} \left\| \begin{bmatrix} \Delta \mathbf{B} \\ \Delta \mathbf{A} \end{bmatrix} \right\|^2$. We thus claim that $f_i^{\mathcal{S}}(\mathbf{B}, \mathbf{A})$ (i.e., $f_i^{\mathcal{S}}(\mathbf{X})$) is $L \frac{r}{k_i}$ -smooth.

iii) Finally, for $f^{\mathcal{S}}(\mathbf{X}) = \frac{1}{N} \sum_{i=1}^N f_i^{\mathcal{S}}(\mathbf{X})$, we have

$$\nabla f^{\mathcal{S}}(\mathbf{X}) = \frac{1}{N} \sum_{i=1}^N \nabla f_i^{\mathcal{S}}(\mathbf{X}).$$

Since $f_i^{\mathcal{S}}(\mathbf{X})$ is $L \frac{r}{k_i}$ -smooth, we thus have

$$\|\nabla f_i^{\mathcal{S}}(\mathbf{X}) - \nabla f_i^{\mathcal{S}}(\mathbf{Y})\| \leq L \frac{r}{k_i} \|\mathbf{X} - \mathbf{Y}\|, \quad \forall \mathbf{X}, \mathbf{Y}.$$

To find the Lipschitz constant of $f^{\mathcal{S}}(\mathbf{X})$, we analyze the difference between the gradients at two points \mathbf{X} and \mathbf{Y} :

$$\begin{aligned} \|\nabla f^{\mathcal{S}}(\mathbf{X}) - \nabla f^{\mathcal{S}}(\mathbf{Y})\| &= \left\| \frac{1}{N} \sum_{i=1}^N (\nabla f_i^{\mathcal{S}}(\mathbf{X}) - \nabla f_i^{\mathcal{S}}(\mathbf{Y})) \right\| \\ &\leq \frac{1}{N} \sum_{i=1}^N \|\nabla f_i^{\mathcal{S}}(\mathbf{X}) - \nabla f_i^{\mathcal{S}}(\mathbf{Y})\| \\ &\leq \left(\frac{1}{N} \sum_{i=1}^N \frac{r}{k_i} L \right) \|\mathbf{X} - \mathbf{Y}\|. \end{aligned} \quad (25)$$

Therefore, $f^{\mathcal{S}}(\mathbf{X})$ is $\left(\frac{1}{N} \sum_{i=1}^N \frac{r}{k_i} L \right)$ -smooth.