#### 000 EXPLORING SELECTIVE LAYER FINE-TUNING IN FED-001 002 ERATED LEARNING 003

Anonymous authors

Paper under double-blind review

### ABSTRACT

Federated learning (FL) has emerged as a promising paradigm for fine-tuning foundation models using distributed data in a privacy-preserving manner. Under limited computational resources, clients often find it more practical to fine-tune a selected subset of layers, rather than the entire model, based on their task-specific data. In this study, we provide a thorough theoretical exploration of selective layer fine-tuning in FL, emphasizing a flexible approach that allows the clients to adjust their selected layers according to their local data and resources. We theoretically demonstrate that the layer selection strategy has a significant impact on model convergence in two critical aspects: the importance of selected layers and the heterogeneous choices across clients. Drawing from these insights, we further propose a strategic layer selection method that utilizes local gradients and regulates layer selections across clients. Extensive experiments on both image and text datasets demonstrate the effectiveness of the proposed strategy compared with several baselines, highlighting its advances in identifying critical layers that adapt to the client heterogeneity and training dynamics in FL.

- INTRODUCTION 1
- 027 028

004

010 011

012

013

014

015

016

017

018

019

021

024

025 026

029 Foundation models (Bommasani et al., 2021), including BERT (Devlin et al., 2019), GPT (Radford et al., 2019; Brown et al., 2020), CLIP (Radford et al., 2021; Dosovitskiy et al., 2021), LLaMA (Touvron et al., 2023), and so on (Ramesh et al., 2021; Chowdhery et al., 2023), have attracted considerable 031 attention due to their exceptional ability in handling complex tasks (Eloundou et al., 2023). When 032 it comes to practical deployments of these models in specialized fields, fine-tuning with domain-033 specific data becomes critical. Nevertheless, the distributed nature of data across various users and 034 organizations presents a challenge for centralized storage and training, as it may lead to severe privacy concerns and incur additional transmission costs. Such issues have positioned federated learning (FL) (McMahan et al., 2017) as a promising paradigm to fine-tune foundation models, aligning model 037 enhancement with privacy preservation (Chen et al., 2023a; Kuang et al., 2023).

038 Generally, FL aims to learn a global model through a collaborative process where clients perform local training and upload the parameter updates to a central server for aggregation. Given that clients 040 have limited resources (Bonawitz et al., 2019; Imteaj et al., 2022), such as computational power, 041 communication bandwidth, and available memory, it becomes impractical for them to fine-tune the 042 entire foundation model. Two kinds of solutions have recently emerged to tackle this challenge. 043 The first solution employs *parameter-efficient fine-tuning* techniques (Houlsby et al., 2019; Gao 044 et al., 2021; Hu et al., 2022; Li & Liang, 2021), which introduces additional modules integrated into foundation models and updates these modules with domain-specific data while keeping the parameters of the foundation model frozen. The second one is *selective model fine-tuning* (Lee et al., 046 2019a; Xu et al., 2021; Zhang et al., 2022a; Shen et al., 2021), which only selects an impactful subset 047 of parameters for optimization to streamline the fine-tuning process under resource constraints. 048

This study focuses on selective model fine-tuning as it is particularly well-suited to address the inherent heterogeneity in FL, i.e., the data heterogeneity and device heterogeneity (Yang et al., 051 2021; Chai et al., 2019; Li et al., 2022). Specifically, clients involved in FL have non-independent and identically distributed (non-IID) data and different system resources, leading to the need to 052 customize fine-tuning strategies to such discrepancies. For example, clients with limited computation resources may opt to update only a fraction of the model, while those with sufficient resources and

high-quality data prefer fine-tuning a large portion of the model to enhance performance. Selective
 model fine-tuning enables clients to adjust the chosen part of the model to be updated based on their
 own capabilities, providing a flexible and advanced solution to mitigating sub-optimal issues induced
 by the heterogeneity in FL.

058 The exploration of selective model fine-tuning within the context of FL, is still in its early stages. Previous studies (Shen et al., 2021; Xu et al., 2021; Lee et al., 2022; Dun et al., 2022) have concentrated 060 on designing static strategies for subnetwork selection to improve model fine-tuning performance, 061 without adequately considering heterogeneity among clients. To fulfill this gap, in this paper, we 062 provide a comprehensive theoretical analysis on selective model fine-tuning in FL, focusing on a 063 general scenario where clients are allowed to choose different layers for local training and vary their 064 choices across different training epochs, called *selective layer fine-tuning*. Specifically, we formulate the optimization objective of selective layer fine-tuning in FL, and provide insights on effectively de-065 termining critical layers to achieve model convergence. Building on these insights, we further propose 066 a novel layer selection strategy that leverages local gradients and the regulation of unified selections. 067

068 069

071

073

074

075 076

077

078

079

081

082

084

085

087

Our main contributions are summarized as follows:

- We study a practical FL setup where clients choose to fine-tune some layers of the model, with arbitrary layer selection that may vary among clients and across different training epochs. We theoretically analyze such a training scheme and investigate the impact of layer selection. The analytical results show that the selected layers affect the convergence performance with two critical aspects, namely the importance of layers and heterogeneous choices across clients.
  - Building on the theoretical analysis, we formulate the optimization problem of selective layer fine-tuning considering the limited and diverse resource budgets of clients. Inspired by the solution to this optimization problem, we propose an effective strategy for selecting layers for fine-tuning that are well-suited for the local data and available resources at clients.
    - We conduct experiments to compare the proposed layer selection strategy with baseline methods on both image and text datasets. Experimental results demonstrate the superior performance of the proposed strategy in achieving better model performance, highlighting that the proposed strategy can find critical layers for fine-tuning while considering the client heterogeneity and training dynamics in FL<sup>1</sup>.

## 2 RELATED WORKS

Various approaches have been proposed to properly select a subset of parameters for fine-tuning foundation models within centralized training, including optimizing a non-structured mask matrix (Lee et al., 2019a; Xu et al., 2021; Zhang et al., 2022a; Shen et al., 2021; Zaken et al., 2022; Zhang et al., 2023; Kovaleva et al., 2019; Lee et al., 2019b) and adopting layer-wise selection strategies (Kovaleva et al., 2019; Lee et al., 2019b; 2022; Kaplun et al., 2023). For example, Lee et al. (2019a) suggest updating the parameters in a stochastic manner based on the Bernoulli distribution, while Kovaleva et al. (2019); Lee et al. (2019b) showcase that fine-tuning the top few layers achieves competitive model performance in downstream tasks. Moreover, Lee et al. (2022) propose to select layers according to their gradient statistics.

097 Recent studies have extended the selective fine-tuning techniques to FL scenarios (Nguyen et al., 098 2022a; Chen et al., 2022; Hilmkil et al., 2021; Zhang et al., 2022b). Specifically, researchers (Lee 099 et al., 2023; Dun et al., 2022) investigate layer-wise network decomposition to achieve selective 100 model fine-tuning on clients. However, these works fail to offer methodologies for adaptive and 101 dynamic layer selection that takes into account the heterogeneous characteristics of clients. In 102 addition, personalized FL algorithms (Pillutla et al., 2022; Chen et al., 2023b) propose to train 103 different subnetworks on clients towards better local models. Different from previous studies, we 104 focus on providing an in-depth understanding of selective layer fine-tuning in FL, considering the heterogeneity from the perspective of client resources and local data distributions. 105

<sup>&</sup>lt;sup>1</sup>The source codes are available at https://anonymous.4open.science/r/fed\_selected\_tune/.



Figure 1: An overview of selective layer fine-tuning in FL. The colored layers are selected for fine-tuning.

#### **PROBLEM FORMULATION**

**Federated learning** We consider an FL system with a central server and N clients (denoted by the set  $\mathcal{N} = \{1, \dots, N\}$ ), where each client has a private dataset  $\mathcal{D}_i$  consisting of  $d_i = |\mathcal{D}_i|$  data instances. The server owns a pretrained foundation model  $\theta \in \mathbb{R}^{P}$ , containing P trainable parameters and L layers with the index set  $\mathcal{L} = \{1, 2, \dots, L\}$ . The server aims to fine-tune this foundation model based on the clients' datasets  $\mathcal{D} = \{\mathcal{D}_i\}_{i \in \mathcal{N}}$  but does not directly access these datasets. The learning goal is formally given as: 

$$\min_{\theta \in \mathbb{R}^P} f(\theta) = \sum_{i=1}^{N} \alpha_i f_i(\theta), \tag{1}$$

where  $\alpha_i = \frac{d_i}{\sum_{j=1}^N d_j}$  denotes the relative sample size, and  $f_i(\theta) = \frac{1}{d_i} \sum_{\xi \in \mathcal{D}_i} F_i(\theta; \xi)$  denotes the local training objective of client i. Here we use  $F_i(\theta;\xi)$  to define the (possibly non-convex) loss function computed by the model  $\theta$  on data sample  $\xi$ . The training process of FL is divided into T training epochs. In each epoch  $t \in [T]$ , the server chooses a subset of clients  $S^t$ , and sends the up-to-date global model  $\theta^t$  to these clients for local training. 

Selective layer fine-tuning in FL An overview of selective layer fine-tuning in FL is illustrated in Figure 1. Due to resource limitations, clients tend to update some of the layers in the local training process rather than the entire global model. Formally, we define a masking vector  $\mathbf{m}_i^t \in \{0, 1\}^L$  for each client i. The l-th element  $\mathbf{m}_{i}^{t}(l)$  equals 1 if the l-th layer is selected to be updated in the t-th training epoch, and  $\mathbf{m}_i^t(l) = 0$  otherwise. Accordingly, the selected layer set of client i is denoted by  $\mathcal{L}_i^t \triangleq \{l \in \mathcal{L} | \mathbf{m}_i^t(l) = 1\}$ , and the set for all selected layers in the t-th training epoch is denoted by  $\mathcal{L}_t = \bigcup_{i \in S^t} \mathcal{L}_i^t$ . The choice of selected layer sets has a substantial effect on training performance, which will be discussed in detail later. 

After determining the selected layer set  $\mathcal{L}_i^t$ , clients initialize the local model according to the global model sent by the server, i.e.,  $\theta_i^{t,0} = \theta^t$ , and train the local model for  $\tau$  local steps using the mini-batch SGD algorithm (McMahan et al., 2017; Wang et al., 2020; Karimireddy et al., 2020). For local step  $k \in [\tau]$ , client *i* samples a batch of data instances  $\xi_i^{t,k}$ , and calculates the gradients for the selected layers, which is given as<sup>2</sup>: 

$$\sum_{l \in \mathcal{L}_i^t} g_{i,l}(\theta_i^{t,k}; \xi_i^{t,k}) = \sum_{l \in \mathcal{L}_i^t} \nabla_l F_i(\theta_i^{t,k}; \xi_i^{t,k}).$$
(2)

Notably, the local gradient calculation pertains solely to the layers within the subset  $\mathcal{L}_i^t$ . Afterward, the local model is updated with the learning rate  $\eta$ : 

159  
160  
161  

$$\theta_{i}^{t,k} = \theta_{i}^{t,k-1} - \eta \sum_{l \in \mathcal{L}_{i}^{t}} g_{i,l}(\theta_{i}^{t,k-1};\xi_{i}^{t,k-1}), \forall k \in \{1, 2, \dots, \tau\}.$$
(3)

 $<sup>{}^{2}\</sup>nabla_{l}F(\theta)$  represents the gradient of a function  $F(\theta)$  w.r.t. the parameters of the *l*-th layer in model  $\theta$ .

Algorithm 1 Selective Layer Fine-tuning in FL	
<b>Input:</b> The pre-trained global model $\theta^0$	
for $t = 0, 1, \dots, T - 1$ do	
Sample a set of clients $S^t$ ;	
Broadcast the up-to-date global model $\theta^t$ and selected layer set $\mathcal{L}_i^t$ to clients	$\mathcal{S}^t;$
for each client $i$ in $S^t$ do	
Compute the gradients w.r.t. layers $\mathcal{L}_i^t$ and update the model for $\tau$ steps;	$\{\triangleright$ Equation (3) $\}$
Upload the accumulated updates $\Delta_i^t$ to the server;	$\{ \triangleright \text{ Equation } (4) \}$
end for	,
Compute the global update $\Delta^t$ ;	$\{\triangleright$ Equation (5) $\}$
Update the global model $\theta^t$ ;	$\{\triangleright$ Equation (6) $\}$
end for	
<b>Return:</b> The global model $\theta^T$	

The accumulated model update in local training is summarized as:

$$\Delta_{i}^{t} = \frac{1}{\eta} (\theta_{i}^{t,0} - \theta_{i}^{t,\tau}) = \sum_{k=0}^{\tau-1} \sum_{l \in \mathcal{L}_{i}^{t}} g_{i,l}(\theta_{i}^{t,k}; \xi_{i}^{t,k}).$$
(4)

After local training, clients upload their model updates  $\Delta_i^t$ ,  $i \in S^t$  to the server. The server performs federated aggregation among these model updates and optimizes the global model accordingly:

$$\Delta^t = \sum_{l \in \mathcal{L}_t} \sum_{i \in \mathcal{S}^t} w_{i,l}^t \sum_{k=0}^{\tau-1} g_{i,l}(\theta_i^{t,k}; \xi_i^{t,k}),$$
(5)

and

175 176

177 178 179

181

182

183

185

187 188

189

190

191 192 193

194 195

196 197

206

207

208

210

$$\theta^{t+1} = \theta^t - \eta \Delta^t. \tag{6}$$

Inspired by previous studies (McMahan et al., 2017; Li et al., 2020), the aggregation weights in selective layer fine-tuning are defined based on the data ratio and the masking vectors, which are formally given as:

$$w_{i,l}^{t} = \begin{cases} \frac{d_i}{\sum_{\{j \in \mathcal{S}^t \mid \mathbf{m}_j^t(l) = 1\}} d_j}, & \text{if } \mathbf{m}_i^t(l) = 1, \\ 0, & \text{otherwise.} \end{cases}$$
(7)

The details of the training process are summarized in Algorithm 1.

#### 4 WHICH LAYERS SHOULD BE SELECTED FOR FINE-TUNING?

The aforementioned training process provides substantial flexibility in selective layer fine-tuning, namely, clients are allowed to select different layers for local training and adjust their choices in 199 different training epochs. Such flexibility enables clients to tailor their local training to their data and 200 resources, providing feasible solutions for handling the heterogeneity in FL. 201

202 However, without a well-designed strategy for layer selection, the optimization of the global model 203 in FL could be severely hindered, potentially leading to a suboptimal solution or even failure in convergence. As a result, researchers have proposed several useful strategies for layer selection in 204 recent years, including: 205

- All clients select the same layer set for fine-tuning (Pillutla et al., 2022; Lee et al., 2019a; Zhang et al., 2022a; 2023; Lee et al., 2019b), i.e.,  $\mathcal{L}_i^t = \mathcal{L}_i^t, \forall i \neq j$ ;
- Clients fix their selections across different training epochs (Arivazhagan et al., 2019; Chen et al., 2023b), i.e.,  $\mathcal{L}_i^{t_1} = \mathcal{L}_i^{t_2}, \forall t_1, t_2 \in [T]$ .

211 These strategies for layer selection are proposed based on the insights drawn from experts' experience, 212 serving as special instantiations of the selected layer sets  $\mathcal{L}_i^t$ . It is worth noting that these experience-213 driven strategies might not consistently yield optimal results in various FL applications, particularly considering client heterogeneity. This leads to an essential question: How to effectively determine the 214 task-specified layer selection strategy among a large search space of possible options? In the rest of 215 this section, we provide a theoretical analysis to answer this question.

#### 216 4.1 THEORETICAL ANALYSIS 217

218 Following previous theoretical analysis in FL (Wang et al., 2020; Karimireddy et al., 2020; Li et al., 219 2020), we begin with some necessary assumptions.

220 **Assumption 4.1.** ( $\gamma$ -Smoothness) There exists a constant  $\gamma > 0$  such that for any  $\theta, \theta' \in \mathbb{R}^{P}$ , 221  $\|\nabla f_i(\theta) - \nabla f_i(\theta')\|_2 \le \gamma \|\theta - \theta'\|_2, \forall i \in \mathcal{N}.$ 222

223 For analyzing the effect of each layer on the model convergence, we give several assumptions for the 224 gradient with respect to each layer l.

225 Assumption 4.2. (Unbiased and variance-bounded stochastic gradient) The stochastic gradient 226  $g_{i,l}(\theta^t; \bar{\xi}_i^t)$  on a randomly sampled batch of data  $\xi_i^t$  is an unbiased estimate of the full-batch gradient, 227 *i.e.*,  $\mathbb{E}[g_{i,l}(\theta^t; \xi_i^t)] = \nabla_l f_i(\theta^t)$ . Besides, there exist constants  $\sigma_l > 0, \forall l \in \mathcal{L}$  such that  $||g_{i,l}(\theta^t; \xi_i^t) - \nabla_l f_i(\theta^t)| \leq 1$ . 228  $\nabla_l f_i(\theta^t) \|^2 \leq \sigma_l^2, \forall i \in \mathcal{N} \text{ and } \sum_{l \in \mathcal{L}_t} \sigma_l^2 \leq \sigma^2.$ 229

230 The non-IID data owned by clients causes diverse gradients. In the following assumption, we state the diversity of each layer's gradient. 231

232 Assumption 4.3. (Gradient diversity) There exist constants  $\kappa_l > 0, \forall l \in \mathcal{L}$  such that 233  $\left\|\nabla_{l}f(\theta^{t}) - \nabla_{l}f_{i}(\theta^{t})\right\|^{2} \leq \kappa_{l}^{2}, \forall i \in \mathcal{N}.$ 234

235 Here we first consider a case where  $\tau = 1$  to simplify the analysis without affecting the insights on 236 layer selection. The detailed analysis for the generalized case, i.e.,  $\tau > 1$ , is provided in Appendix A.3.

237 Compared with the theoretical analysis for the standard FL settings (Wang et al., 2021; Li et al., 238 2020), there exist three additional challenges in selective layer fine-tuning. Firstly, since each client 239 only updates some layers during the local training process, the aggregated gradient is no longer an 240 unbiased estimate of the local gradient  $\nabla f_i(\theta^t)$ , i.e., 241

$$\mathbb{E}[\Delta_i^t] = \sum_{l \in \mathcal{L}_i^t} \nabla_l f_i(\theta^t) \neq \nabla f_i(\theta^t), \tag{8}$$

where the inequality holds unless all layers are selected for fine-tuning, i.e.,  $\mathcal{L}_i^t = \mathcal{L}$ . Secondly, since a certain layer may not be selected by all the clients, the aggregated gradient of this layer is not equivalent to the gradient computed based on the global loss function  $(\sum_{l \in \mathcal{L}_t} \nabla_l f(\theta^t))$ , which is given as:

$$\mathbb{E}[\Delta^t] = \sum_{l \in \mathcal{L}_t} \sum_{i \in \mathcal{S}^t} w_{i,l}^t \nabla_l f_i(\theta^t) \neq \sum_{l \in \mathcal{L}_t} \nabla_l f(\theta^t),$$
(9)

where the inequality holds unless all clients select the same subset of layers. Last but not least, the aforementioned gaps vary across different training epochs, making it rather complicated in the 253 theoretical analysis.

To link the aggregated and desired gradients, we define a surrogate objective function representing the underlying loss function optimized by the clients, which is given as:

$$h_l^t(\theta) \triangleq \sum_{i \in \mathcal{S}^t} w_{i,l}^t f_i(\theta).$$
<sup>(10)</sup>

In essence, the layer-wise gradient of this objective function represents the update of the aggregated 260 global update  $\Delta^t$ . This relationship is elaborated in the following lemma.

Lemma 4.4. With Assumption 4.2, we have:

$$\mathbb{E}[\Delta^t] = \sum_{l \in \mathcal{L}_t} \nabla_l h_l^t(\theta^t), \tag{11}$$

266 where the expectation is with respect to mini-batch data sampling.

267 268

242

243 244

245

246

254

255

261

*Proof.* We rewrite both sides of Equation (11) by using the definitions in Equations (5) and (10), and 269 apply Assumption 4.2 to obtain the result. 

As aforementioned, the underlying loss function  $h_l^t(\theta)$  deviates from the desired global loss function  $f(\theta)$ , which hinders the optimization of the global model and may lead to suboptimal model performance. Such deviation can be quantified by the difference between the underlying update  $\sum_{l \in \mathcal{L}_t} \nabla_l h_l^t(\theta^t)$  and the global gradient  $\nabla f(\theta^t)$ , i.e.,

$$\mathcal{E}_t \triangleq \left\| \nabla f(\theta^t) - \sum_{l \in \mathcal{L}_t} \nabla_l h_l^t(\theta^t) \right\|^2.$$
(12)

The term  $\mathcal{E}_t$  can be further decomposed using the Jensen's inequality into two parts:

$$\mathcal{E}_t \le 2 \left\| \nabla f(\theta^t) - \sum_{l \in \mathcal{L}_t} \nabla_l f(\theta^t) \right\|_2^2 + 2 \left\| \sum_{l \in \mathcal{L}_t} \nabla_l f(\theta^t) - \nabla_l h_l^t(\theta^t) \right\|_2^2.$$
(13)

*Remark* 4.5. These two terms in the right-hand side (RHS) of (13) can be interpreted as follows: (i) The first term is the difference between the gradient w.r.t. all layers and the gradient w.r.t. the selected layers. The value of this term becomes smaller when the selected layers have *larger gradients*; (ii) The second term represents the mismatch between the desired gradient computed by all clients (i.e.,  $\nabla_l f(\theta^t) = \sum_{i \in \mathcal{N}} \alpha_i \nabla_l f_i(\theta^t)$ ) and the underlying update computed by partial clients (i.e.,  $\nabla_l h_l^t(\theta^t) = \sum_{i \in \mathcal{S}_t} w_{i,l}^t \nabla_l f_i(\theta^t)$ ), resulting from *different layer choices* among clients. If some layer is selected by all clients, its corresponding term in this term can be diminished.

For a better understanding, the following lemma shows an upper bound for the value of  $\mathcal{E}_t$ .

**Lemma 4.6.** With Assumption 4.3, we have:

$$\mathcal{E}_{t} \leq 2 \underbrace{\left[ \left\| \sum_{l \notin \mathcal{L}_{t}} \nabla f(\theta^{t}) \right\|^{2} \right]}_{\mathcal{E}_{t,1}} + 2 \underbrace{\sum_{l \in \mathcal{L}_{t}} \chi_{\mathbf{w}_{t,l} \parallel \alpha} \kappa_{l}^{2}}_{\mathcal{E}_{t,2}}, \tag{14}$$

295 296

297

305

306 307 308

310

319

320

279 280 281

283

284

285

286 287

288 289

290

291 292 293

where  $\chi_{\mathbf{w}_{t,l}\parallel\alpha} \triangleq \sum_{i \in \mathcal{N}} \frac{(w_{i,l}^t - \alpha_i)^2}{\alpha_i}$ .

<sup>298</sup> The proofs can be found in Appendix A.1.

Next we aim to analyze the impact of layer selection on the convergence of the global model. Following previous studies (Bottou et al., 2018; Wang et al., 2020), we consider an algorithm to have achieved convergence if it converges to a stationary point of the global loss function, namely, if its expected squared gradient norm  $\min_{t \in [T]} \mathbb{E} \left[ \|\nabla f(\theta^t)\|_2^2 \right]$  is zero. The following theorem and corollary show the convergence of the proposed selective layer fine-tuning framework for FL.

**Theorem 4.7.** Define a constant  $C \triangleq 1 - 4\eta L > 0$ . With Assumptions 4.1-4.3, we have:

$$\min_{t \in [T]} \mathbb{E}\left[\left\|\nabla f(\theta^{t})\right\|_{2}^{2}\right] \leq \frac{2}{\eta C T} \left[f(\theta^{0}) - f(\theta^{*})\right] + \frac{2\gamma\eta}{C} \sigma^{2} + \frac{1}{T} \sum_{t=1}^{T} \left(\frac{1}{\gamma\eta C} + 2\right) \left(\mathcal{E}_{t,1} + \mathcal{E}_{t,2}\right), \quad (15)$$

where  $\theta^* = \arg \min_{\theta \in \mathbb{R}^p} f(\theta)$  is the best model with the minimal loss.

The proofs can be found in Appendix A.2.

Corollary 4.8. With the commonly selected learning rate  $\eta = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$ , the RHS of (15) except the last term becomes zero as  $T \to \infty$ . Therefore, FL with selective layer fine-tuning may only oscillate around a stationary point of the global loss function with a non-zero error floor  $\mathcal{O}(\mathcal{E}_{t,1} + \mathcal{E}_{t,2})$ .

According to Theorem 4.7 and Corollary 4.8, the training performance of the global model is degraded by the increase of  $\mathcal{E}_{t,1} + \mathcal{E}_{t,2}$ , in the following aspects:

- The term  $\mathcal{E}_{t,1}$  indicates that it might lead to a suboptimal global model if layers with large gradient norms were not selected for fine-tuning.
- The term  $\mathcal{E}_{t,2}$  shows that the consistent selections among clients promote the convergence of the global model. Specifically, if the *l*-th layer has large gradient diversity  $\kappa_l$ , implying significant objective bias among clients, reducing the weight divergence  $\chi_{\mathbf{w}_{t,l}\parallel\alpha}$  helps to alleviate this term.

These findings highlight that the layer selection strategy that minimizes both terms can achieve better convergence of the global model. However, minimizing these two terms simultaneously may lead to contradictory selection decisions. Moreover, the optimal solution for minimizing the sum of  $\mathcal{E}_{t,1} + \mathcal{E}_{t,2}$  is inaccessible, since the ground-truth values are intractable in practice. To resolve these challenges, in the next subsection, we propose a strategy to adaptively select layers for clients and promote the learning performance of the global model.

#### 331 4.2 LAYER SELECTION STRATEGY

Based on the above analysis, we need to determine the selected layer sets  $\{\mathcal{L}_t^i\}$  that minimize the values of  $\mathcal{E}_{t,1}$  and  $\mathcal{E}_{t,2}$ . As both terms are hard to compute directly, we first design an approach to estimate their values.

To minimize  $\mathcal{E}_{t,1}$ , we prefer selecting layers with larger gradients, which can be achieved by maximizing the value of  $\sum_{l \in \mathcal{L}_t} \|\nabla_l f(\theta^t)\|_2^2$ . Since the norm of the global gradient  $\nabla_l f(\theta^t)$  is unknown, we estimate it by using the sum of stochastic local gradients, expressed as  $\sum_{i \in S^t} \|g_{i,l}(\theta^t; \xi_i^t)\|_2^2$ . Meanwhile, forcing the same layer selection among clients can reduce the value of  $\chi_{\mathbf{w}_{t,l}}\|_{\alpha}$  and thus alleviate the term  $\mathcal{E}_{t,2}$ . For this purpose, we introduce the regularization term  $\sum_{j \neq i} \|\mathbf{m}_i^t - \mathbf{m}_j^t\|_1$  into the optimization objective. Therefore, the selection of layers is determined by solving the following optimization problem:

343

330

332

344 345

346

362

363

371 372

375

376

(P1)  $\max_{\{\mathbf{m}_i^t\}} \sum_{i \in \mathcal{S}^t} \sum_{l \in \mathcal{L}_i^t} \|g_{i,l}(\theta^t; \xi_i^t)\|_2^2 - \frac{\lambda}{2} \sum_{i \in \mathcal{S}^t} \sum_{j \neq i} \|\mathbf{m}_i^t - \mathbf{m}_j^t\|_1^2,$ s.t.  $\mathcal{R}(\mathbf{m}_i^t) \leq R_i^t, \quad \forall i \in \mathcal{S}^t.$ 

Here  $\lambda \ge 0$  is a weighting constant. Specifically, a large  $\lambda$  forces consistent selection across clients, while  $\lambda = 0$  allows for independent choices among clients. The constraint in Problem (P1) ensures that the total cost of the selected layers meets the clients' local resource budgets  $R_i^t$ , and the cost function  $\mathcal{R}(\cdot)$  is typically a linear function of  $\mathbf{m}_i^t$ .

Solving Problem (P1) further gives an effective layer selection strategy for clients. At the beginning of a training epoch, each participating client  $i \in S^t$  evaluates the current global model  $\theta^t$  on a batch of local data and obtains the layer-wise gradient  $g_{i,l}(\theta^t; \xi_i^t), \forall l \in \mathcal{L}$ . Subsequently, clients upload the norms of these gradients  $||g_{i,l}(\theta^t; \xi_i^t)||_2, \forall l \in \mathcal{L}$ , which are *L*-dimensional vectors, to the server. With these values, the server can optimize the selected layer sets for clients by solving Problem (P1).

In general, the proposed layer selection strategy leads to client-specific layer sets, determined based on the estimated gradient norms. Meanwhile, a hyper-parameter  $\lambda$  is used to regulate the extent to which clients are encouraged to select the same layer. In the next section, we empirically demonstrate the benefits of the proposed strategy in effectively identifying critical layers and achieving better model performance than existing methods.

#### 4.3 DISCUSSIONS ON COMPUTATIONAL AND COMMUNICATION COSTS

In this section, we provide discussions on the computational and communication costs, considering a case where each client selects R layers to fine-tune a model with a total of L layers.

Computational costs Since both the proposed method and full model fine-tuning require the same forward operations, we focus on comparing the computational costs of backward operations among different methods. For simplicity, we assume each layer requires *b* FLOPs of backward operations. The average computational costs of the proposed layer selection method are calculated as:

$$\operatorname{Cost}_{\operatorname{ours}} = \underbrace{b(L-1)}_{\operatorname{Select}} + \underbrace{bR\tau}_{\operatorname{Fine-tune}} = b(R\tau + L - 1), \tag{16}$$

where  $\tau$  represents the local training steps. For comparison, fully fine-tuning a model requires the computational costs of:

$$\operatorname{Cost}_{\operatorname{full}} = bL\tau = \frac{L\tau}{R\tau + L - 1}\operatorname{Cost}_{\operatorname{ours}}.$$
(17)

As a result, the proposed method takes a much lower computational cost than full model fine-tuning, and the cost reduction is proportional to the number of layers and local training steps.

378 Meanwhile, the layer selection step in the proposed method introduces slightly additional costs of 379  $\frac{L-R}{\tau L}$ Cost<sub>full</sub>, which can be further reduced by evaluating the model on a smaller volume of data or making the selection decision at a lower frequency.

**Communication costs** The communication costs are determined by the transmitted bits during the training process. The proposed method only needs to transmit the selected layers, whose communication costs are much lower than those of full model fine-tuning that needs to upload the entire model. For example, assuming that different layers have the same number of parameters, the communication cost of the proposed method is  $\frac{R}{L}$  of full model fine-tuning.

To summarize, the computational and communication costs of the proposed method are much lower than those of full model fine-tuning. More empirical evidence can be found in Section 5.3.

388 389 390

387

5 EXPERIMENTS

391 392

393

5.1 Settings

Datasets & Models We conduct a series of experiments on several widely-used image classification datasets, including CIFAR-10 (Krizhevsky & Hinton, 2009) and DomainNet (Peng et al., 2019), text classification dataset, i.e., XGLUE-NC (Liang et al., 2020), and five benchmark question-answering (QA) datasets, including SCIQ (Welbl et al., 2017), OpenbookQA (Mihaylov et al., 2018), PIQA (Bisk et al., 2020), ARC-Easy and ARC-Challenge (Bhakthavatsalam et al., 2021) datasets. More details of the adopted datasets can be found in in Appendix B.1.

As for the splitting of datasets, inspired by previous studies (Zhu et al., 2021; Kim et al., 2023), 400 we consider two commonly observed data heterogeneity among clients: (i) Label skew: adopting 401 Dirichlet distribution to allocate data samples of the CIFAR-10 dataset; (ii) *Feature skew*: adopting 402 naturally domain shift on the DomainNet, XGLUE-NC, and QA datasets. Specifically, each QA 403 dataset is equally divided into two subsets, with each client possessing one subset of samples from 404 one of the five datasets. We adopt the CLIP model (Radford et al., 2021) for image classification 405 tasks and the multi-lingual XLM-Roberta-Base (Conneau et al., 2019) model on the XGLUE-NC 406 dataset. In addition, we train a LLaMA-2-7B (Touvron et al., 2023) model on the QA dataset. 407

Server & Clients In the experiments, we set up an FL system with a central server and N = 100clients. In each training epoch, the server randomly selects a subset of 20 clients, and broadcasts the up-to-date model to these clients for local training. Besides, there are N = 10 clients in the QA task and five clients are randomly chosen for training in each epoch. The resource budgets of clients are limited, which are quantified as the maximum number of layers they can fine-tune in the local training. For example, we use  $R_i = 1$  to indicate that the resource of client *i* cannot afford fine-tuning more than 1 layer. The resource budgets can be identical or heterogeneous among clients.

414 **Baselines** We compare the proposed layer selection strategy with several competitive baselines, 415 including: (i) **Top** (Kovaleva et al., 2019; Lee et al., 2019b): Clients only fine-tune the top few layers 416 (near the output) based on their task-specific data; (ii) **Bottom** (Lee et al., 2022): Clients only fine-417 tune the bottom few layers (near the input) based on their task-specific data, which can be beneficial 418 for the tasks with input-shift; (iii) Both (Xiao et al., 2023): Clients fine-tune an equal proportion 419 of both the top and bottom layers, which shows the effectiveness for large language models; (iv) 420 SNR (Mahsereci et al., 2017): Clients fine-tune the layers with higher signal-to-noise ratio (SNR) values, defined as the ratio of the mean of gradient elements to their variance; (v) RGN (Cheng et al., 421 2023; Lee et al., 2022): Clients fine-tune the layers with higher relative gradient norm (RGN) values, 422 defined as the ratio of gradient norm to the parameter norm; (vi) Moreover, we consider Full model 423 fine-tuning, i.e., training the entire model, as the performance benchmark. More implementation 424 details can be found in Appendix B.2. 425

426

5.2 Comparisons

We conduct experiments with the *identical resource* scenario and the *heterogeneous resource* scenario.

430 **Identical resource scenario** We first consider clients with identical computational resources, i.e., 431 clients select the same number of layers  $(R_i = R, \forall i \in \mathcal{N})$  for fine-tuning. The experimental results are shown in Table 1. From the model performance (accuracy) on the CIFAR-10 and DomainNet

434		CIFA	R-10	Doma	inNet	XGLU	JE-NC	Q	A
435 436		R = 1	R=2	R = 1	R=2	R = 1	R=2	R = 1	R=2
437	Full	95.	.43	90.	.27	82.	.11	65.	.98
438	Top	93.09	93.61	87.86	88.32	69.86	77.05	63.90	64.44
439	Bottom	27.38	32.81	13.80	18.63	40.43	40.60	64.18	64.60
440	Both	-	94.96	-	85.48	-	74.65	-	64.41
441	SNR	94.47	90.49	86.38	87.67	69.11	79.92	63.80	64.58
442	RGN	92.69	89.48	88.80	87.19	74.06	79.48	63.73	64.70
443	Ours	95.47	96.05	89.37	89.64	74.95	80.39	64.71	65.03
444	010	2000	2 0.000	0, 101	0,101		00109	÷, I	

432 Table 1: Test accuracy (%) on both image and text datasets, where each client selects R layers for fine-tuning. 433 The best results are highlighted in **bold**.

Table 2: Test accuracy (%) on both image and text datasets, where clients have different resources  $(R_i \in [1, 4])$ . The best results are highlighted in **bold**.

	CIFAR-10	DomainNet	XGLUE-NC	QA
Full	95.43	90.27	82.11	65.98
Тор	91.22	89.29	78.17	64.10
Bottom	27.38	23.10	50.92	64.51
Both	89.91	86.27	73.01	64.64
SNR	75.72	87.34	78.24	64.51
RGN	93.83	88.19	79.36	64.56
Ours	95.57	89.39	80.18	65.80

456 457

445

446

458 datasets, we observe that the proposed strategy demonstrates notable superiority over partial layer 459 fine-tuning baselines. Specifically, fine-tuning only one layer of the CLIP model achieves comparable 460 performance with tuning the entire model, since the CLIP model is sufficiently powerful to extract 461 useful features and thus requires less training on task-specific data. This also reveals that selective 462 layer fine-tuning well meets the performance requirement within the resources of clients.

463 On text datasets, including XGLUE-NC and QA, the proposed layer selection strategy and RGN 464 demonstrate similar performance, both surpassing other baseline methods (especially Top and Both) 465 by noticeable margins. One potential explanation for this phenomenon could be that they result in 466 similar layer selections, indicating that updating layers with higher relative gradient norms is more 467 beneficial than other strategies, which is consistent with previous study (Lee et al., 2022). Moving a 468 forward step, the proposed method adopts a flexible and dynamic layer selection strategy instead of fixed strategies, which leads to competitive performance. 469

470 Heterogeneous resource scenario Further, we conduct experiments with heterogeneous clients, i.e., 471 clients have different local resources and thus tend to select different numbers of layers for fine-tuning. 472 Such a heterogeneous resource scenario is more practical (Yang et al., 2021; Chai et al., 2019) and 473 brings additional challenges for selective layer fine-tuning. Inspired by previous studies (Wang et al., 474 2020; Nguyen et al., 2022b), the number of layers to be fine-tuned, denoted as  $R_i$  for client i, is sampled from a truncated half Normal distribution within [1, 4]. 475

476 The experimental results are shown in Table 2, from which we observe that the proposed strategy 477 consistently shows superiority over all the baseline methods on all the datasets. Compared with 478 baselines, the proposed strategy allows clients to flexibly determine the proper number of layers to 479 be tuned and effectively find the most important layers. This advantage arises from enabling clients 480 with sufficient resources to prioritize the selection of more critical layers instead of being restricted 481 to layers in fixed positions. Overall, these experimental results demonstrate the advantage of the proposed strategy when handling heterogeneity in real-world FL applications. 482

483 484

485

5.3 FURTHER DISCUSSIONS

**Visualization of selected layers** For a better understanding on the proposed layer selection strategy,



Figure 2: Visualization of selected layers  $(R_i \in [1, 4])$ .

Table 3: Comparisons of computational and communication costs when fine-tuning the CLIP model on the CIFAR-10 (R = 1). The numbers in brackets represent the costs of the proposed layer selection strategy.

	Computational cost (TFLOPs)	Ratio	Transmission (MBits)	Ratio
Full Model Fine-tuning	8.47	100%	2,811	100%
Proposed Method	2.24 (1.51)	26% (17%)	234	8.33%
Proposed (Sel. Period=2)	1.46 (0.75)	17% (9.5%)	234	8.33%
Proposed (Sel. Batch=1)	0.99 (0.30)	12% (3.4%)	234	8.33%

we visualize the selected layers on different datasets in Figure 2. When fine-tuning the CLIP model 507 on the CIFAR-10 dataset, it can be observed that the focus is primarily on updating a few top layers, 508 indicating that the low-level features (related to middle and bottom layers) are transferable from 509 pre-trained data to downstream tasks. In comparison, the DomainNet dataset, characterized by 510 a significant domain shift, necessitates extensive tuning of the middle layers in the CLIP model. 511 Furthermore, on the XGLUE-NC dataset, we observe a clear progression of selected layers for 512 fine-tuning, with a shift from the top layers progressively down to the bottom layers. Such a pattern 513 is markedly different from the trend observed in the image datasets. One possible reason lies in 514 the intrinsic differences between the modalities of text and image data. These results highlight the 515 necessity for adaptive layer selection and adjustment strategies in FL to accommodate varying dataset properties and domain shifts. 516

517 **Comparisons regarding computational and communication costs** We compare the computational 518 costs (in TFLOPs) and communication costs (in transmitted MBits) of the proposed method with 519 full model fine-tuning when adopting the CLIP model on the CIFAR-10. For the proposed method, 520 we consider fine-tuning only one layer, as it is sufficient to achieve comparable accuracy with full model fine-tuning according to Table 1. The results in Table 3 evidence a substantial decrease in 521 both computational and communication requirements when utilizing the proposed method. Besides, 522 we can observe that the layer selection strategy takes as low as 3.4% of the computational costs. 523 These experimental results demonstrate that the proposed method significantly reduces both the 524 computational and communication costs compared to full model fine-tuning. 525

- **CONCLUSIONS** 6
- 527 528

526

486

487

488

489 490

495

496

497 498

504 505 506

In this paper, we study a practical FL setting for fine-tuning foundation models, where clients are 529 allowed to optimize a subset of layers using their task-specific data. We carefully consider the 530 impact of both data heterogeneity and device heterogeneity across clients, providing a comprehensive 531 theoretical analysis of the optimization objective of selective layer fine-tuning and global model 532 convergence. The theoretical analysis offers insights into how the selected layers influence global 533 model training and highlights the role of layer importance and client heterogeneity. We further 534 propose a novel strategy for layer selection that considers the local data and available resources 535 at clients. The experimental results demonstrate that the proposed strategy outperforms baseline 536 strategies in improving the global model training performance and even matches full model fine-537 tuning performance in some scenarios, showing the potential for more efficient and tailored real-world 538 FL applications of the proposed layer selection strategy.

540 **Reproducibility statement** The assumptions and proofs of theoretical results in this work are 541 given in Section 4.1 and Appendix A.1. The experimental settings are described in Section 5.1 542 and Appendix B. The source codes are available at https://anonymous.4open.science/ 543 r/fed\_selected\_tune/. 544

### References

545

546 547

549

- Manoj Ghuhan Arivazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunav Choudhary. Feder-548 ated learning with personalization layers. arXiv preprint arXiv:1912.00818, 2019.
- Sumithra Bhakthavatsalam, Daniel Khashabi, Tushar Khot, Bhavana Dalvi Mishra, Kyle Richardson, 550 Ashish Sabharwal, Carissa Schoenick, Oyvind Tafjord, and Peter Clark. Think you have solved 551 direct-answer question answering? Try ARC-DA, the direct-answer AI2 reasoning challenge. 552 arXiv preprint arXiv:2102.03315, 2021. 553
- 554 Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. PIQA: Reasoning about 555 physical commonsense in natural language. In The Thirty-Fourth AAAI Conference on Artificial 556 Intelligence, pp. 7432–7439, 2020.
- 558 Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ B. Altman, Simran Arora, Sydney von Arx, 559 Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen 560 Creel, Jared Quincy Davis, Dorottya Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, 561 Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, 562 Lauren Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori 563 Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, 564 Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, 565 Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi, 566 and et al. On the opportunities and risks of foundation models. arXiv preprint arXiv:2108.07258, 567 2021. 568
- 569 Keith Bonawitz et al. Towards federated learning at scale: System design. In Proceedings of Machine 570 Learning and Systems, 2019.
- Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine 572 learning. SIAM Review, 60(2):223-311, Aug. 2018. 573
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in neural information processing systems (NeurIPS), pp. 1877–1901, 2020. 578
- 579 Zheng Chai, Hannan Fayyaz, Zeshan Fayyaz, Ali Anwar, Yi Zhou, Nathalie Baracaldo, Heiko 580 Ludwig, and Yue Cheng. Towards taming the resource and data heterogeneity in federated learning. In 2019 USENIX conference on operational machine learning (OpML 19), pp. 19–21, 2019. 581
  - Chaochao Chen, Xiaohua Feng, Jun Zhou, Jianwei Yin, and Xiaolin Zheng. Federated large language model: A position paper. arXiv preprint arXiv:2307.08925, 2023a.
  - Daoyuan Chen, Liuyi Yao, Dawei Gao, Bolin Ding, and Yaliang Li. Efficient personalized federated learning via sparse model-adaptation. In Proceedings of the 40th International Conference on Machine Learning (ICML), pp. 5234–5256, 2023b.
- 589 Jinyu Chen, Wenchao Xu, Song Guo, Junxiao Wang, Jie Zhang, and Haozhao Wang. FedTune: 590 A deep dive into efficient federated fine-tuning with pre-trained transformers. arXiv preprint 591 arXiv:2211.08025, 2022.
- 592

571

574

575

576

577

582

583

584 585

586

587

588

Hongrong Cheng, Miao Zhang, and Javen Qinfeng Shi. A survey on deep neural network pruningtaxonomy, comparison, analysis, and recommendations. arXiv preprint arXiv:2308.06767, 2023.

594 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam 595 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, 596 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam 597 Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James 598 Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, 600 Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. 601 Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon 602 Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark 603 Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, 604 Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways. Journal of 605 Machine Learning Research, 24(240):1–113, 2023. 606 Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Fran-607 cisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised 608 cross-lingual representation learning at scale. arXiv preprint arXiv:1911.02116, 2019. 609 610 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep 611 bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of 612 the North American Chapter of the Association for Computational Linguistics: Human Language 613 Technologies (NAACL-HLT), pp. 4171-4186, 2019. 614 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 615 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, 616 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. 617 In 9th International Conference on Learning Representations (ICLR), 2021. 618 Chen Dun, Cameron R Wolfe, Christopher M Jermaine, and Anastasios Kyrillidis. Resist: Layer-wise 619 decomposition of resnets for distributed training. In Uncertainty in Artificial Intelligence (UAI), 620 pp. 610-620, 2022. 621 622 Tyna Eloundou, Sam Manning, Pamela Mishkin, and Daniel Rock. GPTs are GPTs: An early look 623 at the labor market impact potential of large language models. arXiv preprint arXiv:2303.10130, 624 2023. 625 Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, 626 and Yu Qiao. CLIP-Adapter: Better vision-language models with feature adapters. arXiv preprint 627 arXiv:2110.04544, 2021. 628 Agrin Hilmkil, Sebastian Callh, Matteo Barbieri, Leon René Sütfeld, Edvin Listo Zec, and Olof 629 Mogren. Scaling federated learning for fine-tuning of large language models. In International 630 *Conference on Applications of Natural Language to Information Systems*, pp. 15–23, 2021. 631 632 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea 633 Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for NLP. In 634 Proceedings of the 36th International Conference on Machine Learning (ICML), pp. 2790–2799, 635 2019. 636 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 637 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In The Tenth Interna-638 tional Conference on Learning Representations (ICLR), 2022. 639 640 Ahmed Imteaj, Urmish Thakker, Shiqiang Wang, Jian Li, and M. Hadi Amini. A survey on federated 641 learning for resource-constrained IoT devices. *IEEE Internet of Things Journal*, 9(1):1–24, 2022. 642 Gal Kaplun, Andrey Gurevich, Tal Swisa, Mazor David, Shai Shalev-Shwartz, and Eran Malach. 643 Less is more: Selective layer finetuning with subtuning. arXiv preprint arXiv:2302.06354, 2023. 644 645 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In 646 Proceedings of the 37th International conference on machine learning (ICML), pp. 5132–5143, 647 2020.

648 649 650 651	Yeachan Kim, Junho Kim, Wing-Lam Mok, Jun-Hyung Park, and SangKeun Lee. Client-customized adaptation for parameter-efficient federated learning. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (eds.), <i>Findings of the Association for Computational Linguistics</i> , pp. 1159–1172, 2023.
652 653 654	Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. Revealing the dark secrets of bert. <i>arXiv preprint arXiv:1908.08593</i> , 2019.
655 656	Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. [Online]. Available: https://www.cs.toronto.edu/~kriz/cifar.html, 2009.
657 658 659 660	Weirui Kuang, Bingchen Qian, Zitao Li, Daoyuan Chen, Dawei Gao, Xuchen Pan, Yuexiang Xie, Yaliang Li, Bolin Ding, and Jingren Zhou. Federatedscope-llm: A comprehensive package for fine-tuning large language models in federated learning. arXiv preprint arXiv:2309.00363, 2023.
661 662 663	Cheolhyoung Lee, Kyunghyun Cho, and Wanmo Kang. Mixout: Effective regularization to finetune large-scale pretrained language models. In <i>International Conference on Learning Representations (ICLR)</i> , 2019a.
664 665	Jaejun Lee, Raphael Tang, and Jimmy Lin. What would elsa do? Freezing layers during transformer fine-tuning. <i>arXiv preprint arXiv:1911.03090</i> , 2019b.
667 668 669	Sunwoo Lee, Tuo Zhang, and A Salman Avestimehr. Layer-wise adaptive model aggregation for scalable federated learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)</i> , pp. 8491–8499, 2023.
670 671 672	Yoonho Lee, Annie S Chen, Fahim Tajwar, Ananya Kumar, Huaxiu Yao, Percy Liang, and Chelsea Finn. Surgical fine-tuning improves adaptation to distribution shifts. In <i>International Conference on Learning Representations (ICLR)</i> , 2022.
673 674 675 676	Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. Federated learning on non-iid data silos: An experimental study. In <i>38th IEEE International Conference on Data Engineering (ICDE)</i> , pp. 965–978, 2022.
677 678 679	Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. In 8th International Conference on Learning Representations (ICLR), Addis Ababa, Ethiopia, 2020.
680 681 682	Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing</i> , pp. 4582–4597, 2021.
683 684 685 686	Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. FedBN: Federated learning on non-iid features via local batch normalization. In <i>9th International Conference on Learning Representations (ICLR)</i> , 2021.
687 688 689 690 691 692	Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun Shou, Daxin Jiang, Guihong Cao, Xiaodong Fan, Ruofei Zhang, Rahul Agrawal, Edward Cui, Sining Wei, Taroon Bharti, Ying Qiao, Jiun-Hung Chen, Winnie Wu, Shuguang Liu, Fan Yang, Daniel Campos, Rangan Majumder, and Ming Zhou. XGLUE: A new benchmark datasetfor cross-lingual pre-training, understanding and generation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pp. 6008–6018, 2020.
693 694	Maren Mahsereci, Lukas Balles, Christoph Lassner, and Philipp Hennig. Early stopping without a validation set. <i>arXiv preprint arXiv:1703.09580</i> , 2017.
695 696 697 698	Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In <i>Proceedings of the</i> 20th International Conference on Artificial Intelligence and Statistics (AISTATS), 2017.
699 700 701	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? A new dataset for open book question answering. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pp. 2381–2391, 2018.

702 703 704	John Nguyen, Kshitiz Malik, Maziar Sanjabi, and Michael Rabbat. Where to begin? Exploring the impact of pre-training and initialization in federated learning. <i>arXiv preprint arXiv:2206.15387</i> , 2022a.
705 706 707 708	John Nguyen, Kshitiz Malik, Hongyuan Zhan, Ashkan Yousefpour, Mike Rabbat, Mani Malek, and Dzmitry Huba. Federated learning with buffered asynchronous aggregation. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera (eds.), <i>International Conference on Artificial</i> <i>Intelligence and Statistics (AISTATS)</i> , pp. 3581–3607, 2022b.
709 710 711 712	Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In 2019 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 1406–1415, 2019.
713 714 715 716	Krishna Pillutla, Kshitiz Malik, Abdel-Rahman Mohamed, Mike Rabbat, Maziar Sanjabi, and Lin Xiao. Federated learning with partial model personalization. In <i>Proceedings of the 39th International Conference on Machine Learning (ICML)</i> , pp. 17716–17758. PMLR, 2022.
717 718	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9, 2019.
719 720 721 722 723	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In <i>Proceedings of the 38th International Conference on Machine Learning (ICML)</i> , pp. 8748–8763, 2021.
724 725 726	Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In <i>Proceedings of the 38th International Conference on Machine Learning (ICML)</i> , pp. 8821–8831, 2021.
727 728 720	Zhiqiang Shen, Zechun Liu, Jie Qin, Marios Savvides, and Kwang-Ting Cheng. Partial is better than all: Revisiting fine-tuning strategy for few-shot learning. <i>arXiv preprint arXiv:2102.03983</i> , 2021.
729 730 731 732	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> , 2023.
733 734 735 736	Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H. Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2020.
737 738 739 740	Jianyu Wang, Zachary Charles, Zheng Xu, Gauri Joshi, H Brendan McMahan, Maruan Al-Shedivat, Galen Andrew, Salman Avestimehr, Katharine Daly, Deepesh Data, et al. A field guide to federated optimization. arXiv preprint arXiv:2107.06917, 2021.
741 742 743	Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. In Leon Derczynski, Wei Xu, Alan Ritter, and Tim Baldwin (eds.), <i>Proceedings of the 3rd Workshop</i> on Noisy User-generated Text (NUT@EMNLP), pp. 94–106, 2017.
744 745 746	Guangxuan Xiao, Ji Lin, and Song Han. Offsite-tuning: Transfer learning without full model. <i>arXiv</i> preprint arXiv:2302.04870, 2023.
747 748 749	Runxin Xu, Fuli Luo, Zhiyuan Zhang, Chuanqi Tan, Baobao Chang, Songfang Huang, and Fei Huang. Raise a child in large language model: Towards effective and generalizable fine-tuning. <i>arXiv</i> preprint arXiv:2109.05687, 2021.
750 751 752	Chengxu Yang, Qipeng Wang, Mengwei Xu, Zhenpeng Chen, Kaigui Bian, Yunxin Liu, and Xuanzhe Liu. Characterizing impacts of heterogeneity in federated learning upon large-scale smartphone data. In <i>Proceedings of the Web Conference 2021</i> , pp. 935–946, 2021.
754 755	Elad Ben Zaken, Yoav Goldberg, and Shauli Ravfogel. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pp. 1–9, 2022.

756 757 758 750	Haojie Zhang, Ge Li, Jia Li, Zhongjin Zhang, Yuqi Zhu, and Zhi Jin. Fine-tuning pre-trained language models effectively by optimizing subnetworks adaptively. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , pp. 21442–21454, 2022a.
760	Kaiyan Zhang, Ning Ding, Biging Oi, Xuekai Zhu, Xinwei Long, and Bowen Zhou. CRaSh:
760	Clustering, removing, and sharing enhance fine-tuning without full large language model. <i>arXiv</i>
761	preprint arXiv:2310.15477, 2023.
763	Lin Zhang Li Shen Liang Ding Dacheng Tao and Ling-Yu Duan Fine-tuning global model via
764	data-free knowledge distillation for non-iid federated learning. In <i>Proceedings of the IEEE/CVF</i>
765	conference on computer vision and pattern recognition (CVPR), pp. 10174–10183, 2022b.
766	
767	Hangyu Zhu, Jinjin Xu, Shiqing Liu, and Yaochu Jin. Federated learning on non-IID data: A survey. <i>Neurocomputing</i> , 465:371–390, 2021.
768	
769	
770	
771	
772	
773	
774	
775	
776	
777	
778	
779	
780	
781	
782	
703	
704	
786	
787	
788	
789	
790	
791	
792	
793	
794	
795	
796	
797	
798	
799	
800	
801	
802	
803	
804	
805	
806	
807	
808	
809	

## A CONVERGENCE ANALYSIS: FULL PROOFS

### A.1 USEFUL LEMMAS

A.1.1 ONE-ROUND LOSS DECAY

**Lemma A.1.** With Assumption 4.1, we have:

$$\mathbb{E}[f(\theta^{t+1})] - \mathbb{E}[f(\theta^{t})] \leq \frac{1}{2\gamma} \mathcal{E}_t + \underbrace{\mathbb{E}\left\langle \sum_{l \in \mathcal{L}_t} \nabla_l h_l^t(\theta^t), \theta^{t+1} - \theta^t \right\rangle}_{\mathcal{T}_1} + \gamma \underbrace{\mathbb{E}\left[ \left\| \theta^{t+1} - \theta^t \right\|^2 \right]}_{\mathcal{T}_2}.$$
 (18)

*Proof.* We begin with analyzing the loss decay by using  $\gamma$ -smoothness in Assumption 4.1 as follows:  $\mathbb{E}[f(\theta^{t+1})] - \mathbb{E}[f(\theta^{t})]$ (19)

$$\leq \mathbb{E}\langle \nabla f(\theta^{t}), \theta^{t+1} - \theta^{t} \rangle + \frac{\gamma}{2} \mathbb{E}[\left\| \theta^{t+1} - \theta^{t} \right\|^{2}]$$
(20)

$$= \mathbb{E}\left\langle \nabla f(\theta^{t}) - \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) + \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}), \theta^{t+1} - \theta^{t} \right\rangle + \frac{\gamma}{2} \mathbb{E}[\left\|\theta^{t+1} - \theta^{t}\right\|^{2}]$$
(21)

$$=\underbrace{\mathbb{E}\left\langle \nabla f(\theta^{t}) - \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}), \theta^{t+1} - \theta^{t} \right\rangle}_{\mathcal{T}_{0}} + \underbrace{\mathbb{E}\left\langle \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}), \theta^{t+1} - \theta^{t} \right\rangle}_{\mathcal{T}_{1}} + \underbrace{\mathbb{E}\left[ \left\| \theta^{t+1} - \theta^{t} \right\|^{2} \right]}_{\mathcal{T}_{2}} + \underbrace{\mathbb{E}\left[ \left\| \theta^{t} + \theta^{t} \right\|^{2} \right]}_{\mathcal{T}_{2}} + \underbrace{$$

By Young's inequality, we upper bound the term  $T_0$  as:

$$\mathcal{T}_{0} = \mathbb{E}\left\langle \nabla f(\theta^{t}) - \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}), \theta^{t+1} - \theta^{t} \right\rangle$$
(23)

$$\leq \frac{1}{2\gamma} \underbrace{\mathbb{E}\left[ \left\| \nabla f(\theta^{t}) - \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) \right\|^{2} \right]}_{\mathcal{E}_{t}} + \frac{\gamma}{2} \mathbb{E}\left[ \left\| \theta^{t+1} - \theta^{t} \right\|^{2} \right]$$
(24)

$$=\frac{1}{2\gamma}\mathcal{E}_t + \frac{\gamma}{2}\mathcal{T}_2.$$
(25)

Plugging (25) back into (22) gives the result in (18).

### A.1.2 ANALYZING $\mathcal{E}_t$ : PROOF OF LEMMA 4.6

In this subsection, we prove the result in Lemma 4.6.

We begin with decomposing the term  $\mathcal{E}_t$  using the Jensen's inequality as:

$$\mathcal{E}_{t} \leq 2 \underbrace{\|\nabla f(\theta^{t}) - \sum_{l \in \mathcal{L}_{t}} \nabla_{l} f(\theta^{t})\|_{2}^{2} + 2}_{\tilde{\mathcal{E}}_{t,1}} \underbrace{\|\sum_{l \in \mathcal{L}_{t}} \nabla_{l} f(\theta^{t}) - \nabla_{l} h_{l}^{t}(\theta^{t})\|_{2}^{2}}_{\tilde{\mathcal{E}}_{t,2}}.$$
(26)

For the first term  $\tilde{\mathcal{E}}_{1,t}$ , we directly obtain:

$$\tilde{\mathcal{E}}_{1,t} = \mathbb{E}\left[\left\|\nabla f(\theta^t) - \sum_{l \in \mathcal{L}_t} \nabla_l f(\theta^t)\right\|^2\right] = \mathbb{E}\left[\left\|\sum_{l \notin \mathcal{L}_t} \nabla_l f(\theta^t)\right\|^2\right].$$
(27)

Afterwards, we derive the value of  $\tilde{\mathcal{E}}_{2,t}$  as follows: 

$$\tilde{\mathcal{E}}_{2,t} = \mathbb{E}\left[ \left\| \sum_{l \in \mathcal{L}_t} \nabla_l f(\theta^t) - \sum_{l \in \mathcal{L}_t} \nabla_l h_l^t(\theta^t) \right\|^2 \right]$$
(28)

$$= \mathbb{E}\left[\left\|\sum_{l\in\mathcal{L}_{t}}\sum_{i\in\mathcal{N}}\alpha_{i}\nabla_{l}f_{i}(\theta^{t}) - \sum_{l\in\mathcal{L}_{t}}\sum_{i\in\mathcal{S}_{t}}w_{i,l}^{t}\nabla_{l}f_{i}(\theta^{t})\right\|^{2}\right]$$
(29)

$$= \mathbb{E}\left[ \left\| \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} \alpha_{i} \nabla_{l} f_{i}(\theta^{t}) - \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} w_{i,l}^{t} \nabla_{l} f_{i}(\theta^{t}) \right\|^{2} \right]$$
(30)

$$= \sum_{l \in \mathcal{L}_{t}} \mathbb{E}\left[ \left\| \sum_{i \in \mathcal{N}} \frac{w_{i,l}^{t} - \alpha_{i}}{\sqrt{\alpha_{i}}} \sqrt{\alpha_{i}} \left( \nabla_{l} f_{i}(\theta^{t}) - \nabla_{l} f(\theta^{t}) \right) \right\|^{2} \right]$$
(31)

$$\leq \sum_{l \in \mathcal{L}_t} \left[ \sum_{i \in \mathcal{N}} \frac{(w_{i,l}^t - \alpha_i)^2}{\alpha_i} \right] \left[ \sum_{i \in \mathcal{N}} \alpha_i \mathbb{E} \left[ \left\| \nabla_l f_i(\theta^t) - \nabla_l f(\theta^t) \right\|^2 \right] \right]$$
(32)

$$\leq \sum_{l \in \mathcal{L}_t} \chi_{\mathbf{w}_{t,l} \parallel \alpha} \kappa_l^2.$$
(33)

where (32) follows the Cauchy-Schwarz inequality and (33) applies Assumption 4.3.

By substituting the RHS of (27) and (33) into (26), we complete the proof. 

### A.2 CONVERGENCE ANALYSIS FOR SINGLE-STEP CASE

In this subsection, we consider  $\tau = 1$  and prove Theorem 4.7.

We derive the value of  $\mathcal{T}_1$  as follows:

$$\mathcal{T}_{1} = \mathbb{E}\left\langle \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}), -\eta \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) \right\rangle = -\eta \mathbb{E}\left[ \left\| \sum_{l \in L_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) \right\|^{2} \right].$$
(34)

Afterwards, we give an upper bound for the term  $T_2$  as follows:

$$\mathcal{T}_{2} = \left[ \left\| \eta \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{S}^{t}} w_{i,l}^{t} g_{i,l}(\theta^{t}; \xi_{i}^{t}) \right\|^{2} \right]$$
(35)

$$\leq \eta^{2} \mathbb{E} \left[ \left\| \sum_{l \in L_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) \right\|^{2} \right] + \eta^{2} \sigma^{2}, \tag{36}$$

where (36) follows Assumption 4.2.

Using the result in Lemma A.1, we have:  $m \left[ e(at+1) \right]$  $\pi r c(at)$ 

$$\mathbb{E}[f(\theta^{t+1})] - \mathbb{E}[f(\theta^{t})] \\ \leq \frac{1}{2\gamma} \mathcal{E}_t - \eta \mathbb{E}\left[\left\|\sum_{l \in L_t} \nabla_l h_l^t(\theta^t)\right\|^2\right] + \gamma \left[\eta^2 \mathbb{E}\left[\left\|\sum_{l \in L_t} \nabla_l h_l^t(\theta^t)\right\|^2\right] + \eta^2 \sigma^2\right]$$
(37)

$$= \frac{1}{2\gamma} \mathcal{E}_t - \eta (1 - \gamma \eta) \mathbb{E} \left[ \left\| \sum_{l \in L_t} \nabla_l h_l^t(\theta^t) \right\|^2 \right] + \gamma \eta^2 \sigma^2.$$
(38)

We define a constant  $C \triangleq 1 - \gamma \eta > 0$  and arrange the terms in (38) as follows:

916  
917 
$$\mathbb{E}\left[\left\|\sum_{l\in L_t} \nabla_l h_l^t(\theta^t)\right\|^2\right] \le \frac{1}{\eta C} \left[\mathbb{E}[f(\theta^t)] - \mathbb{E}[f(\theta^{t+1})]\right] + \frac{1}{2\gamma \eta C} \mathcal{E}_t + \frac{\gamma \eta}{C} \sigma^2.$$
(39)

918 By Jensen's inequality, we have: 

$$\mathbb{E}\left[\left\|\nabla f(\theta^{t})\right\|^{2}\right] \tag{40}$$

$$= \mathbb{E}\left[ \left\| \nabla f(\theta^{t}) - \sum_{l \in L_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) + \sum_{l \in L_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) \right\|^{2} \right]$$
(41)

$$\leq 2\mathbb{E}\left[\left\|\nabla f(\theta^{t}) - \sum_{l \in L_{t}} \nabla_{l} h_{l}^{t}(\theta^{t})\right\|^{2}\right] + 2\mathbb{E}\left[\left\|\sum_{l \in L_{t}} \nabla_{l} h_{l}^{t}(\theta^{t})\right\|^{2}\right]$$
(42)

$$= 2\mathcal{E}_t + 2\mathbb{E}\left[ \left\| \sum_{l \in L_t} \nabla_l h_l^t(\theta^t) \right\|^2 \right].$$
(43)

Combining (39) and (43) gives:

$$\mathbb{E}\left[\left\|\nabla f(\theta^{t})\right\|^{2}\right] \leq \frac{2}{\eta C} \left[\mathbb{E}[f(\theta^{t})] - \mathbb{E}[f(\theta^{t+1})]\right] + \left(\frac{1}{\gamma \eta C} + 2\right) \mathcal{E}_{t} + \frac{2\gamma \eta}{C} \sigma^{2}.$$
 (44)

We sum up both sides of (44) over t = 0, 1, ..., T - 1 and divide them by T to obtain the following result:

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E} \left[ \left\| \nabla f(\theta^{t}) \right\|^{2} \right]$$

$$\leq \frac{2}{\eta CT} \left[ \mathbb{E} [f(\theta^{0})] - \mathbb{E} [f(\theta^{T})] \right] + \frac{1}{T} \sum_{t=1}^{T} \left( \frac{1}{\gamma \eta C} + 2 \right) \mathcal{E}_{t} + \frac{2\gamma \eta}{C} \sigma^{2}$$
(45)

$$\leq \frac{2}{\eta CT} \left[ f(\theta^0) - f(\theta^*) \right] + \frac{1}{T} \sum_{t=1}^T \left( \frac{1}{\gamma \eta C} + 2 \right) \mathcal{E}_t + \frac{2\gamma \eta}{C} \sigma^2 \tag{46}$$

$$\leq \frac{2}{\eta CT} \left[ f(\theta^0) - f(\theta^*) \right] + \frac{2\gamma\eta}{C} \sigma^2 + \frac{1}{T} \sum_{t=1}^T \left( \frac{1}{\gamma\eta C} + 2 \right) \left( \mathcal{E}_{t,1} + \mathcal{E}_{t,2} \right). \tag{47}$$

### A.3 CONVERGENCE ANALYSIS FOR MULTI-STEP CASE

Consider the general case where  $\tau > 1$ . We characterize the convergence in the following theorem and note that the impact of  $\mathcal{E}_{t,1} + \mathcal{E}_{t,2}$  is similar to that in Theorem 4.7.

Theorem A.2. Let 
$$C' \triangleq 1 - 4\eta\tau - 8\eta^2\gamma^2\tau(\tau-1) - 32\eta^3\gamma^2\tau^2(\tau-1) > 0$$
 and  $A_{\tau} \triangleq \eta + 2\eta^2\gamma^2\tau(\tau-1) + 8\eta^3\gamma^2\tau^2(\tau-1)$ . With Assumptions 4.1-4.3, we have:  

$$\frac{1}{T}\sum_{\tau}^{T}\mathbb{E}\left[\left\|\nabla f(\theta^t)\right\|^2\right] \leq \frac{2}{\eta\tau C'T}\left[f(\theta^0) - f(\theta^*)\right] + \frac{4A_{\tau}}{C'}\sigma^2 + \frac{1}{T}\sum_{\tau}^{T}\left(\frac{1}{\eta\tau\gamma C'} + 2\right)\left(\mathcal{E}_{t,1} + \mathcal{E}_{t,2}\right)$$

$$\overline{T}\sum_{t=1}^{\infty} \mathbb{E}\left[\left\|\nabla f(\theta^{*})\right\|\right] \leq \frac{1}{\eta\tau C'T} \left[f(\theta^{*}) - f(\theta^{*})\right] + \frac{1}{C'}\sigma^{2} + \frac{1}{T}\sum_{t=1}^{\infty} \left(\frac{\eta\tau\gamma C'}{\eta\tau\gamma C'} + 2\right)(\mathcal{E}_{t,1} + \mathcal{E}_{t,2}).$$
(48)

*Proof.* In Lemma A.1, the term  $T_1$  is related to client drift caused by multiple local SGD steps, which can be upper bounded as follows:

 $\mathcal{T}_1$ 

$$= -\eta \sum_{k=0}^{\tau-1} \mathbb{E} \left\langle \sum_{l \in \mathcal{L}_t} \nabla_l h_l^t(\theta^t), \sum_{l \in \mathcal{L}_t} \sum_{i \in \mathcal{N}} w_{i,l}^t \nabla_l f_i(\theta_i^{t,k}) \right\rangle$$

$$= -\eta \sum_{k=0}^{\tau-1} \mathbb{E} \left\langle \sum_{l \in \mathcal{L}_t} \nabla_l h_l^t(\theta^t), \sum_{l \in \mathcal{L}_t} \nabla_l h_l^t(\theta^t) \right\rangle$$
(49)

$$+ \eta \sum_{k=0}^{\tau-1} \mathbb{E} \left\langle \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}), \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) - \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} w_{i,l}^{t} \nabla_{l} f_{i}(\theta_{i}^{t,k}) \right\rangle$$

$$\leq - \frac{\eta}{2} \sum_{k=0}^{\tau-1} \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) \right\|^{2} \right] + \frac{\eta}{2} \sum_{k=0}^{\tau-1} \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) - \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} w_{i,l}^{t} \nabla_{l} f_{i}(\theta_{i}^{t,k}) \right\|^{2} \right]$$

$$= - \frac{\eta\tau}{2} \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) \right\|^{2} \right] + \frac{\eta}{2} \sum_{k=0}^{\tau-1} \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} w_{i,l}^{t} \nabla_{l} f_{i}(\theta^{t}) - \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} w_{i,l}^{t} \nabla_{l} f_{i}(\theta_{i}^{t,k}) \right\|^{2} \right]$$

$$(51)$$

$$= - \frac{\eta\tau}{2} \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) \right\|^{2} \right] + \frac{\eta}{2} \sum_{k=0}^{\tau-1} \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} w_{i,l}^{t} \nabla_{l} f_{i}(\theta^{t}) - \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} w_{i,l}^{t} \nabla_{l} f_{i}(\theta_{i}^{t,k}) \right\|^{2} \right]$$

$$(52)$$

$$\leq -\frac{\eta\tau}{2}\mathbb{E}\left[\left\|\sum_{l\in\mathcal{L}_{t}}\nabla_{l}h_{l}^{t}(\theta^{t})\right\|^{2}\right] + \frac{\eta\gamma^{2}}{2}\underbrace{\sum_{k=0}^{\tau-1}\mathbb{E}\left[\left\|\sum_{l\in\mathcal{L}_{t}}\sum_{i\in\mathcal{N}}w_{i,l}^{t}\left(\theta^{t}-\theta_{i}^{t,k}\right)\right\|^{2}\right]}_{\mathcal{T}_{4}},\tag{53}$$

where (51) uses the inequality  $\langle \mathbf{a}, \mathbf{b} \rangle \leq \frac{\|\mathbf{a}\|^2}{2} + \frac{\|\mathbf{b}\|^2}{2}$ , and (53) follows Assumption 4.1. Then we analyze the term  $\mathcal{T}_2$  as follows:  $\mathcal{T}_2$ 

$$\leq \eta^{2} \mathbb{E} \left[ \left\| \sum_{k=0}^{\tau-1} \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} w_{i,l}^{t} \nabla_{l} f_{i}(\theta_{i}^{t,k}) \right\|^{2} \right] + \eta^{2} \tau \sigma^{2}$$

$$\leq \eta^{2} \tau \sum_{k=0}^{\tau-1} \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} w_{i,l}^{t} \nabla_{l} f_{i}(\theta_{i}^{t,k}) - \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} w_{i,l}^{t} \nabla_{l} f_{i}(\theta^{t}) + \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} w_{i,l}^{t} \nabla_{l} f_{i}(\theta^{t}) \right\|^{2} \right]$$

$$+ \eta^{2} \tau \sigma^{2}$$

$$(54)$$

$$\leq 2\eta^{2}\tau \sum_{k=0}^{\tau-1} \mathbb{E}\left[\left\|\sum_{l\in\mathcal{L}_{t}}\sum_{i\in\mathcal{N}}w_{i,l}^{t}\nabla_{l}f_{i}(\theta_{i}^{t,k}) - \sum_{l\in\mathcal{L}_{t}}\sum_{i\in\mathcal{N}}w_{i,l}^{t}\nabla_{l}f_{i}(\theta^{t})\right\|^{2}\right] + 2\eta^{2}\tau \sum_{k=0}^{\tau-1} \mathbb{E}\left[\left\|\sum_{l\in\mathcal{L}_{t}}\nabla_{l}h_{l}^{t}(\theta^{t})\right\|^{2} + \eta^{2}\tau\sigma^{2}\right]$$
(56)

 $\frac{\overline{k=0}}{\leq 2\eta^2 \gamma^2 \tau} \sum_{k=0}^{\tau-1} \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_t} \sum_{i \in \mathcal{N}} w_{i,l}^t \left( \theta_i^{t,k} - \theta^t \right) \right\|^2 \right] + 2\eta^2 \tau^2 \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_t} \nabla_l h_l^t(\theta^t) \right\|^2 \right] + \eta^2 \tau \sigma^2 \quad (57)$ 

$$=2\eta^{2}\gamma^{2}\tau\mathcal{T}_{4}+2\eta^{2}\tau^{2}\mathbb{E}\left[\left\|\sum_{l\in\mathcal{L}_{t}}\nabla_{l}h_{l}^{t}(\theta^{t})\right\|^{2}\right]+\eta^{2}\tau\sigma^{2},$$
(58)

where (54) follows Assumption 4.2, (55)-(56) apply the Jensen's inequality, and (57) follows Assumption 4.1.

Following Lemma 22 in (Pillutla et al., 2022),  $T_4$  can be upper bounded as: 

$$\mathcal{T}_{4} \leq \sum_{k=0}^{\tau-1} \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_{t}} \sum_{i \in \mathcal{N}} w_{i,l}^{t} \theta^{t} - \theta_{i}^{t,k} \right\|^{2} \right]$$
(59)

$$\leq 8\eta^{2}\tau^{2}(\tau-1)\mathbb{E}\left[\left\|\sum_{l\in\mathcal{L}_{t}}\sum_{i\in\mathcal{N}}w_{i,l}^{t}\nabla_{l}f_{i}(\theta^{t})\right\|^{2}\right] + \sum_{l\in\mathcal{L}_{t}}\sum_{i\in\mathcal{N}}w_{i,l}^{t}4\eta^{2}\tau^{2}(\tau-1)\sigma_{l}^{2} \qquad (60)$$

1024  
1025 
$$= 8\eta^2 \tau^2 (\tau - 1) \mathbb{E}$$

$$=8\eta^{2}\tau^{2}(\tau-1)\mathbb{E}\left[\left\|\sum_{l\in\mathcal{L}_{t}}\nabla_{l}h_{l}^{t}(\theta^{t})\right\|^{2}\right]+4\eta^{2}\tau^{2}(\tau-1)\sigma^{2}.$$
(61)

Plugging (54),(58) and (61) back into (18), we have the following result:

$$\begin{aligned} & \mathbb{E}[f(\theta^{t+1})] - \mathbb{E}[f(\theta^{t})] & \qquad (62) \\ & 1029 \\ & 1030 \\ & 1031 \\ & 1032 \end{aligned} \leq \frac{1}{2\gamma} \mathcal{E}_t - \frac{\eta\tau}{2} \mathbb{E}\left[ \left\| \sum_{l \in \mathcal{L}_t} \nabla_l h_l^t(\theta^t) \right\|^2 \right] + \frac{\eta\gamma^2}{2} \mathcal{T}_4 + 2\eta^2 \gamma^2 \tau \mathcal{T}_4 + 2\eta^2 \tau^2 \mathbb{E}\left[ \left\| \sum_{l \in \mathcal{L}_t} \nabla_l h_l^t(\theta^t) \right\|^2 \right] + \eta^2 \tau \sigma^2 \end{aligned}$$

 $((\alpha))$ 

$$=\frac{1}{2\gamma}\mathcal{E}_{t} - \frac{\eta\tau}{2}\left(1 - 4\eta\tau\right)\mathbb{E}\left[\left\|\sum_{l\in\mathcal{L}_{t}}\nabla_{l}h_{l}^{t}(\theta^{t})\right\|^{2}\right] + \eta^{2}\tau\sigma^{2} + \left(\frac{\eta\gamma^{2}}{2} + 2\eta^{2}\gamma^{2}\tau\right)\mathcal{T}_{4}$$
(64)

$$= \frac{1}{2\gamma} \mathcal{E}_{t} - \frac{\eta\tau}{2} (1 - 4\eta\tau) \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) \right\|^{2} \right] + \eta^{2} \tau \sigma^{2} \\ + \left( \frac{\eta\gamma^{2}}{2} + 2\eta^{2} \gamma^{2} \tau \right) \left\{ 8\eta^{2} \tau^{2} (\tau - 1) \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) \right\|^{2} \right] + 4\eta^{2} \tau^{2} (\tau - 1) \sigma^{2} \right\}$$

$$= \frac{1}{2} \mathcal{E}_{t} - \frac{\eta\tau}{2} \left[ 1 - 4\eta\tau - 8\eta^{2} \gamma^{2} \tau (\tau - 1) - 32\eta^{3} \gamma^{2} \tau^{2} (\tau - 1) \right] \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) \right\|^{2} \right]$$

$$(65)$$

$$= \frac{1}{2\gamma} \mathcal{E}_{t} - \frac{\eta\tau}{2} \left[ 1 - 4\eta\tau - 8\eta^{2}\gamma^{2}\tau(\tau-1) - 32\eta^{3}\gamma^{2}\tau^{2}(\tau-1) \right] \mathbb{E} \left[ \left\| \sum_{l \in \mathcal{L}_{t}} \nabla_{l} h_{l}^{t}(\theta^{t}) \right\| \right] + \left( \eta^{2}\tau + 2\eta^{3}\gamma^{2}\tau^{2}(\tau-1) + 8\eta^{4}\gamma^{2}\tau^{3}(\tau-1) \right) \sigma^{2}.$$
(66)

Let  $C' \triangleq 1 - 4\eta\tau - 8\eta^2\gamma^2\tau(\tau-1) - 32\eta^3\gamma^2\tau^2(\tau-1) > 0$  and  $A_\tau \triangleq \eta + 2\eta^2\gamma^2\tau(\tau-1) + 2\eta^2\tau(\tau-1) +$  $8\eta^3\gamma^2\tau^2(\tau-1)$ . We have:

$$\mathbb{E}\left[\left\|\sum_{l\in\mathcal{L}_t}\nabla_l h_l^t(\theta^t)\right\|^2\right] \le \frac{2}{\eta\tau C'} \left[\mathbb{E}[f(\theta^t)] - \mathbb{E}[f(\theta^{t+1})]\right] + \frac{1}{\eta\tau\gamma C'}\mathcal{E}_t + \frac{2}{C'}A_\tau\sigma^2.$$
(67)

Using the result in (43), we have:

$$\mathbb{E}\left[\left\|\nabla f(\theta^{t})\right\|^{2}\right] \leq \frac{4}{\eta\tau C'} \left[\mathbb{E}[f(\theta^{t})] - \mathbb{E}[f(\theta^{t+1})]\right] + \left(\frac{1}{\eta\tau\gamma C'} + 2\right)\mathcal{E}_{t} + \frac{4A_{\tau}}{C'}\sigma^{2}.$$
 (68)

We sum up both sides of (68) over t = 0, 1, ..., T - 1 and divide them by T to obtain the following result: 

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E} \left[ \left\| \nabla f(\theta^{t}) \right\|^{2} \right]$$

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E} \left[ \left\| \nabla f(\theta^{t}) \right\|^{2} \right]$$

$$\leq \frac{2}{\eta \tau C'T} \left[ \mathbb{E} [f(\theta^{0})] - \mathbb{E} [f(\theta^{T})] \right] + \frac{1}{T} \sum_{t=1}^{T} \left( \frac{1}{\eta \tau \gamma C'} + 2 \right) \mathcal{E}_{t} + \frac{4A_{\tau}}{C'} \sigma^{2}$$
(69)

$$\leq \frac{2}{\eta \tau C'T} \left[ f(\theta^0) - f(\theta^*) \right] + \frac{1}{T} \sum_{t=1}^T \left( \frac{1}{\eta \tau \gamma C'} + 2 \right) \mathcal{E}_t + \frac{4A_\tau}{C'} \sigma^2 \tag{70}$$

$$\leq \frac{2}{\eta \tau C' T} \left[ f(\theta^0) - f(\theta^*) \right] + \frac{4A_\tau}{C'} \sigma^2 + \frac{1}{T} \sum_{t=1}^T \left( \frac{1}{\eta \tau \gamma C'} + 2 \right) (\mathcal{E}_{t,1} + \mathcal{E}_{t,2}).$$
(71)

#### В **EXPERIMENTAL DETAILS**

We implement all methods with PyTorch and run experiments on Nvidia V100 GPUs. For fair comparisons, we adopt the same training epochs and hyper-parameters for all methods.

# 1080 B.1 TRAINING TASKS

To evaluate the methods in various scenarios with different non-IID patterns, we consider two imageclassification tasks and two text classification tasks.

1084						
1085			Table 4: Summ	hary of datasets.		
1086					D	-
1087		Dataset	Data Type	Non-IID Type	Partition	-
1088		CIFAR-10	Image	Label skew	Dir(0.1)	
1089		DomainNet	Image	Feature skew	Domain	
1090		XGLUE-NC	Text	Feature skew	Domain	
1091		QA	Text	Feature skew	Domain	-
1092						
1093 1094	The image classifica	tion tasks includ	de:			
1095	CIEAD 10	(Vrighauslay & I	(Linton 2000)	. In this training t	alt wa aana	iden the lebel alterned
1096	• CIFAK-10	$(\mathbf{K} \mathbf{\Pi} \mathbf{Z} \mathbf{\Pi} \mathbf{e} \mathbf{v} \mathbf{s} \mathbf{K} \mathbf{y} \boldsymbol{\alpha} \mathbf{I}$	nt among alia	the following raining to	ask, we cons	$(L_i \text{ ot al} 2022)$ we
1097	adopt Diric	$F(y_i)$ is unlete	with concentry	ation parameter of	= 0.1 amor	s (LI et al., $2022$ ), we
1098						
1099	DomainNe	t (Peng et al., 20	(19): Domaini	Net contains six d	omains of da	ta samples, including
1100	the varying	i, sketch, infogra	apii, painting, $P(x,   y)$	) is different am	in uns traini	Following (Li et al
1101	2021) each	client is allocat	ted random sa	mples from only	one domain	Tonowing (Li et al.,
1102	2021), eder	i chent is unocu	ieu rundoni se	imples from only	one domain	•
1103	On both tasks, we fi	ne-tune a CLIP	Vision Transf	Former (CLIP) mc	del (Radford	d et al., 2021).
1104	Resides we fine tur	a an VI M Dah	arta Basa mo	del on the followi	ng tayt data	set:
1105	Desides, we fille-tuil		enta-Dase mo	der off the followi	ing text uatas	set.
1106	• XGLUE-N	<b>C</b> (Liang et al.,	2020): This is	s a news classifica	tion task co	nsisting of 10 classes.
1107	The news to	exts comprise fi	ve languages	(English, Spanish	n, French, G	erman, and Russian).
1108	We allocate	e one random la	nguage to eac	ch client, which r	aturally intr	oduces domain shift
1109	among clie	nts.				
1110 1111	In addition, we eval	uate a LLaMA-2	2-7B model o	on the QA datasets	S.	
1112						
1113	• $\mathbf{Q}\mathbf{A}$ : The $\mathbf{Q}$	A datasets con	SIST OF TOUR C	Commonly used q	uestion-ans	Wering datasets, i.e.,
1114	APC Free	ond APC Challe	OpenbookQA	(Nilliaylov et al. 20	, 2018), PIQ	A (DISK et al., 2020),
1115	into classif	ication tasks wh	ere the mode	determines the	correct answ	ver for each question
1116	and corresp	onding choices.	We allocate t	he samples from	one random	dataset to each client.
1117	indicating t	he domain shift	among client	ts.		
1118	-		-			
1119	B.2 IMPLEMENTA	TION DETAILS				
1120						
1121	The CLIP model is	pre-trained on th	ne DataComp	dataset and is ad	apted from h	<pre>nttps://github.</pre>
1122	com/openai/CL	IP; the XLM-Ro	berta-Base m	odel is adapted fr	om https:	://huggingface.
1123	co/xlm-roberta	a-base; the LL	LaMA-2-/B m	nodel is adapted fr	om https:	//huggingface.
1124	of the model and fix	/LLaMA-/D-C	s commonly	all training tasks	, we neeze t	(10b) The values of
1125	adopted hyperparam	eters are summa	arized in Table	e 5 For the propo	sed method	we tune the value of
1126	$\lambda$ from {1, 5, 10, 10	0.500.1000.		e 5. i or the prope	sea methoa,	we take the value of
1127	(-,-,-,-)	,,j.				
1128						
1129						
1130						
1131						
1132						
1133						

		Table 5: Impler	mentation details.	
Dataset	CIFAR-10	DomainNet	XGLUE-NC	QA
Model	CLIP	CLIP	XLM-Roberta-Base	LLaMA-2-7B
Batch size	64	64	32	16
Learning rate	0.01	0.01	0.01	2e-5
Local steps*	5	1	-1	-1
Total epochs	50	30	20	2

\*The local steps -1 means clients iterate all training samples (i.e., one single local training epoch).

 $\lambda$