# A CLOSER LOOK AT PERSONALIZED FINE-TUNING IN HETEROGENEOUS FEDERATED LEARNING

Anonymous authors

Paper under double-blind review

### Abstract

Federated Learning (FL) enables privacy-preserving, decentralized model training but faces significant challenges in balancing global generalization and local personalization due to non-identical data distributions across clients. While Personalized Fine-Tuning (PFT) adapts models to local data, excessive personalization often degrades global performance. In this work, we present a comprehensive empirical study encompassing seven diverse datasets, multiple model architectures, and various fine-tuning methods under both covariate and concept shift scenarios. Our extensive evaluation reveals critical limitations in existing PFT methods, which struggle with overfitting and exhibit inconsistent performance across distribution shifts, even with careful hyperparameter tuning and regularization. To address these issues, we identify LP-FT, a simple yet effective strategy that combines Linear Probing with full Fine-Tuning, adapted to the FL setting. LP-FT consistently outperforms existing methods, achieving an optimal balance between local personalization and global generalization across all tested scenarios. By investigating the feature change after PFT, we hypothesize the a phenomena dubbed as federated feature distortion is linked to the global generalization. Motivated by the observation, we provide a theoretical analysis of two-layer linear networks, offering novel insights into the conditions under which LP-FT excels, thereby enhancing our understanding of personalization dynamics in FL. This work contributes in three key areas: (1) a rigorous and comprehensive evaluation of PFT methods under diverse distribution shifts, (2) the introduction of LP-FT as a robust and versatile solution to FL personalization challenges, and (3) theoretical foundations that explain LP-FT's superior effectiveness. Our findings set a new venue for PFT research and provide valuable insights to the broader FL community.

034 035

004

006

008 009

010 011

012

013

014

015

016

017

018

019

020

021

024

025

026

027

028

029

031

### 1 INTRODUCTION

Federated Learning (FL) (McMahan et al., 2017) has emerged as a promising paradigm for collab orative learning from decentralized data, enabling multiple clients to train a shared global model,
 known as General FL (GFL), without compromising data privacy. A significant challenge in FL is
 the differing data distributions among clients, which undermines the effectiveness of conventional FL
 methods on individual clients.

042 In contrast to GFL, Personalized FL (PFL) (Kairouz et al., 2021) emerges as an approach to address 043 this issue by customizing global models for each client. Clients can enhance the global model 044 by applying local updates using their own datasets, a process known as Personalized Fine-Tuning (PFT) (Wu et al., 2022), as shown in Fig. 1 (a). However, this approach often causes models to overfit on local data, thereby compromising the benefits of FL and the generalization performance of the 046 resulting FL models. This is particularly concerning in critical real-world applications, such as FL 047 across multiple hospitals for disease diagnosis, where a local model must not only perform well 048 on hospital patient data, but also generalize effectively to diverse patient populations that may be encountered on-site in the future (Xu et al., 2021). Therefore, balancing the optimization of individual client performance (personalization) with strong global performance (generalization across all or 051 unseen clients) is crucial (Wu et al., 2022; Huang et al., 2024). 052

To address the personalization and generalization trade-off in PFT, a simple and commonly used strategy is personalized regularized fine-tuning, where each client fine-tunes a global model using



Figure 1: Overview of the problem setting and FL strategies investigated in this paper. (a) PFT framework, where each client fine-tunes a global model previously trained via GFL. (b) Three different PFT models: the global FL model, the full FT model, and the LP-FT model; their parameter updating patterns and local/global performance (perf.) under data heterogeneity; The fire icon indicates the actively tuned parameter, the frozen icon represents the fixed weight, and the mixed fire-frozen icon denotes the weight that is not actively tuned.). (c) Visualization of feature distortion under PFL and its possible link to global generalization.

their local data, with no further communication with the server. This method is appealing for its simplicity and broad applicability across datasets (Kumar et al., 2022), eliminating the need for complex model modifications or extensive data preprocessing. The final local fine-tuning phase plays a crucial role, as it influences both the model's adaptability to individual client data and its overall generalizability. Effective PFT can improve local performance without compromising global performance, whereas poorly executed PFT can result in severe overfitting (Wu et al., 2022).

In this work, we conduct comprehensive evaluation of various regularization-based PFT methods in heterogeneous FL environments under different distribution shift scenarios, categorized as covariate shift (Peng et al., 2019a; Hendrycks & Dietterich, 2019) and concept shift (Izmailov et al., 2022).
Covariate and concept shifts are more prevalent in real-world scenarios where client data often varies in feature distributions or task definitions (Hendrycks & Dietterich, 2019). Despite meticulously tuning the hyper-parameters in conventional regularized fine-tuning methods, we observe persistent issues of local overfitting, wherein localized performance gains are achieved at the significant cost of compromised global generalization.

082 By observing this, we identify a straightforward alternative strategy—combining linear probing 083 (LP) followed by full fine-tuning (FT), referred to as LP-FT (Kumar et al., 2022)—that consistently 084 balances personalization and generalization across diverse distribution shifts, as shown in Fig. 1 (b). We hypothesize that LP-FT achieves this by preserving the representations learned in the last layer 085 of the neural network (before the linear head) during local adaptation, thereby reducing the risk of 086 federated feature distortion, as shown in Fig 1 (c). This phenomenon occurs when local fine-tuning 087 disrupts these last-layer representations (features), which are critical for global performance, leading 088 to performance degradation across the global client set. Although feature distortion was introduced 089 in standard LP-FT approaches within the centralized domain (Kumar et al., 2022), the setting differs 090 significantly from the one considered here. 091

To further support our findings on the superior performance of LP-FT, we conduct theoretical analysis by decompling feature extractor and linear head in neural networks. We examine how LP-FT adapts to client-specific local data and its resulting impact on generalization performance. The developped theoretical frameworksheds light on the conditions when LP-FT outperforms full FT-based PFT<sup>1</sup>, offering new insights into its efficacy within FL.

In summary, our analysis justifies the superiority of LP-FT over standard fine-tuning methods under different distribution shift scenarios. In this paper, we make the following key contributions: 098 **Evaluation:** We present comprehensive empirical evaluations, spanning seven diverse datasets, various models, and different fine-tuning methods within the PFL framework on various distribution 100 shift settings. Our evaluation exposes key limitations in existing PFT methods, highlighting their 101 tendency to overfit and exhibit inconsistent performance across distribution shifts, despite careful 102 hyperparameter tuning and regularization. Insight: We introduce LP-FT, a simple yet effective 103 PFT strategy that enhances both local and global performance. We reveal, for the first time, that 104 preserving pre-trained global features is linked to improved global performance in PFT settings, 105 and we attribute FT's degraded performance to catastrophic feature distortion. Theory: We 106 offer a rigorous theoretical analysis of LP-FT using two-layer linear networks, demonstrating its

107

<sup>&</sup>lt;sup>1</sup>In the following content, we use LP-FT and FT in reference to PFT with LP-FT and full FT, respectively.

108 superior ability to preserve global performance compared to FT in both concept shift and combined 109 concept-covariate shift scenarios.

110 111 112

2 **RELATED WORK** 

113 Heterogeneous FL. Heterogeneous FL refers to a decentralized training paradigm that accommodates 114 diverse and disparate data sources or devices participating in a collaborative model-building process. 115 Examples include FedAvg (McMahan et al., 2017) and various improvements in terms of aggregation 116 optimization and local optimization. It is shown to experience challenges in heterogeneous scenarios. 117 Thus, various literature proposes alternative strategies. Here, we summarize these strategies into 118 aggregation optimization and local optimization. FedNova (Wang et al., 2020) belongs to the 119 category of Aggregation optimization, which normalizes and scales the local updates. Examples 120 of local optimization include FedProx (Li et al., 2020a) and Scaffold (Karimireddy et al., 2020), 121 where FedProx adds a  $L_2$  regularization for each client and Scaffold adds a variance reduction term. However, these methods often exhibit limited personalization capabilities and may not adequately 122 meet the performance requirements of different clients. Consequently, various personalized FL 123 approaches have been proposed, with a primary emphasis on enhancing local client performance to 124 the greatest extent possible. We can group these personalized FL strategies into clustering-based 125 methods (Ghosh et al., 2020), transfer learning (Yu et al., 2020), and interpolating the local and 126 global models (Mansour et al., 2020; Deng et al., 2021). Some FL methods (Guo et al., 2024; Son 127 et al., 2024) highlight that existing federated learning methods often fail under feature shift despite 128 addressing label shift, proposing clustering and regularization strategies respectively to tackle diverse 129 distribution shifts in non-IID data settings. 130

**Regularized Fine-Tuning** Fine-tuning pre-trained models has become increasingly popular with the 131 rise of foundation models (Bommasani et al., 2021). However, fine-tuning with limited data often 132 lead to overfitting. Several strategies can mitigate this issue, such as using optimizers that promote a 133 flatter loss landscape (Li et al., 2018; Kaddour et al., 2022). Notably, Sharpness-Aware Minimization 134 (SAM) (Foret et al., 2021) and Stochastic Weight Averaging (SWA) (Izmailov et al., 2018) are two 135 popular methods that help achieve this. Additionally, a recent technique called *model soups* (Wortsman 136 et al., 2022b), uses a simple greedy weight averaging approach similar to SWA, shown significant 137 improvements in fine-tuning. An interesting perspective focuses on minimizing the linear mode 138 connectivity barrier between the pre-trained and fine-tuned models, helping maintain consistency in decision-making mechanisms from a loss landscape perspective (Vlaar & Frankle, 2022). Partial 139 *fine-tuning* is another common method to prevent overfitting, which involves selectively fine-tuning 140 specific layers of the model to better adapt to variations in data distribution (Lee et al., 2023). Recent 141 studies have introduced the concept of LP-FT (Kumar et al., 2022), highlighting potential distortions 142 in pre-trained features and their underperformance in scenarios involving previously unseen data. 143 Further research on LP-FT provides a deeper analysis of model adaptation (Trivedi et al., 2023), 144 focusing on feature distortion and simplicity bias, thereby enhancing our understanding of fine-tuning 145 mechanisms and safe model adaptation. A detailed literature review is introduced in App. B.

146 147

148

151

153

#### 3 **OBSERVATION: OVERFITTING IN PERSONALIZED FINE-TUNING**

149 In this section, we describe our experimental settings and evaluation methods to analyze personalized 150 overfitting and the effects of LP-FT. This section is organized as follows: (1) defining the problem overview and the distribution shifts considered within a data heterogeneity context, (2) introducing 152 the experimental setups, including the datasets and PFT strategies under investigation, (3) demonstrating the tendency of personalized overfitting with full FT, even with careful regularization, across 154 distribution shifts, (4) proposing LP-FT and evaluating its performance against state-of-the-art full 155 FT strategies in FL, and (5) examining federated feature distortion.

156 157 158

3.1 OVERVIEW AND DEFINITIONS

159 **Problem Setting.** In a FL setting, each client  $i \in [C]$  has a local dataset  $(\mathbf{X}_i, \mathbf{Y}_i)$  generated from a potentially distinct distribution, which may differ across clients due to distribution shifts. The 160 goal is to perform personalized fine-tuning for each client by optimizing local model parameters  $\theta_L$ , 161 initialized from a well-trained global model  $\theta_G$ . The objective is to minimize the local loss  $\mathcal{L}_L(\theta_L)$ 

for improved local performance while ensuring that the global loss  $\mathcal{L}_G(\theta_L)$  remains close to that of a pre-trained global model. This creates a trade-off between personalization (minimizing local loss) and maintaining generalization (minimizing global loss) across clients. The global data distribution  $\mathcal{D}_G$  is defined as a mixture of the local distributions  $\mathcal{D}_i$ , given by  $\mathcal{D}_G = \frac{1}{C} \sum_{i \in [C]} \mathcal{D}_i$ .

We formally define distributions of interests, concept shift and covariate shift in heterogeneous FL context, following Li et al. (2021b).

**Covariate Shift.** Covariate shift refers to variations in the input feature distribution across clients while keeping the conditional distribution of the output given the input consistent. Formally, for any pair of clients *i* and *j* with  $i \neq j$ , the data-generating process is characterized by:

$$P_i(x) \neq P_j(x)$$
, but  $P_i(y \mid x) = P_j(y \mid x)$  for all  $i \neq j$ .

This means that while clients *i* and *j* may have different input distributions  $P_i(x)$  and  $P_j(x)$ , the conditional distribution  $P(y \mid x)$  remains consistent across all clients.

**Concept Shift.** Concept shift, in contrast, occurs when the conditional relationship between input features and outputs varies across clients, while the input feature distribution remains unchanged. The data-generating process for any pair of clients *i* and *j* with  $i \neq j$  is given by:

$$P_i(y \mid x) \neq P_i(y \mid x)$$
, but  $P_i(x) = P_i(x)$  for all  $i \neq j$ .

In this scenario, all clients share the same input distribution P(x), but the conditional distribution  $P_i(y \mid x)$  varies, indicating different relationships between features and labels across clients.

### 183184 3.2 Empirical Analysis Setting

185 Datasets with Cvariate Shift. We include Digit5, DomainNet, CIFAR10-C, and 186 CIFAR100-C. Digit5 and DomainNet belong to the *feature-shift* subgroup, where the data 187 features represent different subpopulations within the same classes. For example, Digit5 contains 188 10-digit images collected from various sources with different backgrounds, such as black-and-white for MNIST and colorful digits for synthetic datasets. CIFAR10-C and CIFAR100-C belong to the 189 input-level shift subgroup, where noise or distortion is introduced in the input data, such as blurred 190 images or sensor errors, degrading input quality. A detailed explanation of the data splitting and its 191 introduction is provided in Table 3 in the Appendix. Visualization is provided in Fig. 5. 192

193 Datasets with Concept Shift. We include CheXpert, and CelebA. CheXpert and CelebA 194 belong to the *spurious correlation-based shift* subgroup, which involves misleading relationships 195 in the training data that models may exploit (also known as spurious correlations), despite being 196 unrelated to the actual target. This reliance can lead to poor performance when such correlations are 197 absent in new data, classifying it as a form of concept shift (Izmailov et al., 2022). Similarly, the 198 detailed information is provided in Table 3 in the Appendix. Visualization is provided in Fig. 5.

199 Fine-tuning Strategies. In this study, we explore several common fine-tuning strategies in PFL: Full-200 parameter FT, a naive FT strategy that adjusts all model parameters. Proximal FT (Li et al.,  $20\overline{20b}$ ) aims to preserve the pre-trained model's original knowledge by applying proximal regularization, 201 which penalizes large deviations from the initial model parameters to maintain generalization. Soup 202 FT (Wortsman et al., 2022a) improves robustness by averaging the weights of multiple fine-tuned 203 model instances, each trained with different initializations, creating a "model soup" that integrates 204 their strengths. Lastly, Sparse FT (Lee et al., 2018) promotes sparsity in parameter updates, adjusting 205 only the most relevant weights to enhance efficiency and interpretability, particularly for deployment 206 in resource-constrained environments. Each strategy is designed to balance model performance with 207 different priorities, such as preserving knowledge, enhancing robustness, or improving efficiency. A 208 more detailed experiment setting is presented in App. C.

209 210 211

172

179

181

182

### 3.3 GLOBAL AND LOCAL PERFORMANCE TRENDS IN COMMON PFT METHODS

In practice, PFT is susceptible to overfitting to local data, due to the relatively small amount of
 data available at local clients. Unlike conventional overfitting, the overfitting referred to in the FL
 context is characterized by a consistent improvement in local performance while global performance
 noticeably deteriorates (Wu et al., 2022; Chen et al., 2023) – the average gain in local performance
 can be smaller than the loss in global performance. To measure the model's overall local and global



Figure 2: Visualization of the prevalence of personalization overfitting across different distribution shift scenarios, where (a) shows the global and local accuracy under different learning rates for full-parameter fine-tune; (b) shows the different sparsity rate for sparse fine-tune; (c) shows the different regularization strength under the proximal fine-tune. In all subfigures, the global accuracy is shown as the solid line, and the local accuracy is shown as the dashed line. As shown, the global accuracy is kept constant while the local accuracy drops in all hyper-parameter settings. This indicates that the PFT model overgeneralizes due to feature distortion.

performance, we used the averaged client-wise local and global accuracy. The metric's decreasing
 trend with increasing local training epochs indicates personalized overfitting. Notably, this trend
 persists even when considering only global performance metrics, as local performance tends to show
 increases in PFT under overfitting conditions.

238 In all subplots of Fig. 2, we observe the tendency of common PFT strategies in FL to overfit under various distribution shift conditions, including input-level shifts (CIFAR100-C), feature-level 239 shifts (Digit5), and spurious correlation-based distribution shifts (CheXpert). We systematically 240 adjusted hyperparameters to evaluate their impact on performance. Notably, we tuned the parameters 241 to avoid a significant increase in local performance, ensuring that any decline in the y-axis primarily 242 reflects a drop in global performance. Fig. 2a illustrates that overfitting persists even when fine-tuning 243 methods use SGD optimizers and reduced learning rates. Fig. 2b shows that, in scenarios where the 244 sparsity rate is adjusted (with sparsity involving masking operations on parameter gradients, meaning 245 that higher sparsity rates result in fewer parameters being updated), overfitting still occurs as the 246 number of training epochs increases in sparsity fine-tuning methods. Fig. 2c depicts the average local 247 and global performance of fine-tuning methods using proximal regularization terms. Similarly, even 248 after adjusting the regularization terms to bias updates towards the initial global model, overfitting 249 remains evident. More empirical studies can be found in App. D.

250 251

252

269

### 3.4 LP-FT: Performance Comparison Across Fine-Tuning Strategies

**LP-FT.** Observing the challenges of personalized overfitting in common regularized fine-tuning 253 methods in PFT, we propose a straightforward approach called Linear Probing and Subsequent 254 Fine-tuning (LP-FT), which demonstrates strong personalization and generalizability across diverse 255 datasets. In practice, clients utilize LP-FT by performing linear probing on the global model's linear 256 layers with local data, followed by fine-tuning the full model. Although the concept of LP-FT was 257 first proposed in centralized setting (Kumar et al., 2022), describing the fine-tuning of a pre-trained 258 model from data in domain A to downstream data in domain B, leading to performance degradation 259 on out-of-distribution testing due to domain misalignment, we refer LP-FT differently in PFT. 260 Specifically, the initial pre-trained model is trained on the combined data from all clients (global 261 distribution) and each client has a distinct ground-truth function. This is in contrast to the setting in 262 the ground-truth function, which is the same for all the data points in Kumar et al. (2022). Therefore, 263 it is novel to explore the efficacy of LP-FT in the PFT setting.

Experimental Settings. To demonstrate the effectiveness of LP-FT in PFT, we focus on comparing different FT methods at the local PFT stage (see Fig. 1 (a)). We use the FedAvg algorithm to train this shared global model in a GFL manner. After the GFL stage, all the clients further fine-tune the obtained global model using local data for 15 epochs for personalization. The final models are

<sup>&</sup>lt;sup>1</sup>The worst group accuracy on CIFAR10 is not noted as all baseline FT results are very close to each other, where comparing STD, in this case, is less meaningful.

303

304

305

Table 1: Performance of various PFT strategies. Red represents the *input shift* subgroup; green from the *feature-shift* subgroup; blue the *spurious correlation-based shift* subgroup. Each experiment is performed three times independently with different random seeds, and the standard deviation of the results is presented in parentheses.  $\uparrow$  indicates that higher values are better, while  $\downarrow$  indicates that lower values are better.

| Dataset    | Method      | Local ↑      | Global ↑     | C-Std.↓      | Worst ↑      | Average ↑    |
|------------|-------------|--------------|--------------|--------------|--------------|--------------|
|            | FT          | 54.50 (0.64) | 44.16 (0.13) | 10.04 (0.06) | 19.83 (0.18) | 39.50 (0.33) |
|            | Proximal FT | 61.76 (0.13) | 53.58 (0.14) | 11.61 (0.08) | 25.82 (0.12) | 47.05 (0.07) |
| CIFAR10-C  | Soup FT     | 56.36 (0.23) | 44.94 (0.06) | 10.22 (0.06) | 20.47 (0.35) | 40.59 (0.09) |
|            | Sparse FT   | 61.31 (0.01) | 50.21 (0.17) | 11.10 (0.11) | 24.56 (0.09) | 45.36 (0.04) |
|            | LP-FT       | 63.55 (0.04) | 55.35 (0.01) | 12.45 (0.01) | 26.33 (0.06) | 48.41 (0.03) |
|            | FT          | 20.05 (0.05) | 14.45 (0.04) | 5.37 (0.02)  | 3.37 (0.06)  | 12.62 (0.03) |
|            | Proximal FT | 27.38 (0.15) | 19.96 (0.11) | 6.90 (0.04)  | 4.84 (0.04)  | 17.41 (0.05) |
| CIFAR100-C | Soup FT     | 20.99 (0.24) | 14.81 (0.04) | 5.48 (0.03)  | 3.56 (0.01)  | 13.12 (0.06) |
|            | Sparse FT   | 28.93 (0.04) | 20.66 (0.02) | 7.75 (0.02)  | 5.05 (0.09)  | 18.15 (0.10) |
|            | LP-FT       | 32.60 (0.14) | 25.44 (0.10) | 9.66 (0.04)  | 5.92 (0.06)  | 21.32 (0.04) |
|            | FT          | 91.17 (0.90) | 67.87 (0.74) | 22.93 (0.28) | 42.03 (0.48) | 67.02 (0.70) |
|            | Proximal FT | 92.09 (0.18) | 81.40 (0.03) | 15.04 (0.15) | 61.71 (0.16) | 78.40 (0.09) |
| Digit5     | Soup FT     | 91.82 (0.34) | 70.82 (0.43) | 21.99 (0.67) | 45.10 (1.27) | 69.02 (0.65) |
|            | Sparse FT   | 91.43 (0.31) | 76.89 (0.72) | 17.90 (0.38) | 54.21 (0.56) | 74.21 (0.35) |
|            | LP-FT       | 91.20 (0.04) | 82.78 (0.05) | 13.75 (0.02) | 65.80 (0.02) | 79.92 (0.02) |
|            | FT          | 64.90 (1.18) | 42.48 (0.58) | 17.49 (0.75) | 22.31 (0.93) | 43.23 (0.52) |
|            | Proximal FT | 67.20 (1.39) | 56.05 (0.27) | 16.68 (0.36) | 33.20 (1.79) | 52.60 (0.35) |
| DomainNet  | Soup FT     | 67.48 (0.61) | 44.27 (0.46) | 18.44 (0.42) | 23.73 (1.24) | 44.49 (0.54) |
|            | Sparse FT   | 69.62 (0.53) | 50.24 (0.44) | 18.14 (0.17) | 27.89 (0.15) | 49.14 (0.45) |
|            | LP-FT       | 68.50 (0.19) | 57.52 (0.20) | 17.36 (0.21) | 34.53 (0.44) | 53.52 (0.19) |
|            | FT          | 76.18 (0.41) | 76.25 (0.56) | 0.35 (0.13)  | 76.31 (0.76) | 76.25 (0.44) |
|            | Proximal FT | 76.44 (0.07) | 76.63 (0.09) | 0.71 (0.09)  | 76.81 (0.07) | 76.63 (0.07) |
| CheXpert   | Soup FT     | 77.51 (0.15) | 77.49 (0.31) | 0.48 (0.07)  | 77.46 (0.43) | 77.49 (0.26) |
|            | Sparse FT   | 77.29 (0.13) | 77.20 (0.14) | 0.31 (0.11)  | 77.11 (0.25) | 77.20 (0.14) |
|            | LP-FT       | 77.64 (0.37) | 77.54 (0.37) | 0.53 (0.41)  | 77.43 (0.71) | 77.54 (0.37) |
|            | FT          | 90.55 (1.20) | 73.76 (2.15) | 18.79 (3.64) | 53.52 (5.51) | 72.39 (2.84) |
|            | Proximal FT | 93.74 (0.59) | 81.11 (0.82) | 13.39 (1.14) | 67.50 (2.10) | 80.78 (0.90) |
| CelebA     | Soup FT     | 89.42 (2.16) | 75.28 (1.11) | 16.29 (1.19) | 57.79 (2.90) | 74.17 (1.50) |
|            | Sparse FT   | 91.43 (0.48) | 77.32 (1.46) | 14.16 (2.57) | 62.94 (4.34) | 77.65 (1.65) |
|            | LP-FT       | 93.24 (0.17) | 83.32 (0.31) | 11.18 (0.14) | 71.89 (0.75) | 82.82 (0.64) |

evaluated using the *metrics* described below. Detailed descriptions of the datasets, preprocessing steps, data splitting, and models used are provided in Appendix C.3, Tab. 3.

Metrics. We report five metrics in this part of the experiment: (1) Local Accuracy (Local) measures 306 the performance of the PFT model on the client's local test set. Higher Local Acc indicates better 307 personalization. (2) Global Accuracy (Global) measures the PFT model's average test accuracy over 308 all other clients' test sets. Higher Global Acc indicates better generalization. (3) Client-wise Standard 309 Deviation (C-Std.) calculates the standard deviation of local test accuracies across all clients. Lower 310 C-Std. indicates less variance in performance among clients. (4) Worst Accuracy (Worst) reports the 311 lowest test accuracy among all clients. The closer this value is to Local Acc, the better the worst-case 312 generalization. (5) Average reports the average of both Local Acc and Global Acc, providing a better 313 understanding of the tradeoff between personalization (local performance) and generalization (global 314 performance). All metrics, except C-Std., are averaged over the number of clients, and higher values 315 are preferable. For the C-Std. metric, lower values are better.

316 Results. Our results are presented in Tab. 1, where the best method is highlighted in **bold**. Datasets 317 with the same distribution shift pattern are grouped into the same colors as detailed in the caption. 318 Tab. 1 shows that LP-FT consistently achieves the highest global and average accuracy across 319 most datasets, demonstrating strong generalization and personalization performances, particularly 320 in challenging conditions like CIFAR100-C and CIFAR10-C. Sparse FT also performs well, 321 especially in Digits5 and DomainNet, but generally lags behind LP-FT. Soup FT and Proximal FT show mixed results, with stronger performance in specific datasets such as CheXpert but 322 weaker overall compared to LP-FT. Standard fine-tuning consistently underperforms, highlighting 323 the limitations of basic fine-tuning methods in heterogeneous data scenarios.



331 332

333

334

362

364

366

367

368

369

370

371

372

373

Figure 3: Illustration of federated feature distortion (FD) and decision boundaries.

### 3.5 INSIGHT AND EXPLANATION ON THE OBSERVATIONS

Federated Feature Distortion. We hypothesize that personalized overfitting in PFT and the associ-335 ated performance degradation are linked to federated feature distortion, which refers to the *alteration* 336 of intermediate data representations from the global model (i.e., network features before the linear 337 head) during the final local training phase in PFL, as illustrated in Fig. 1. This distortion, driven by 338 data heterogeneity across clients, occurs when local models overfit to client-specific distributions. 339 To illustrate how feature distortion affects performance, we present Fig. 3, showcasing model deci-340 sion boundaries and features. Specifically, after GFL, the global server model is well-trained with 341 well-separated decision boundaries for global features. However, it may misclassify some local data 342 points. While fine-tuning local clients with full FT can improve local performance, it severely distorts 343 their feature representations compared to those from global model, potentially resulting in shifted decision boundaries that harm global performance. In contrast, LP-FT uses linear probing first, where 344 only the linear head is tuned, preserving the pre-trained feature representations. The subsequent 345 fine-tuning after linear probing causes only minor adjustments to full-model parameters and features, 346 as the model is already close to optimal for the local data. This process can result in moderate feature 347 distortion, simultaneously fitting the model to local data while maintaining global performance. 348

349 Experimental Validation on the Relationship Between Feature Distortion and Performance. In this part, we utilize multiple datasets (*i.e.*, DomainNet and Digit5) to investigate the relationship 350 between global performance and feature distortion. We measure the similarity of representations in 351 the feature extraction layer (*i.e.*, the input to the classification head) in the following way. Consider a 352 feature extraction function  $f : \mathcal{X} \to \mathbb{R}^k$ , which maps inputs from the input space  $\mathcal{X}$  to a representation 353 space  $\mathbb{R}^k$ . Let  $\theta_G$  denote the global pre-trained model and  $\theta_i$  the fine-tuned model after local fine-354 tuning for client i. Assume there are C clients in total, each with n samples. Let  $x_{c,j}$  represent the j-th 355 data point of the c-th client. The federated feature distortion  $\Delta_c(f)$  quantifies the change in features 356 after fine-tuning for the c-th client, defined as the average  $\ell_2$  distance between the representations 357 produced by the global model and the locally fine-tuned model over all data points across all clients. 358 Formally, it is expressed as:  $\Delta_c(f) = \frac{1}{n} \sum_{j=1}^n \|f(\theta_G; x_{c,j}) - f(\theta_c; x_{c,j})\|_2$ , where  $\|\cdot\|_2$  is the 359  $\ell_2$  distance in the representation space  $\mathbb{R}^k$ . We compute the average of  $\Delta_c(f)$  across all clients to 360 represent the feature distortion in the PFT setting, as shown in Fig. 4. 361

<u>Observation 1.</u> Our findings reveal that feature distortion is correlated with a drop in global performance. Specifically, in Fig. 4(a), common FT methods cause significant feature distortion, leading



Figure 4: Observations of the feature distortion in our PFT setting, where (a) presents the positive correlation between global performance drops and feature distortion intensity on DomainNet and
(b) presents the ablation study on preserving federated features with controlled local train loss on Digit5. We set local loss thresholds (0.1, 0.5, and 1.0) and used gradient ascent when the loss fell below, ensuring training loss fluctuated around these points.

to a substantial decline in global performance. In contrast, LP-FT maintains a high level of global
performance. This stability can be attributed to the initial phase of LP-FT (the first 5 epochs), during
which only the classifier is fine-tuned while the parameters of the feature extraction layer remain
fixed. Compared with FT, LP-FT suffers from significantly moderate feature distortion.

382 *Observation 2.* In Fig. 4(b), we further analyze the relationship between global performance and feature distortion by eliminating the influence of local loss magnitude. We achieve this by controlling 384 the level of local training loss using the loss flooding technique (Ishida et al., 2020). Specifically, we set thresholds for local loss (*i.e.*, 0.1, 0.5, and 1.0) and apply gradient ascent when the loss drops 386 below these thresholds, ensuring that the training loss fluctuates around the set points. Examining a 387 single FT strategy (*i.e.*, either FT or LP-FT), we observe that global accuracy remains stable despite 388 variations in local training loss magnitudes. At fixed local training loss levels, we compare the effects of LP-FT (moderate feature distortion) and FT (severe feature distortion) on global performance. 389 LP-FT consistently delivers better global performance than FT across different loss levels. These 390 comparisons enable us to dismiss the alternative explanation that FT leads to higher global loss simply 391 due to achieving lower local loss, while LP-FT exhibits the opposite pattern. 392

- 393
- 394
- 394 395
- 306

### 4 THEORETICAL ANALYSIS OF THE LP-FT METHOD

397 In Sec. 3, we presented a series of experiments demonstrating the effectiveness of LP-FT compared to common FT strategies in PFL. Our results indicate that FT performs poorly relative to LP-FT, 398 and we hypothesize that this suboptimal global performance arises from federated feature distortion. 399 To further understand how feature learning impacts generalization error in PFT, we decompose 400 the data-generating function and the model into two components: a feature extractor and a linear 401 head. This decomposition allows us to distinguish between the learned features and their influence 402 on performance. To further explain the empirical observations in Sec. 3, we provide a theoretical 403 analysis of LP-FT's global performance under various distribution shift scenarios. Specifically, in 404 Sec. 4.1 and Sec. 4.2, we formalize concept and covariate shifts within a two-layer linear network 405 and examine how LP-FT effectively adapts to these shifts, outperforming FT in the PFL setting.

Overview of Theoretical Analysis: To compare the performance of LP-FT and FT, we make
 assumptions about the data-generating function for clients (Assumption 4.1) and a specific model
 structure (Assumption 4.2). Based on these assumptions, we analyze the global performance of LP-FT
 and FT under concept shift (Theorem 4.4) and combined concept-covariate shift (Theorem 4.5).

410 411 412

413

### 4.1 LP-FT'S GLOBAL PERFORMANCE UNDER CONCEPT SHIFT

In this section, we analyze LP-FT's performance compared to FT under concept shift. To facilitate a rigorous theoretical study, we define the data-generating process and model structure across clients, assuming both are represented by two-layer linear networks, as in (Kumar et al., 2022).

417 Assumption 4.1 (Data-Generating Process). The data-generating function for client *i* is given by 418  $y_i = V_i^{*T} B_* x_i$  for all  $i \in [C]$ , where  $y_i \in \mathbb{R}$ , *C* is the number of clients,  $x_i \in \mathbb{R}^d$ ,  $B_* \in \mathbb{R}^{k \times d}$ , and 419  $V_i^* \in \mathbb{R}^k$ . All clients share a common feature extractor  $B_*$ , assumed to have orthonormal rows, while 420 their linear heads  $V_i^*$  differ. Each  $V_i^*$  decomposes as  $V_i^* = \begin{bmatrix} V_{com}^* & \lambda e_i^T \end{bmatrix}^T$ , where  $V_{com}^* \in \mathbb{R}^m$  is 422 shared across clients,  $e_i \in \mathbb{R}^C$  is a unit vector, and  $\lambda$  controls heterogeneity. Here, m + C = k.

This assumption distinguishes between a shared and client-specific component in the data-generating functions, allowing analysis of both global and local performance of PFT methods after fine-tuning.

426 Assumption 4.2 (Model Structure). The training model is a two-layer linear network defined as 427  $y = V^T B x$ , where  $V \in \mathbb{R}^k$  is the linear head and  $B \in \mathbb{R}^{k \times d}$  is the feature extractor. The dimensions 428 of V and B match Assumption 4.1, allowing the model to learn both shared and client-specific data 429 components.

430

423

431 In PFT settings, our objective is to evaluate the performance of a model on both global and local data. By local data, we refer to the data of a specific client undergoing fine-tuning (e.g., client *i*). The local and global losses are defined using the Mean Squared Error (MSE) as follows:

433 434 435

436 437

453 454 455

456

432

$$\mathcal{L}_{L}(V,B) = \mathbb{E}_{(x,y)\sim\mathcal{D}_{i}}\left[\frac{1}{2}(V^{T}Bx - y)^{2}\right] = \mathbb{E}_{x\sim\mathcal{D}_{i}}\left[\frac{1}{2}(V^{T}Bx - V_{i}^{*T}B_{*}x)^{2}\right],$$
$$\mathcal{L}_{G}(V,B) = \mathbb{E}_{(x,y)\sim\mathcal{D}_{G}}\left[\frac{1}{2}(V^{T}Bx - y)^{2}\right] = \frac{1}{C}\sum_{i\in[C]}\mathbb{E}_{x\sim\mathcal{D}_{i}}\left[\frac{1}{2}(V^{T}Bx - V_{i}^{*T}B_{*}x)^{2}\right]$$

Since this section focuses on concept shift, we assume all clients' data is drawn from similar distributions. Accordingly, we assume for every client  $i \in [C]$ , the input features satisfy  $\mathbb{E}_{x \sim \mathcal{D}_i}[xx^T] = I_d$ .

With the theoretical framework established by Assumptions 4.1 and 4.2, we compare the global performance of LP-FT and FT, highlighting cases where LP-FT outperforms FT. In a PFL setting, the initial model is trained on data from all clients to capture their shared components. Thus, we initialize the model parameters as  $B_0 = B_*$  and  $V_0 = \begin{bmatrix} V_{com}^* & \mathbf{0} \end{bmatrix}^T$ . In LP-FT, a step of linear probing first updates  $V_0$  using local data while keeping  $B_0$  fixed, followed by full fine-tuning to update both V and B. In contrast, FT performs only the second step. The following lemma characterizes B after one gradient descent step in FT, forming the basis for our comparison.

**Lemma 4.3.** Under Assumptions 4.1 and 4.2, and assuming that  $\mathbb{E}_{x \sim D_i}[xx^T] = I_d$  for all clients  $i \in [C]$ , let the initial parameters before starting FT be  $B_0 = B_*$  and  $V_0 = \begin{bmatrix} V_{com}^* & \mathbf{0} \end{bmatrix}^T$ . Assume fine-tuning is performed locally on the data of the *i*-th client. Let  $B_{FT}$  denote the feature extractor matrix after a single gradient descent step (processing the entire dataset once) with learning rate  $\eta$ . If  $(b_j^{FT})^T$  is the *j*-th row of  $B_{FT}$ , then:

$$\mathbb{E}\left[(b_j^{FT})^T\right] = (b_j^*)^T + \eta \lambda(V_0)_j (b_{m+i}^*)^T,$$

where  $(b_j^*)^T$  is the *j*-th row of  $B_*$ , and  $(V_0)_j$  is the *j*-th element of  $V_0$  for  $j \in [k]$ .

This lemma examines the impact of FT on the feature extractor  $B_{FT}$ , highlighting the deviations from the pre-trained matrix  $B_0 = B_*$ . Given that all clients share the same  $B_*$  in their labeling functions, substantial changes to the feature extractor can lead to a decline in global performance. Since the matrix B functions as the feature extractor in our framework, significant feature distortion occurs when  $B_{FT}$  deviates considerably from  $B_*$ . Building on Lemma 4.3, Theorem 4.4 offers a comparative analysis of the global performance of LP-FT versus FT in the context of concept shift.

**Theorem 4.4.** Under Assumptions 4.1 and 4.2, and assuming  $\mathbb{E}_{x \sim \mathcal{D}_i}[xx^T] = I_d$  for all clients  $i \in [C]$ , let the initial model parameters be  $B_0 = B_*$  and  $V_0 = \begin{bmatrix} V_{com}^* & \mathbf{0} \end{bmatrix}^T$ . Let  $B_{FT}$  and  $V_{FT}$ denote the parameters of the FT method after one gradient descent step (processing the entire dataset once). For LP-FT, let  $B_{LPFT}$  and  $V_{LPFT}$  denote the parameters after (i) linear probing, which optimizes V with B fixed at  $B_*$ , and (ii) one gradient descent step with learning rate  $\eta$ . Then:

$$\mathcal{L}_G(V_{LPFT}, B_{LPFT}) \le \mathcal{L}_G(V_{FT}, B_{FT})$$

This theorem characterizes the global performance of LP-FT, suggesting that under concept shift,
LP-FT achieves better performance on global data than FT. When starting from a model initialized
to capture the shared feature extractor and linear head among clients, LP-FT is more effective in
minimizing global loss, aligning with common FL scenarios where the initial model leverages shared
client structure.

475 476

477

468 469

#### 4.2 LP-FT'S GLOBAL PERFORMANCE UNDER COMBINED CONCEPT AND COVARIATE SHIFTS

In the previous section, we assumed all clients' data came from the same distribution with  $\mathbb{E}_{x \sim \mathcal{D}_i}[xx^T] = I_d$ . However, this may not hold in many practical scenarios. To address this, we introduce covariate shift, where each client's data is generated as  $x_i = e_i + \epsilon n$ , with  $n \sim \mathcal{N}(0, I)$ ,  $e_i$  as a client-specific shift, and  $\epsilon$  controlling the noise level. The model structure and data-generating assumptions remain consistent with Sec. 4.1. This section thus considers both concept and covariate shifts. Theorem 4.5 analyzes the impact of heterogeneity on the global performance of LP-FT and FT.

Theorem 4.5. Under Assumptions 4.1 and 4.2, let each client's data be  $x_i = e_i + \epsilon n$ , where  $n \sim \mathcal{N}(0, I)$  and  $e_i$  is a client-specific shift. Assume the initial parameters are  $B_0 = B_*$  and  $V_0 = \begin{bmatrix} V_{com}^* & \mathbf{0} \end{bmatrix}^T$ . Let  $B_{FT}, V_{FT}$  be the FT parameters after one gradient descent step, and  $B_{LPFT}, V_{LPFT}$  be the LP-FT parameters after linear probing and one gradient descent step (with learning rate  $\eta$ ). Then, there exists a threshold  $\lambda^*$  such that for all  $\lambda \leq \lambda^*$ :

$$\mathcal{L}_G(V_{LPFT}, B_{LPFT}) \le \mathcal{L}_G(V_{FT}, B_{FT}).$$

*Remark* 4.6. In Theorem 4.5, the parameter  $\lambda$  characterizes the level of heterogeneity among clients. The theorem shows that under both covariate and concept shifts, LP-FT outperforms FT in low heterogeneity settings ( $\lambda \leq \lambda^*$ ), highlighting its advantage in maintaining generalization. To further reinforce the theoretical insights and cover more extensive settings, Sec. 5 provides empirical validation of our findings, confirming the global superiority of LP-FT over FT. While our theoretical analysis of Theorem 4.5 focuses on the low heterogeneity regime, the experiments in Sec. 5 explore a broader range, simulating both high and low heterogeneity levels. These results validate and extend our theoretical insights, demonstrating that LP-FT consistently outperforms FT across all heterogeneity levels (see also Sec. G in the appendix).

498 499 500

486

487

488 489 490

491

492

493

494

495

496

497

### 5 EXPERIMENT: FURTHER VALIDATIONS FOR THEORETICAL FINDINGS

501

Despite being based on simplified data and model assumptions, our theoretical results demonstrate significant practical relevance. In this section, we empirically validate the contributions in Sec. 4, exploring the performance implications of controllable heterogeneities in neural networks and datasets.

508 **Experimental Settings.** To validate the impact of  $\lambda$  in 509 Theorem 4.5, we simulate a controllable concept shift setting on the Digit5 dataset with label-flipping under 510 PFT for both FT and LP-FT. For each client, a proportion 511 of labels is randomly flipped, referred to as the flipping 512 ratio. For example, class one is flipped with class two 513 for the first client, and class two with class three for the 514 second, using a randomized mechanism. A higher flipping 515 ratio indicates greater heterogeneity  $\lambda$ . The settings align 516 with prior studies: the model is pre-trained within the FL 517 framework and used to initialize both FT and LP-FT. This 518 simulates the combined concept-covariate shift discussed

| Table  | 2:   | Performance       | under    | Label-  |
|--------|------|-------------------|----------|---------|
| Flippi | ng f | or FT and LPI     | FT, with | h LF.R. |
| as the | labe | 1-flipping ratio. |          |         |

| LF.R. (%) | Metric            | FT    | LPFT  |
|-----------|-------------------|-------|-------|
|           | Avg. ↑            | 67.73 | 79.83 |
| 20        | Global ↑          | 68.76 | 83.08 |
|           | Local ↑           | 91.32 | 91.23 |
|           | Avg. ↑            | 60.04 | 72.95 |
| 30        | Global ↑          | 55.18 | 72.75 |
|           | Local ↑           | 91.12 | 89.20 |
|           | Avg. ↑            | 58.27 | 71.55 |
| 40        | Global ↑          | 53.70 | 69.89 |
|           | Local ↑           | 90.84 | 90.02 |
|           | Avg. ↑            | 60.06 | 73.26 |
| 50        | Global $\uparrow$ | 56.17 | 72.32 |
|           | Local ↑           | 91.88 | 90.87 |

in Sec. 4.2. Flipping labels reflects different labeling functions, where higher flipping rates indicate
 stronger concept shifts. The Digit5 dataset also introduces covariate shift, as outlined in Sec. 3.2.

Results. As shown in Tab. 2, LP-FT consistently outperforms FT in global performance across various flipping ratios. This aligns with our theoretical results in Sec. 4.2, especially for deep neural networks under realistic PFT settings. The flipping rate controls concept shift heterogeneity, with higher rates indicating greater heterogeneity, while varying data distributions introduce covariate shift. These experiments simulate the combined concept-covariate shift, as analyzed in our framework. Notably, LP-FT outperforms FT in all heterogeneity levels, validating its advantage in both low and high heterogeneity regimes (larger flipping ratios).

### 6 CONCLUSION

529 530

528

531 In this work, we tackled the key challenge of balancing local personalization and global generalization 532 in PFL. Through an extensive empirical evaluation across seven datasets, multiple model architectures, 533 and various distribution shifts, we revealed critical limitations in existing PFT methods, which often 534 suffer from overfitting and inconsistent performance across scenarios. To address these issues, we introduced a simple yet effective strategy combining Linear Probing with full Fine-Tuning (LP-FT), 536 which consistently outperforms other methods by preserving pre-trained global features and mitigating 537 the adverse effects of excessive personalization. Through in-depth analysis, we attribute LP-FT's strong performance to its ability to prevent feature distortion, which we linked to the degradation of 538 global performance. Furthermore, we provided a theoretical analysis using two-layer linear networks, explaining LP-FT's superior performance, particularly under concept and covariate shift conditions.

## 540 REFERENCES 541

| 541<br>542<br>543   | Manoj Ghuhan Arivazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunav Choudhary. Feder-<br>ated learning with personalization layers. <u>CoRR</u> , abs/1912.00818, 2019.   |
|---|--|
| 544<br>545<br>546<br>547<br>548<br>549<br>550<br>551<br>552<br>552<br>553 | Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ B. Altman, Simran Arora, Sydney von Arx,<br>Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson,<br>Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen<br>Creel, Jared Quincy Davis, Dorottya Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus,<br>Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale,<br>Lauren Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori<br>Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang,<br>Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling,<br>Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi,<br>and et al. On the opportunities and risks of foundation models. <u>CoRR</u> , abs/2108.07258, 2021. |
| 54<br>555<br>56   | Minghui Chen, Zhiqiang Wang, and Feng Zheng. Benchmarks for corruption invariant person re-identification. <u>arXiv preprint arXiv:2111.00880</u> , 2021.  |
| 557<br>558<br>559   | Minghui Chen, Meirui Jiang, Qi Dou, Zehua Wang, and Xiaoxiao Li. Fedsoup: Improving general-<br>ization and personalization in federated learning via selective model interpolation. In <u>MICCAI (2)</u> ,<br>volume 14221 of <u>Lecture Notes in Computer Science</u> , pp. 318–328. Springer, 2023.   |
| i60<br>i61<br>i62   | Liam Collins, Hamed Hassani, Aryan Mokhtari, and Sanjay Shakkottai. Exploiting shared repre-<br>sentations for personalized federated learning. In <u>ICML</u> , volume 139 of <u>Proceedings of Machine</u><br><u>Learning Research</u> , pp. 2089–2099. PMLR, 2021.  |
| 63<br>64<br>65  | Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Adaptive personalized federated learning. <u>CoRR</u> , abs/2003.13461, 2020.   |
| 566<br>567  | Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Distributionally robust federated averaging. <u>CoRR</u> , abs/2102.12660, 2021.  |
| 69<br>70  | Alireza Fallah, Aryan Mokhtari, and Asuman E. Ozdaglar. Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach. In <u>NeurIPS</u> , 2020.  |
| 571<br>572  | Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In <u>ICLR</u> . OpenReview.net, 2021.   |
| 74<br>75  | Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. An efficient framework for clustered federated learning. In <u>NeurIPS</u> , 2020.  |
| 76<br>77<br>78  | Yongxin Guo, Xiaoying Tang, and Tao Lin. Fedrc: Tackling diverse distribution shifts challenge in federated learning by robust clustering. In <u>ICML</u> . OpenReview.net, 2024.  |
| 579<br>580  | Dan Hendrycks and Thomas G. Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In <u>ICLR (Poster)</u> . OpenReview.net, 2019.  |
| i81<br>i82<br>i83<br>i84  | Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for NLP. In <u>ICML</u> , volume 97 of <u>Proceedings of Machine Learning Research</u> , pp. 2790–2799. PMLR, 2019.  |
| i85<br>i86<br>i87   | Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In <u>ICLR</u> . OpenReview.net, 2022.  |
| 588<br>589<br>590<br>591  | Chun-Yin Huang, Ruinan Jin, Can Zhao, Daguang Xu, and Xiaoxiao Li. Federated virtual learning on heterogeneous data with local-global distillation. <u>CoRR</u> , abs/2303.02278, 2023. doi: 10.48550/<br>ARXIV.2303.02278. URL https://doi.org/10.48550/arXiv.2303.02278.   |
| 592<br>593  | Chun-Yin Huang, Kartik Srinivas, Xin Zhang, and Xiaoxiao Li. Overcoming data and model heterogeneities in decentralized federated learning via synthetic anchors, 2024. URL https://openreview.net/forum?id=PcBJ4pA6bF.  |

| 594<br>595<br>596<br>597  | Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik<br>Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest<br>radiograph dataset with uncertainty labels and expert comparison. In <u>Proceedings of the AAAI</u><br><u>conference on artificial intelligence</u> , volume 33, pp. 590–597, 2019.  |
|---|--|
| 598<br>599<br>600<br>601  | Takashi Ishida, Ikko Yamane, Tomoya Sakai, Gang Niu, and Masashi Sugiyama. Do we need zero training loss after achieving zero training error? In <u>ICML</u> , volume 119 of <u>Proceedings of Machine Learning Research</u> , pp. 4604–4614. PMLR, 2020.  |
| 602<br>603<br>604   | Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry P. Vetrov, and Andrew Gordon Wilson.<br>Averaging weights leads to wider optima and better generalization. In <u>UAI</u> , pp. 876–885. AUAI<br>Press, 2018.  |
| 606<br>607  | Pavel Izmailov, Polina Kirichenko, Nate Gruver, and Andrew Gordon Wilson. On feature learning in the presence of spurious correlations. In <u>NeurIPS</u> , 2022.  |
| 608<br>609<br>610<br>611  | Ruinan Jin, Wenlong Deng, Minghui Chen, and Xiaoxiao Li. Debiased noise editing on foundation models for fair medical image classification. In <u>International Conference on Medical Image</u> <u>Computing and Computer-Assisted Intervention</u> , pp. 164–174. Springer, 2024.   |
| 612<br>613  | Jean Kaddour, Linqing Liu, Ricardo Silva, and Matt J. Kusner. When do flat minima optimizers work? In <u>NeurIPS</u> , 2022.   |
| 614<br>615<br>616<br>617<br>618<br>619<br>620<br>621<br>622<br>623<br>624 | Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin<br>Bhagoji, Kallista A. Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, Rafael<br>G. L. D'Oliveira, Hubert Eichner, Salim El Rouayheb, David Evans, Josh Gardner, Zachary Garrett,<br>Adrià Gascón, Badih Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaïd Harchaoui, Chaoyang<br>He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri Joshi,<br>Mikhail Khodak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar, Sanmi Koyejo,<br>Tancrède Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer Özgür, Rasmus<br>Pagh, Hang Qi, Daniel Ramage, Ramesh Raskar, Mariana Raykova, Dawn Song, Weikang Song,<br>Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr, Praneeth Vepakomma,<br>Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu, and Sen Zhao. Advances<br>and open problems in federated learning. Found. Trends Mach. Learn., 14(1-2):1–210, 2021. |
| 625<br>626<br>627   | Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Sebastian U. Stich, and<br>Ananda Theertha Suresh. SCAFFOLD: stochastic controlled averaging for federated learning. In<br><u>ICML</u> , volume 119 of <u>Proceedings of Machine Learning Research</u> , pp. 5132–5143. PMLR, 2020.  |
| 628<br>629<br>630   | Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. <u>corr</u> , 2009.  |
| 631<br>632<br>633   | Ananya Kumar, Aditi Raghunathan, Robbie Matthew Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. In <u>ICLR</u> . OpenReview.net, 2022.  |
| 635<br>636  | Namhoon Lee, Thalaiyasingam Ajanthan, and Philip HS Torr. Snip: Single-shot network pruning based on connection sensitivity. <u>arXiv preprint arXiv:1810.02340</u> , 2018.  |
| 637<br>638<br>639<br>640  | 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 <u>ICLR</u> . OpenReview.net, 2023.  |
| 641<br>642  | Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. In <u>NeurIPS</u> , pp. 6391–6401, 2018.   |
| 643<br>644<br>645   | Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith.<br>Federated optimization in heterogeneous networks. In <u>MLSys</u> . mlsys.org, 2020a.   |
| 646<br>647  | Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith.<br>Federated optimization in heterogeneous networks. <u>Proceedings of Machine learning and systems</u> , 2:429–450, 2020b.  |

648 Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated 649 learning through personalization. In ICML, volume 139 of Proceedings of Machine Learning 650 Research, pp. 6357–6368. PMLR, 2021a. 651 Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. Fedbn: Federated learning 652 on non-iid features via local batch normalization. In ICLR. OpenReview.net, 2021b. 653 654 Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. Fedbn: Federated learning 655 on non-iid features via local batch normalization. arXiv preprint arXiv:2102.07623, 2021c. 656 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In 657 Proceedings of International Conference on Computer Vision (ICCV), December 2015. 658 659 Yishay Mansour, Mehryar Mohri, Jae Ro, and Ananda Theertha Suresh. Three approaches for 660 personalization with applications to federated learning. CoRR, abs/2002.10619, 2020. 661 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. 662 Communication-efficient learning of deep networks from decentralized data. In AISTATS, vol-663 ume 54 of Proceedings of Machine Learning Research, pp. 1273–1282. PMLR, 2017. 664 665 Eric Mintun, Alexander Kirillov, and Saining Xie. On interaction between augmentations and 666 corruptions in natural corruption robustness. Advances in Neural Information Processing Systems, 667 34:3571-3583, 2021. 668 Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching 669 for multi-source domain adaptation. In ICCV, pp. 1406-1415. IEEE, 2019a. 670 671 Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In Proceedings of the IEEE/CVF international conference on 672 computer vision, pp. 1406–1415, 2019b. 673 674 Ha Min Son, Moon-Hyun Kim, Tai-Myoung Chung, Chao Huang, and Xin Liu. Feduv: Uniformity 675 and variance for heterogeneous federated learning. In CVPR, pp. 5863-5872. IEEE, 2024. 676 Puja Trivedi, Danai Koutra, and Jayaraman J. Thiagarajan. A closer look at model adaptation using 677 feature distortion and simplicity bias. In ICLR. OpenReview.net, 2023. 678 679 Tiffany J. Vlaar and Jonathan Frankle. What can linear interpolation of neural network loss landscapes 680 tell us? In ICML, volume 162 of Proceedings of Machine Learning Research, pp. 22325–22341. 681 PMLR, 2022. 682 Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H. Vincent Poor. Tackling the objective 683 inconsistency problem in heterogeneous federated optimization. In NeurIPS, 2020. 684 685 Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, 686 Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model 687 soups: averaging weights of multiple fine-tuned models improves accuracy without increasing 688 inference time. In International conference on machine learning, pp. 23965–23998. PMLR, 2022a. 689 Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, Rebecca Roelofs, Raphael Gontijo Lopes, 690 Ari S. Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig 691 Schmidt. Model soups: averaging weights of multiple fine-tuned models improves accuracy 692 without increasing inference time. In ICML, volume 162 of Proceedings of Machine Learning 693 Research, pp. 23965–23998. PMLR, 2022b. 694 Shanshan Wu, Tian Li, Zachary Charles, Yu Xiao, Ken Ziyu Liu, Zheng Xu, and Virginia Smith. Motley: Benchmarking heterogeneity and personalization in federated learning. CoRR, 696 abs/2206.09262, 2022. 697 Jie Xu, Benjamin S. Glicksberg, Chang Su, Peter B. Walker, Jiang Bian, and Fei Wang. Federated 699 learning for healthcare informatics. J. Heal. Informatics Res., 5(1):1–19, 2021. 700 Tao Yu, Eugene Bagdasaryan, and Vitaly Shmatikov. Salvaging federated learning by local adaptation. 701 CoRR, abs/2002.04758, 2020.

# A NAVIGATING THE TRADE-OFF BETWEEN LOCAL AND GLOBAL PERFORMANCE IN FL

The Challenge of improving both local and global performance. In FL, a primary challenge in enhancing both local and global performance arises from data heterogeneity. The disparities in data distribution among individual clients result in divergent risk functions for each client, consequently leading to disparate optimal solutions. Additionally, in the context of federated learning, each client's local training process is oblivious to the data of other clients, rendering models prone to overfitting to local data distribution after personalized fine-tuning. This personalization overfitting phenomenon detrimentally impacts global performance.

712 Previous approaches to personalized federated learning, primarily based on local training, can be 713 broadly categorized into two main types: Partial fine-tuning (which fine-tunes only specific layers of 714 the global model on local data to retain globally learned features) (Collins et al., 2021; Arivazhagan 715 et al., 2019), and regularization guided by the global model to constrain local updates (Li et al., 716 2021a; Deng et al., 2020). However, while mitigating overfitting to some extent during local training, 717 these methods often face challenges in significantly improving both local and global performance, especially in scenarios with small local training datasets and increased data heterogeneity, as observed 718 in (Wu et al., 2022). 719

720

721 Potential of improving the trade-off of local and global performance from the perspective of 722 loss landscape. In order to better understand the dilemma of simultaneously improving local and global performance, we provide a novel perspective based on the analysis of personalized overfitting 723 using loss landscape. We believe that the fundamental reason for overfitting in previous PFL methods 724 is the inability to find a wide low-loss basin influenced by both local and global losses. Previous GFL 725 methods focus solely on finding the optimal solution in the global loss landscape, while PFL methods 726 only focus on the optimal solution in the local loss landscape, neglecting the structural information 727 of the combined loss landscape composed of the two losses. Due to lacking consideration of the 728 offset between the global and local loss landscapes caused by data heterogeneity, even if previous 729 PFL methods can find a wide low-loss basin in the local loss landscape, they still cannot guarantee 730 generalization to other clients. That is also the reason why personalization overfitting emerges.

731 732 733

734

### **B** FINE-TUNE DETAILS

This section provides an overview of the baseline techniques utilized in our study. We describe the
characteristics and implementation specifics of three main fine-tuning methods: Proximal FT, Soup
FT, and Sparse FT.

738 739

740

748

749

B.1 PROXIMAL FT

Proximal Fine-Tuning (Proximal FT) (Li et al., 2020b) is a method that emphasizes preserving the
 original knowledge of the pre-trained model while adapting it to new tasks. This technique employs
 proximal regularization, which penalizes large deviations from the initial model parameters during the
 fine-tuning process. The primary advantage of Proximal FT is its ability to maintain the generalization
 capabilities of the pre-trained model, thus reducing the risk of overfitting to the new task's data. In
 our experiments, we used an L2 regularization term to enforce proximity between the pre-trained and
 fine-tuned weights, with a regularization coefficient of 0.01.

B.2 SOUP FT

Soup Fine-Tuning (Soup FT) (Wortsman et al., 2022a) is an innovative approach that leverages the concept of "model soups," where multiple fine-tuned models are combined to create a more robust final model. The key idea is to fine-tune several instances of the pre-trained model on the target task with different random initializations or data shuffling, and then average the resulting weights to form a "soup." This method aims to enhance model robustness and performance by integrating the strengths of various fine-tuning instances. For our implementation, we fine-tuned five versions of the pre-trained model and averaged their parameters to create the final Soup FT model.

### 756 B.3 SPARSE FT

758 Sparse Fine-Tuning (Sparse FT) (Lee et al., 2018) introduces sparsity constraints into the fine-tuning 759 process, encouraging the model to update only a subset of its parameters. This approach aims to 760 improve model efficiency and interpretability by ensuring that only the most relevant weights are 761 adjusted during training. Sparse FT can be particularly beneficial for deploying models in resource-762 constrained environments where computational efficiency is paramount. In our experiments, we 763 applied L1 regularization to enforce sparsity, setting the regularization coefficient to 0.001 to achieve 764 a balance between performance and sparsity.

765 766

767

768

B.4 LP-FT

769 Linear Probing and then Fine-Tuning (LP-FT) (Kumar et al., 2022) is a two-step transfer learning 770 approach designed to balance in-distribution (ID) and out-of-distribution (OOD) performance. In 771 the first step, linear probing trains only the final layer (head) while freezing the pretrained feature extractor to ensure OOD robustness. The second step fine-tunes all model parameters to improve ID 772 accuracy while retaining the benefits of linear probing for OOD generalization. LP-FT addresses the 773 trade-offs inherent in full fine-tuning by initializing with a well-aligned linear head, reducing feature 774 distortion during optimization. Empirically, LP-FT demonstrates superior ID and OOD accuracy 775 across diverse datasets. 776

- 777
- 778 779

780 781

782

### C EXPERIMENTAL DETAILS

C.1 COMPUTING ENVIRONMENT AND HYPER-PARAMETERS

All experiments in this paper are conducted on NVIDIA A40 Graphics cards using PyTorch. The Adam optimizer is employed with a learning rate of  $1 \times 10^{-3}$ . In FL for all datasets, the standard local model update epochs are set to 1. The communication round is set to be 100 epochs, where we validated the model results from FL converged. Unless specified otherwise, the batch size for all benchmarks is standardized at 128. To ensure a fair comparison with various baselines, all methods initiate the FL personalized fine-tuning with models derived from the best-performing global model in terms of overall effectiveness.

790 791 792

793

C.2 VISUALIZATION OF THE ORIGINAL IMAGES

Fig. 5 illustrates a visual representation of the various datasets used in this study, categorized by their
levels of transformation or domain. The figure is divided into three main sections:

(a) *Feature-Level Shift* (*Digit5 and DomainNet*): The left panel displays examples from the Digit5 dataset, showcasing digit images in diverse styles and appearances. These include handwritten digits, digits rendered in varying fonts, and those with unique textures. The right panel features images from the DomainNet dataset, which includes objects and scenes represented in various artistic styles such as clip art, sketches, and realistic photographs. Examples include a strawberry, a zebra, and a cat.

(b) *Input-Level Shift (CIFAR10-C and CIFAR100-C):* This section highlights images from CIFAR10-C and CIFAR100-C, which apply corruptions to standard CIFAR datasets to evaluate robustness.
Corruptions include noise, blurring, and distortions that affect the clarity and quality of the images. Example categories feature animals, vehicles, and natural scenes under different types of degradation.

806

(c) *Output-Level Shift (CheXpert and CelebA):* The left panel presents grayscale X-ray images from
 the CheXpert dataset, widely used in medical imaging tasks. The right panel showcases color images
 of celebrity faces from the CelebA dataset, commonly used for facial recognition and attribute
 prediction.



### C.3 DETAILED DATASET AND MODEL INFORMATION

**Tab. 3** provides a visual overview of the datasets used in this study, categorized by their levels of transformation and data heterogeneity. The table is divided into four sections, corresponding to feature-level, input-level, output-level, and label shift settings:

*Feature-Level Shift (Digit5 and DomainNet):* The **Digit5** dataset, which consists of digit images collected from five distinct domains, including MNIST, SVHN, USPS, SynthDigits, and MNIST-M.

These images exhibit a variety of styles, such as handwritten digits, digits rendered in different fonts, and textured representations, demonstrating substantial visual heterogeneity. The **DomainNet** dataset, a large-scale collection featuring objects and scenes from six domains, including styles like clip art, sketches, and realistic photographs.

Input-Level Shift (CIFAR10-C and CIFAR100-C): This type of distribution shift includes corrupted versions of CIFAR-10 and CIFAR-100 datasets. The CIFAR10-C dataset applies 50 types of corruptions, such as noise, blur, and distortions, to evaluate model robustness under various degradation conditions. Similarly, CIFAR100-C extends the CIFAR-100 dataset by introducing the same set of corruptions, enabling robustness evaluation on a larger and more diverse set of categories.

| Data Heterogeneity<br>Type | Dataset    | Model    | Description  | Clients | Classes |
|----------------------------|------------|----------|--|---------|---------|
| Feature-Level Shift        | Digit5     | ResNet18 | A collection of digit images from five domains, used for<br>domain adaptation and digit recognition tasks (Huang<br>et al., 2023). Datasets include MNIST, SVHN, USPS,<br>SynthDigits, and MNIST-M.  | 5       | 10      |
|                            | DomainNet  | ResNet50 | A large-scale dataset of images from six distinct domains<br>for multi-source domain adaptation (Peng et al., 2019b).<br>Preprocessing follows the strategy in FedBN (Li et al.,<br>2021c).  | 6       | 10      |
| Input-Level Shift          | CIFAR10-C  | ResNet18 | A corruption benchmark dataset for<br>CIFAR-10 (Hendrycks & Dietterich, 2019), augmented<br>with 30 additional corruption types (Mintun et al., 2021)<br>and one extra type (Chen et al., 2021), totaling 50<br>corruption types (including the original 19 types of<br>corruption from CIFAR-10-C, plus 30 additional<br>corruption types and one extra type, resulting in a total of<br>50 distinct corruption types). | 50      | 10      |
|                            | CIFAR100-C | ResNet18 | An extension of CIFAR-100 with common corruptions, following the same strategy as CIFAR-10-C.  | 50      | 100     |
| Output_Level Shift         | CheXpert   | ResNet50 | A chest radiograph dataset labeled for 14 common chest<br>conditions (Irvin et al., 2019). Edema and No Finding<br>labels are grouped as described in (Jin et al., 2024).<br>Clients are spuriously correlated with the attribute gender<br>(e.g., 90% of label 1 examples in a client are male).  | 2       | 2       |
| output-Level Smit          | CelebA     | ResNet50 | Over 200,000 celebrity images with 40 attributes (Liu et al., 2015). Client splitting follows the same strategy as CheXpert, with attributes: male, female, blonde hair, and non-blonde hair.  | 4       | 2       |
| Label Shift                | CIFAR10    | ResNet18 | A benchmark dataset with label shift induced via<br>Dirichlet distribution ( $\alpha = 0.1$ ), distributed across 20<br>clients (Krizhevsky et al., 2009).   | 20      | 10      |

Table 3: Detailed information about the datasets and their splitting strategies used in the study. For all the settings above, each client has an individual data distribution, ensuring the non-IID nature required for heterogeneous federated learning. Feature-level shift, also referred to as subgroup shift, and input-level shift, corresponding to image corruption, are categorized as covariate shift. Output-level shift, representing spurious correlations in our setting, is categorized as concept shift.

Output-Level Shift (CheXpert and CelebA): The CheXpert dataset, a widely used medical imaging dataset labeled for 14 common chest conditions. In this study, the Edema and No Finding labels are grouped, and spurious correlations are introduced at the client level, where attributes (i.e. gender in our client splitting) are disproportionately represented (i.e., 90% of label 1 examples in a client are a certain attribute). The CelebA dataset, which includes over 200,000 celebrity faces annotated with 40 attributes. Client splitting follows the same strategy as CheXpert, with attributes such as male, female, blonde hair, and non-blonde hair used to create spurious correlations across clients. 

Label Shift (CIFAR10): The final one focuses on label shift in the CIFAR10 dataset. This setting simulates non-IID conditions by inducing label distributions across clients using a Dirichlet distribution with  $\alpha = 0.1$ . This creates significant variations in class distributions among the 20 clients, mimicking real-world federated learning scenarios where data availability across clients is inherently imbalanced.

### D FURTHER EMPIRICAL RESULTS

### D.1 OTHER PEFT METHODS DISTORT FEDERATED FEATURE

| PEFT Method       | Local | Global | Avg.  |
|-------------------|-------|--------|-------|
| LoRA (lr=1e-3)    | 41.54 | 26.75  | 26.87 |
| LoRA (lr=1e-4)    | 42.01 | 26.41  | 27.18 |
| Adapter (lr=1e-3) | 48.35 | 39.28  | 38.08 |
| Adapter (lr=1e-4) | 49.07 | 39.83  | 38.36 |
| LP-FT             | 68.50 | 57.52  | 53.52 |

Table 4: Different PEFT compared on DomainNet with ViT.

**Setup.** In this study, we adapted two widely recognized parameter-efficient fine-tuning (PEFT) methods: LoRA (Hu et al., 2022) and Adapter (Houlsby et al., 2019). Both methods were fine-tuned with meticulous adjustments to their learning rates on the ViT, as detailed in Tab. 4. These configurations allowed us to assess the impact of different fine-tuning strategies on federated features. The performance of each method was measured in terms of local, global, and average accuracy.

**Result.** The effectiveness of bias tuning in personalized fine-tuning naturally raises the question of whether other PEFT methods, commonly used for fine-tuning large models, exhibit similar effectiveness in our setting. In this study, we compare the local and global performance of other popular PEFT methods. Our findings reveal that while these methods can achieve high levels of local performance, their global performance still drops significantly, indicating that they distort the federated features to a certain extent. These PEFT methods' local and global performance still fall short compared to LP-FT, indicating that they distort the federated features to a certain extent.

The comparison of PEFT methods in the context of personalized fine-tuning sheds light on the unique
challenges and requirements of this setting. Despite the success of PEFT techniques in fine-tuning
large models for various tasks, our results suggest that their direct application to personalized FL
may not yield optimal results in terms of preserving the global knowledge captured by the federated
features.

### D.2 PFL RESULTS

**Setup.** To further demonstrate the importance of preserving federated features, we compare LP-FT with other popular personalized FL methods. Our primary focus in this paper is on the training during the FT phase. Unless otherwise specified, we consistently use models trained with FedAvg for FL. Unlike other personalized FL methods, which often involve additional operations during the FL phase (such as local training or model aggregation), our LP-FT method relies solely on vanilla FedAvg.

965 Results. From the table, it is evident that even without additional operations, LP-FT, which focuses
 966 on preserving federated features, remains highly competitive compared to other personalized FL
 967 methods. This observation underscores the significance of retaining the global features learned during
 968 the FL process. By directly comparing LP-FT with methods that employ extra techniques during
 969 the federated training phase, we demonstrate that the simple yet effective approach of preserving
 970 federated features can yield competitive results. This finding suggests that the key to achieving strong
 971 personalization lies in maintaining the knowledge acquired through collaboration across multiple
 972 clients, rather than relying on complex modifications to the FL algorithm.

| 972 | DEL M. Al J                   | T 1 A (A)               | $C_{1} = 1 + 1 + \dots + (A)$ | A . A     |
|-----|-------------------------------|-------------------------|-------------------------------|-----------|
| 973 | PFL Method                    | Local Acc. $(\uparrow)$ | Global Acc. (†)               | Avg. Acc. |
| 974 | FedBN (Li et al., 2021b)      | 88.68                   | 61.35                         | 61.00     |
| 975 | PerAvg (Fallah et al., 2020)  | 87.06                   | 67.26                         | 66.66     |
| 976 | FedNova (Wang et al., 2020)   | 88.68                   | 54.26                         | 53.41     |
| 077 | FedRep (Collins et al., 2021) | 86.94                   | 52.99                         | 52.57     |
| 977 | FedSoup (Chen et al., 2023)   | 90.30                   | 75.62                         | 75.21     |
| 978 | LP-FT                         | 93.03                   | 82.46                         | 82.17     |
| 979 |                               |                         |                               |           |

Table 5: Different PFL Methods on CelebA.

### D.3 LABEL SHIFT

We simulated label shift using a Dirichlet distribution with an alpha parameter set to 0.1. The dataset was distributed across 20 clients, and we trained a ResNet18 model for classification. Results is shown in Tab. 6.

| 909 |         | Baseline    | Local        | Global       | C-Std.       | Worst       | Avg          |
|-----|---------|-------------|--------------|--------------|--------------|-------------|--------------|
| 991 |         | FT          | 87.20 (0.24) | 16.67 (0.14) | 26.81 (1.47) | 0.01 (0.00) | 34.62 (0.09) |
| 992 | CIFAR10 | Proximal FT | 89.80 (1.21) | 17.42 (1.46) | 27.31 (0.08) | 0.01 (0.00) | 35.75 (0.09) |
| 993 |         | Soup FT     | 88.16 (0.15) | 16.94 (0.05) | 27.07 (0.11) | 0.01 (0.00) | 35.04 (0.04) |
| 994 |         | Sparse FT   | 89.16 (0.00) | 17.54 (0.14) | 27.27 (0.02) | 0.01 (0.00) | 35.57 (0.04) |
| 995 |         | LP-FT       | 90.15 (0.43) | 17.73 (0.16) | 27.37 (0.09) | 0.01 (0.00) | 35.96 (0.10) |

Table 6: Comparison of PFT methods on CIFAR10 label shift setting.

Tab. 6 compares different fine-tuning methods on CIFAR10 across multiple evaluation metrics, including local performance, global performance, robustness to corruption (C-Std.), worst-case performance, and average performance. Among the methods, LP-FT achieves the best results, excelling in local performance (90.15), global evaluation (17.73), and average performance (35.96). Proximal FT and Sparse FT also perform competitively, with improvements over standard fine-tuning (FT) and Soup FT in most metrics. All methods show near-identical performance in the worst-case scenario (0.01), indicating a shared limitation in extreme cases. Overall, LP-FT demonstrates the most robust and effective fine-tuning approach on CIFAR10. 

#### Ε LIMITATIONS AND BROADER IMPACT.

In our paper, we focused exclusively on vision-related tasks. Extending our empirical findings to language tasks or multimodal scenarios would be a promising direction for future research. Our FL study primarily addresses simulated cross-silo FL settings. Validating our conclusions in real-world FL deployments would also be a worthwhile direction for future exploration. Deploying FL that balances personalization and generalization capabilities in healthcare scenarios holds great promise. However, excessive personalization could lead to issues with the fairness of algorithmic decisions. 

F PROOFS

*Proof of Lemma 4.3.* We want to analyze the Fine-Tuning (FT) method, focusing on the effect of initial parameters. We perform one pass through the entire dataset to simulate the complete fine-tuning process. Consider the Mean Squared Error (MSE) loss function with parameters V and B, where B is represented as follows:

 1026

 1027

 1028

 1029

 1030

 1031

 1032

 1033

1034 where  $b_i^T \in \mathbb{R}^{1 \times d}$  denotes the *i*-th row of matrix *B*, and m + C = k.

To apply one step of gradient descent, we need to compute the gradient of the loss function with respect to  $V, b_1, b_2, \ldots, b_{m+C}$ , and then perform one update step.

W.L.O.G. we assume the local client is client 1. We define the local loss function as follows:

$$\mathcal{L}_L(V,B) = \mathbb{E}_{x \sim \mathcal{D}_1} \left[ \frac{1}{2} (V^T B x - V_1^{*T} B_* x)^2 \right],$$

1042 where  $\mathcal{D}_1$  is the data distribution for client 1.

Now, let  $(\mathbf{X}_1, \mathbf{Y}_1)$  represent the local dataset of client 1, consisting of  $n_1$  data points  $\{(x_{1j}, y_{1j})\}_{j=1}^{n_1}$ . We aim to calculate the gradient of the empirical loss function with respect to the parameters. The empirical loss function is given by:

$$\widehat{\mathcal{L}}_L(V,B) = \frac{1}{n_1} \sum_{j=1}^{n_1} \left[ \frac{1}{2} (V^T B x_{1j} - V_1^{*T} B_* x_{1j})^2 \right].$$

1050 In practice, we take the gradient of this empirical loss function with respect to the parameters  $V, b_1$ , 1051  $b_2, \ldots, b_{m+C}$ . However, since we are particularly interested in computing the expectation  $\mathbb{E}[b_j^{FT}]$ , 1052 we evaluate the expected value of the gradients using one pass through the whole dataset as follows:

$$\mathbb{E}\left[\left.\frac{\partial\widehat{\mathcal{L}_L}}{\partial V}\right|_{\substack{V=V_0\\B=B_*}}\right] = \mathbb{E}\left[\left.\frac{\partial}{\partial V}\left(\frac{1}{n_1}\sum_{j=1}^{n_1}\frac{1}{2}\left(V^TBx_{1j} - {V_1^*}^TB_*x_{1j}\right)^2\right)\right|_{\substack{V=V_0\\B=B_*}}\right]$$

 $= \frac{1}{n_1} \sum_{j=1} \mathbb{E} \left[ (V_0^T B_* x_{1j} - V_1^{*T} B_* x_{1j}) x_{1j}^T B_*^T \right]$  $= \mathbb{E}_{x \sim \mathcal{D}_1} \left[ (V_0^T B_* x - V_1^{*T} B_* x) x^T B_*^T \right].$ 

1063 Therefore, let  $V_0 = \begin{bmatrix} V_{com}^* & \mathbf{0}^T \end{bmatrix}^T$ . It follows that:

$$\mathbb{E}\left[\frac{\partial L}{\partial V}\Big|_{\substack{V=V_0\\B=B_*}}\right] = \mathbb{E}_{x\sim\mathcal{D}_1}\left[\left(V_0^T B_* x - V_1^{*T} B_* x\right) x^T B_*^T\right]$$
$$= \mathbb{E}_{x\sim\mathcal{D}_1}\left[\left(\left(V_0 - V_1^*\right)^T B_* x\right) x^T B_*^T\right]$$
$$= \left(V_0 - V_1^*\right)^T B_* \left(\mathbb{E}_{x\sim\mathcal{D}_1}\left[xx^T\right]\right) B_*^T$$
$$= \left(V_0 - V_1^*\right)^T B_* B_*^T$$

 $= (V_0 - V_1^*)^T$ .

(second moment is identity) (B<sub>\*</sub> has orthonormal rows)

1073  
1074  
1075  
1076  
1077  
1078  
1079  
Let 
$$B_* = \begin{bmatrix} b_1^{*T} \\ \vdots \\ b_m^{*T} \\ b_{m+1}^{*} \\ \vdots \\ b_{m+C}^{*} \end{bmatrix}$$
. Then, similarly, it can be shown that:

 $\mathbb{E}\left[\left.\frac{\partial L}{\partial b_j}\right|_{\substack{V=V_0\\B=B_-}}\right] = \mathbb{E}_{x\sim\mathcal{D}_1}\left[\left(V_0^T B_* x - V_1^{*T} B_* x\right)(V_0)_j x^T\right]$ 

 $= (V_0)_i (V_0 - V_1^*)^T B_*.$ 

 $= \mathbb{E}_{x \sim \mathcal{D}_1} \left[ (V_0)_j \left( (V_0 - V_1^*)^T B_* x \right) x^T \right]$ 

 $= (V_0)_j (V_0 - V_1^*)^T B_* \left( \mathbb{E}_{x \sim \mathcal{D}_1} \left[ x x^T \right] \right)$ 

Here,  $(V_0)_j$  is the *j*-th element of the vector  $V_0$ . For learning rate  $\eta$ , one step of gradient descent is:

(second moment is identity)

1094  
1095  
1096  
1097  
1098  
1099  
1100  

$$V_{FT} = V_0 - \eta \left( \frac{\partial L}{\partial V} \Big|_{\substack{V = V_0 \\ B = B_*}} \right)^T$$

$$b_j^{FT} = b_j^* - \eta \left( \frac{\partial L}{\partial b_j} \Big|_{\substack{V = V_0 \\ B_0 = B_*}} \right)^T.$$

1101 These two equations can be further refined as:

$$\mathbb{E}\begin{bmatrix} V_{FT} \end{bmatrix} = \begin{bmatrix} V_{com}^* & \mathbf{0}^T \end{bmatrix}^T - \eta(V_0 - V_1^*) = \begin{bmatrix} V_{com}^* & \eta\lambda e_1^T \end{bmatrix}^T$$
$$\mathbb{E}\begin{bmatrix} b_j^{FT} \end{bmatrix} = b_j^* - \eta \left( \left. \frac{\partial L}{\partial b_j} \right|_{\substack{V=V_0\\B_0=B_*}} \right)^T = b_j^* - \eta \left( (V_0)_j (V_0 - V_1^*)^T B_* \right)^T$$

$$= b_{j}^{*} - \eta \lambda \left( (V_{0})_{j} \begin{bmatrix} \mathbf{0}^{T} & -e_{1}^{T} \end{bmatrix} B_{*} \right)^{T} = b_{j}^{*} + \eta \lambda (V_{0})_{j} b_{m+1}^{*}.$$

Therefore, we have:

$$\mathbb{E}\left[B_{FT}\right] = \begin{bmatrix} b_{1}^{*T} + \eta\lambda(V_{0})_{1}b_{m+1}^{*}^{T} \\ \vdots \\ b_{m}^{*T} + \eta\lambda(V_{0})_{m}b_{m+1}^{*}^{T} \\ b_{m+1}^{*} + \eta\lambda(V_{0})_{m+1}b_{m+1}^{*}^{T} \\ \vdots \\ b_{m+C}^{*T} + \eta\lambda(V_{0})_{m+C}b_{m+1}^{*}^{T} \end{bmatrix} = \begin{bmatrix} b_{1}^{*T} + \eta\lambda(V_{0})_{1}b_{m+1}^{*}^{T} \\ \vdots \\ b_{m}^{*T} + \eta\lambda(V_{0})_{m+1}b_{m+1}^{*} \\ \vdots \\ b_{m+C}^{*T} + \eta\lambda(V_{com})_{1}b_{m+1}^{*}^{T} \\ \vdots \\ b_{m+1}^{*T} \\ b_{m+1}^{*T} \\ \vdots \\ b_{m+1}^{*T} \end{bmatrix}$$

Similarly, if the fine-tuning is done over the data of client *i*, we would have:

$$\mathbb{E}\left[b_{j}^{FT}\right] = b_{j}^{*} + \eta \lambda(V_{0})_{j} b_{m+i}^{*},$$

which concludes the proof.

*Proof of Theorem 4.4.* We assume that the pre-trained model perfectly captures the feature extractor 1133 matrix  $B_*$ , and its linear head represents the common part shared across all clients, excluding any client-specific components of the ground-truth function. Thus,  $B_0 = B_*$  and  $V_0 = \begin{bmatrix} V_{com}^* & \mathbf{0} \end{bmatrix}^T$ .

In this setting, we analyze the effects of LP-FT and FT on the model parameters. For both LP-FT and FT, we determine the parameters after fine-tuning, compute the global loss, and then compare these global losses.

1137 1138 W.L.O.G. we assume that we are doing the fine-tuning over the local data of client 1. First, we study 1139 LP-FT. Initially, one step of linear probing is conducted with the fixed feature extractor  $B_*$ . After 1140 this step, the linear head  $V_{LP}$  will converge to  $V_1^*$ . This is because we know that:

1141 1142

1143

1146 1147 1148

$$\arg\min_{v} \left\| \mathbf{X} B_0^\top v - \mathbf{X} B_*^\top v_* \right\|_2^2 = \left( B_0 \mathbf{X}^\top \mathbf{X} B_0^\top \right)^{-1} B_0 \mathbf{X}^\top \mathbf{X} B_*^\top v_*$$

where X is the  $n \times d$  matrix including data of n individuals. Since the fine-tuning is on the data of the client 1 (local data), we have:

$$V_{LP} = \left(B_0 \mathbf{X_1}^{\top} \mathbf{X_1} B_0^{\top}\right)^{-1} B_0 \mathbf{X_1}^{\top} \mathbf{X_1} B_*^{\top} V_1^*$$

1149 Therefore, we have:

1150  
1151  
1152  
1153  
1154  
1155  

$$V_{LP} = \left(B_0 \mathbf{X}_1^{\top} \mathbf{X}_1 B_0^{\top}\right)^{-1} B_0 \mathbf{X}_1^{\top} \mathbf{X}_1 B_*^{\top} V_1^*$$

$$= \left(B_* \mathbf{X}_1^{\top} \mathbf{X}_1 B_*^{\top}\right)^{-1} B_* \mathbf{X}_1^{\top} \mathbf{X}_1 B_*^{\top} V_1^*$$

$$= V_1^*.$$

1156 1157 Since at the beginning of the fine-tuning (FT) step, we have the perfect  $B_*$  and  $V_1^*$  for the local client 1, 1158 and FT is performed on the data of the same client, we can conclude that after one step of FT following 1159 LP, the parameters will remain unchanged. Specifically, we have  $V_{LPFT} = V_1^* = \begin{bmatrix} V_{com}^* & e_1^T \end{bmatrix}^T$ 1160 and  $B_{LPFT} = B_*$ .

1161 For the performance on the global data, we have:

$$\begin{aligned} \mathcal{L}_{G}(V_{LPFT}, B_{LPFT}) &= \frac{1}{C} \sum_{i \in [C]} \mathbb{E}_{x \sim \mathcal{D}_{i}} \left[ \frac{1}{2} (V_{LPFT}^{T} B_{LPFT} x - V_{i}^{*T} B_{*} x)^{2} \right] \\ &= \frac{1}{C} \sum_{i \in [C]} \mathbb{E}_{x \sim \mathcal{D}_{i}} \left[ \frac{1}{2} (V_{1}^{*T} B_{*} x - V_{i}^{*T} B_{*} x)^{2} \right] \\ &= \frac{1}{C} \sum_{i \in [C]} \mathbb{E}_{x \sim \mathcal{D}_{i}} \left[ (B_{*}^{T} V_{1}^{*} - B_{*}^{T} V_{i}^{*})^{T} x x^{T} (B_{*}^{T} V_{1}^{*} - B_{*}^{T} V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \mathbb{E}_{x \sim \mathcal{D}_{i}} \left[ (B_{*}^{T} V_{1}^{*} - B_{*}^{T} V_{i}^{*})^{T} \mathbb{E}_{x \sim \mathcal{D}_{i}} \left[ x x^{T} \right] (B_{*}^{T} V_{1}^{*} - B_{*}^{T} V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ (B_{*}^{T} V_{1}^{*} - B_{*}^{T} V_{i}^{*})^{T} \mathbb{E}_{x \sim \mathcal{D}_{i}} \left[ x x^{T} \right] (B_{*}^{T} V_{1}^{*} - B_{*}^{T} V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ (B_{*}^{T} V_{1}^{*} - B_{*}^{T} V_{i}^{*})^{T} (B_{*}^{T} V_{1}^{*} - B_{*}^{T} V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ (V_{1}^{*} - V_{i}^{*})^{T} B_{*} B_{*}^{T} (V_{1}^{*} - B_{*}^{T} V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ (V_{1}^{*} - V_{i}^{*})^{T} I_{k} (V_{1}^{*} - V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ (V_{1}^{*} - V_{i}^{*})^{T} (V_{1}^{*} - V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ (V_{1}^{*} - V_{i}^{*})^{T} (V_{1}^{*} - V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ \left[ (V_{1}^{*} - V_{i}^{*})^{T} (V_{1}^{*} - V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ \left[ (V_{1}^{*} - V_{i}^{*})^{T} (V_{1}^{*} - V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ \left[ (V_{1}^{*} - V_{i}^{*})^{T} (V_{1}^{*} - V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ \left[ (V_{1}^{*} - V_{i}^{*})^{T} (V_{1}^{*} - V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ \left[ \left[ (V_{1}^{*} - V_{i}^{*})^{T} (V_{1}^{*} - V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ \left[ \left[ (V_{1}^{*} - V_{i}^{*})^{T} (V_{1}^{*} - V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ \left[ \left[ (V_{1}^{*} - V_{i}^{*})^{T} (V_{1}^{*} - V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ \left[ \left[ (V_{1}^{*} - V_{i}^{*})^{T} (V_{1}^{*} - V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [C]} \left[ \left[ \left[ (V_{1}^{*} - V_{i}^{*})^{T} (V_{1}^{*} - V_{i}^{*}) \right] \\ &= \frac{1}{2C} \sum_{i \in [$$

$$= \frac{1}{2C} \sum_{\substack{i \in [C] \\ i \neq 1}} \left[ \left\| \left( \begin{bmatrix} V_{com}^* & \lambda e_1^T \end{bmatrix}^T - \begin{bmatrix} V_{com}^* & \lambda e_i^T \end{bmatrix}^T \right) \right\|_2^2 \right]$$

$$= \left(\frac{1}{2C}\right) 2(C-1) = \lambda^2 \frac{C-1}{C}.$$

$$(1)$$

1195 It can be shown that:

$$\mathcal{L}_{G}(V_{FT}, B_{FT}) = \frac{1}{C} \sum_{i \in [C]} \mathbb{E}_{x \sim \mathcal{D}_{i}} \left[ \frac{1}{2} (V_{FT}^{T} B_{FT} x - V_{i}^{*T} B_{*} x)^{2} \right]$$

$$= \frac{1}{C} \sum_{i \in [C]} \mathbb{E}_{x \sim \mathcal{D}_{i}} \left[ (B_{FT}^{T} V_{FT} - B_{*}^{T} V_{i}^{*})^{T} x x^{T} (B_{FT}^{T} V_{FT} - B_{*}^{T} V_{i}^{*}) \right]$$

$$= \frac{1}{2C} \sum_{i \in [C]} (B_{FT}^{T} V_{FT} - B_{*}^{T} V_{i}^{*})^{T} \left[ \mathbb{E}_{x \sim \mathcal{D}_{i}} x x^{T} \right] (B_{FT}^{T} V_{FT} - B_{*}^{T} V_{i}^{*})$$

$$= \frac{1}{2C} \sum_{i \in [C]} (B_{FT}^{T} V_{FT} - B_{*}^{T} V_{i}^{*})^{T} (B_{FT}^{T} V_{FT} - B_{*}^{T} V_{i}^{*})$$

$$= \frac{1}{2C} \sum_{i \in [C]} (B_{FT}^{T} V_{FT} - B_{*}^{T} V_{i}^{*})^{T} (B_{FT}^{T} V_{FT} - B_{*}^{T} V_{i}^{*})$$

$$= \frac{1}{2C} \sum_{i \in [C]} \| (B_{FT}^{T} V_{FT} - B_{*}^{T} V_{i}^{*}) \|_{2}^{2}.$$

$$(2)$$

We have:

$$B_*^T V_i^* = \sum_{j=1}^m (V_{com}^*)_j b_j^* + \lambda b_{m+i}^*$$
$$B_{FT}^T V_{FT} = \sum_{j=1}^m (V_{com}^*)_j b_j^* + \sum_{j=1}^m \eta \lambda (V_{com}^*)_j^2 b_{m+1}^* + \eta \lambda b_{m+1}^*.$$

1218 Therefore, we can obtain:

$$(B_{FT}^T V_{FT} - B_*^T V_i^*) = \lambda \Big( \sum_{j=1}^m \eta (V_{com}^*)_j^2 b_{m+1}^* + \eta b_{m+1}^* - b_{m+i}^* \Big).$$

For  $i \neq 1$ , we have:

$$(B_{FT}^T V_{FT} - B_*^T V_i^*)^T (B_{FT}^T V_{FT} - B_*^T V_i^*)$$
  
=  $\lambda^2 (\sum_{j=1}^m \eta(V_{com}^*)_j^2 b_{m+1}^* + \eta b_{m+1}^* - b_{m+i}^*)^T (\sum_{j=1}^m \eta(V_{com}^*)_j^2 b_{m+1}^* + \eta b_{m+1}^* - b_{m+i}^*)$   
=  $\lambda^2 \left( \left( \eta + \eta \sum_{j=1}^m (V_{com}^*)_j^2 \right)^2 + 1 \right).$  (rows of  $B_*$  are orthonormal)

1233 For i = 1, we have:

$$\begin{aligned} & (B_{FT}^{T}V_{FT} - B_{*}^{T}V_{i}^{*})^{T}(B_{FT}^{T}V_{FT} - B_{*}^{T}V_{i}^{*}) \\ & = \lambda^{2} \Big( \sum_{j=1}^{m} \eta(V_{com}^{*})_{j}^{2}b_{m+1}^{*} + \eta b_{m+1}^{*} - b_{m+i}^{*} \Big)^{T} \Big( \sum_{j=1}^{m} \eta(V_{com}^{*})_{j}^{2}b_{m+1}^{*} + \eta b_{m+1}^{*} - b_{m+i}^{*} \Big) \\ & = \lambda^{2} \Big( \eta + \eta \sum_{j=1}^{m} (V_{com}^{*})_{j}^{2} - 1 \Big)^{2}. \end{aligned}$$
 (rows of  $B_{*}$  are orthonormal)  
1241

Combining these with (2), we can conclude: 

1244  
1245  
1246  
1246  
1247  
1248  
1249  
1249  
1250  

$$\mathcal{L}_{G}(V_{FT}, B_{FT}) = \frac{1}{2C} \sum_{i \in [C]} \left\| (B_{FT}^{T} V_{FT} - B_{*}^{T} V_{i}^{*}) \right\|_{2}^{2}$$

$$= \frac{\lambda^{2}}{2C} \left( \left( \eta + \eta \sum_{j=1}^{m} (V_{com}^{*})_{j}^{2} - 1 \right)^{2} + (C - 1) \left( \left( \eta + \eta \sum_{j=1}^{m} (V_{com}^{*})_{j}^{2} \right)^{2} + 1 \right) \right).$$
(3)

Combining (1) and (3), we have:

$$\mathcal{L}_G(V_{LPFT}, B_{LPFT}) \le \mathcal{L}_G(V_{FT}, B_{FT}).$$

Proof of Theorem 4.5. W.L.O.G. we assume that the local fine-tuning is performed on the data of the first client. Initially, one step of linear probing is conducted with the fixed feature extractor  $B_*$ . After this step, the linear head  $V_{LP}$  will converge to  $V_1^*$ . This is because we know that: 

$$\arg\min_{v} \left\| \mathbf{X}_{1} B_{0}^{\top} v - \mathbf{X}_{1} B_{*}^{\top} v_{*} \right\|_{2}^{2} = \left( B_{0} \mathbf{X}_{1}^{\top} \mathbf{X}_{1} B_{0}^{\top} \right)^{-1} B_{0} \mathbf{X}_{1}^{\top} \mathbf{X}_{1} B_{*}^{\top} v_{*}$$

where  $\mathbf{X}_1$  is the  $n \times d$  matrix including data of n individuals. Since the fine-tuning is on the data of the client 1 (local data), we have: 

$$V_{LP} = \left(B_0 \mathbf{X_1}^{\top} \mathbf{X_1} B_0^{\top}\right)^{-1} B_0 \mathbf{X_1}^{\top} \mathbf{X_1} B_*^{\top} V_1^*$$

Therefore, we have:

1268  
1269  
1270  
1271  
1272  
1273  

$$V_{LP} = \left(B_0 \mathbf{X}_1^{\top} \mathbf{X}_1 B_0^{\top}\right)^{-1} B_0 \mathbf{X}_1^{\top} \mathbf{X}_1 B_*^{\top} V_1^*$$

$$= \left(B_* \mathbf{X}_1^{\top} \mathbf{X}_1 B_*^{\top}\right)^{-1} B_* \mathbf{X}_1^{\top} \mathbf{X}_1 B_*^{\top} V_1^*$$

$$= V_1^*.$$

This part is identical to the initial part of the proof of Theorem 4.4. Since at the beginning of the fine-tuning (FT) step, we have the perfect  $B_*$  and  $V_1^*$  for the local client 1, and FT is performed on the data of the same client, we can conclude that after one step of FT following LP, the parameters will remain unchanged. Specifically, we have  $V_{LPFT} = V_1^* = \begin{bmatrix} V_{com}^* & \lambda e_1^T \end{bmatrix}^T$  and  $B_{LPFT} = B_*$ . e: 

$$\mathcal{L}_G(V_{LPFT}, B_{LPFT}) = \frac{1}{C} \sum_{i \in [C]} \mathbb{E}_{x \sim \mathcal{D}_i} \left[ \frac{1}{2} (V_{LPFT}^T B_{LPFT} x - V_i^{*T} B_* x)^2 \right]$$

$$= \frac{1}{C} \sum_{i \in [C]} \mathbb{E}_{x \sim \mathcal{D}_i} \left[ \frac{1}{2} (V_1^{*T} B_* x - V_i^{*T} B_* x)^2 \right]$$

$$= \frac{1}{2C} \sum_{i \in [C]} \mathbb{E}_{x \sim \mathcal{D}_i} \left[ (B_*^T V_1^* - B_*^T V_i^*)^T x x^T (B_*^T V_1^* - B_*^T V_i^*) \right]$$

1288  
1289  
1290  
1290  

$$= \frac{1}{2C} \sum_{i \in [C]} \left[ (B_*^T V_1^* - B_*^T V_i^*)^T \mathbb{E}_{x \sim \mathcal{D}_i} [xx^T] (B_*^T V_1^* - B_*^T V_i^*) \right]$$

$$= \frac{1}{2C} \sum_{i \in [C]} \left[ (V_1^* - V_i^*)^T B_* \left( \mathbb{E}_{x \sim \mathcal{D}_i} \left[ x x^T \right] \right) B_*^T (V_1^* - V_i^*) \right]$$

1294  
1295 
$$= \frac{1}{2C} \sum_{i \in [C]} \left[ (V_1^* - V_i^*)^T B_* \Big( \mathbb{E}_{n \sim \mathcal{N}(0, I_d)} \big[ (e_i + \epsilon n) (e_i + \epsilon n)^T \big] \Big) B_*^T (V_1^* - V_i^*) \right]$$

$$\begin{aligned}
\begin{aligned}
1296 \\
1297 \\
1298 \\
= \frac{1}{2C} \sum_{i \in [C]} \left[ (V_1^* - V_i^*)^T B_* \left( e_i e_i^T + \epsilon^2 \mathbb{E}_{n \sim \mathcal{N}(0, I_d)} \left[ nn^T \right] \right) B_*^T (V_1^* - V_i^*) \right] \\
1299 \\
= \frac{1}{2C} \sum_{i \in [C]} \left[ (V_1^* - V_i^*)^T B_* \left( e_i e_i^T + \epsilon^2 I_d \right) B_*^T (V_1^* - V_i^*) \right] \\
1301 \\
= \frac{1}{2C} \sum_{i \in [C]} \left[ (V_1^* - V_i^*)^T B_* \left( e_i e_i^T \right) B_*^T (V_1^* - V_i^*) \right] \\
1304 \\
+ \frac{1}{2C} \sum_{i \in [C]} \left[ (V_1^* - V_i^*)^T B_* \left( \epsilon^2 I_d \right) B_*^T (V_1^* - V_i^*) \right] \\
1306 \\
= \frac{1}{2C} \sum_{i \in [C]} \left[ (V_1^* - V_i^*)^T (B_*)_{:,i} (B_*)_{:,i}^T (V_1^* - V_i^*) \right] \\
1307 \\
1308 \\
= \frac{1}{2C} \sum_{i \in [C]} \left[ \left( V_1^* - V_i^* \right)^T (V_1^* - V_i^*) \right] \\
1309 \\
= \frac{1}{2C} \sum_{i \in [C]} \left[ \left( (B_*)_{m+1,i} - (B_*)_{m+i,i} \right)^2 \right] + \frac{1}{2C} \epsilon^2 \sum_{i \in [C]} \left[ \left\| (V_1^* - V_i^*) \right\|_2^2 \right] \\
= \frac{\lambda^2}{2C} \sum_{i \in [C]} \left[ \left( (B_*)_{m+1,i} - (B_*)_{m+i,i} \right)^2 \right] + \frac{\lambda^2 (C-1)}{C} \epsilon^2. \end{aligned}$$
(4)

1319 We want to analyze the fine tuning (FT) method, focusing on the effect of initial parameters. We 1320 perform one pass through the entire dataset to simulate the complete fine-tuning process. Consider the Mean Squared Error (MSE) loss function with parameters V and B, where B is represented as 1321 follows: 1322

where  $b_i^T \in \mathbb{R}^{1 \times d}$  denotes the *i*-th row of matrix B, and m + C = k. 1332

1333 To apply one step of gradient descent, we need to compute the gradient of the loss function with 1334 respect to  $V, b_1, b_2, \ldots, b_{m+C}$ , and then perform one update step. 1335

,

Let 
$$V_0 = \begin{bmatrix} V_{com}^* & \mathbf{0}^T \end{bmatrix}^T$$
. It follows that:  

$$\mathbb{E} \begin{bmatrix} \frac{\partial L}{\partial V} \Big|_{\substack{V=V_0\\B=B_*}} \end{bmatrix} = \mathbb{E}_{x\sim\mathcal{D}_1} \begin{bmatrix} (V_0^T B_* x - V_1^{*T} B_* x) x^T B_*^T \end{bmatrix}$$

$$= \mathbb{E}_{x\sim\mathcal{D}_1} \begin{bmatrix} ((V_0 - V_1^*)^T B_* x) x^T B_*^T \end{bmatrix}$$

$$= (V_0 - V_1^*)^T B_* \left( \mathbb{E}_{x\sim\mathcal{D}_1} \begin{bmatrix} xx^T \end{bmatrix} \right) B_*^T$$

$$= (V_0 - V_1^*)^T B_* \left( \mathbb{E}_{n\sim\mathcal{N}(0,I_d)} \begin{bmatrix} (e_1 + \epsilon n)(e_1 + \epsilon n)^T \end{bmatrix} \right) B_*^T$$

$$= (V_0 - V_1^*)^T B_* \left( e_1 e_1^T + \epsilon^2 I_d \right) B_*^T$$

$$= (V_0 - V_1^*)^T B_* \left( e_1 e_1^T \right) B_*^T + (V_0 - V_1^*)^T B_* \left( \epsilon^2 I_d \right) B_*^T$$

$$= (V_0 - V_1^*)^T B_* \left( e_1 e_1^T \right) B_*^T + \epsilon^2 (V_0 - V_1^*)^T$$

$$(B_* \text{ has orthonormal rows)}$$

$$\begin{aligned} &= (V_0 - V_1^*)^T \left( (B_s)_{;,1} (B_s)_{;,1}^T \right) + \epsilon^2 (V_0 - V_1^*)^T \qquad ((B_s)_{;,1} \text{ is first column of } B_s) \\ &= \left( -\lambda (B_s)_{m+1,1} \right) (B_s)_{;,1}^T + \epsilon^2 (V_0 - V_1^*)^T. \\ &= \left( -\lambda (B_s)_{m+1,1} \right) (B_s)_{;,1}^T + \epsilon^2 (V_0 - V_1^*)^T. \\ &= \left[ \frac{b_1^* T}{b_{m,T}^* T} \right]. \text{ Then, it can be shown that:} \\ &= \left[ \frac{\partial L}{\partial b_j} \right|_{\substack{V=V_0 \\ B_m \in C^*}} \right] = \mathbb{E}_{x \sim \mathcal{D}_1} \left[ (V_0^T B_s x - V_1^{*T} B_s x) (V_0)_j x^T \right] \\ &= \mathbb{E}_{x \sim \mathcal{D}_1} \left[ (V_0)_j ((V_0 - V_1^*)^T B_s x) x^T \right] \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \mathbb{E}_{x \sim \mathcal{D}_1} [xT] \right) \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \mathbb{E}_{x \sim \mathcal{D}_1} [xT] \right) \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \exp^T + e^2 I_s \right) \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \exp^T + e^2 I_s \right) \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \exp^T + e^2 I_s \right) \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \exp^T + e^2 I_s \right) \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \exp^T + e^2 I_s \right) \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \exp^T + e^2 I_s \right) \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \exp^T + e^2 I_s \right) \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \exp^T + e^2 I_s \right) \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \exp^T + e^2 I_s \right) \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \exp^T + e^2 I_s \right) \\ &= (V_0)_j (V_0 - V_1^*)^T B_s \left( \exp^T - \frac{1}{2} \right) \\ &= b_j^* - \eta \left( \frac{\partial L}{\partial U_j} \right) \\ &= \left[ \frac{V_{emm}} T \quad 0^T \right]^T - \eta \left( -\lambda (B_s)_{m+1,1} (B_s)_{;,1} + e^2 (V_0 - V_1^*) \right) \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V_{emm}} B_s \left( \frac{D_s}{U_s} \right) \right] \\ &= \left[ \frac{V$$

Therefore, we have:  $(1 + \eta\lambda(V_0)_1\epsilon^2 b_{m+1}^* + \eta\lambda(V_0)_1(B_*)_{m+1,1}e_1^T)$  $\mathbb{E}[B_{FT}] = \begin{vmatrix} \vdots \\ b_m^* & T + \eta \lambda(V_0)_m \epsilon^2 b_{m+1}^* & T + \eta \lambda(V_0)_m (B_*)_{m+1,1} e_1^T \\ b_{m+1}^* & T + \eta \lambda(V_0)_{m+1} \epsilon^2 b_{m+1}^* & T + \eta \lambda(V_0)_{m+1} (B_*)_{m+1,1} e_1^T \end{vmatrix}$  $\vdots \\ b_{m+C}^* {}^T + \eta \lambda(V_0)_{m+C} \epsilon^2 b_{m+1}^* {}^T + \eta \lambda(V_0)_{m+C} (B_*)_{m+1,1} e_1^T$  $= \begin{bmatrix} b_1^{*T} + \eta \lambda(V_0)_1 \epsilon^2 b_{m+1}^{*T} + \eta \lambda(V_0)_1(B_*)_{m+1,1} e_1^T \\ \vdots \\ b_m^{*T} + \eta \lambda(V_0)_m \epsilon^2 b_{m+1}^{*T} + \eta \lambda(V_0)_m(B_*)_{m+1,1} e_1^T \\ b_{m+1}^{*T} \\ \vdots \\ b_{m+C}^{*T} \end{bmatrix}.$ It can be shown that:  $\mathcal{L}_G(V_{FT}, B_{FT}) = \frac{1}{C} \sum_{i \in [C]} \mathbb{E}_{x \sim \mathcal{D}_i} \left[ \frac{1}{2} (V_{FT}^T B_{FT} x - V_i^{*T} B_* x)^2 \right]$  $= \frac{1}{2C} \sum_{i \in [C]} \mathbb{E}_{x \sim \mathcal{D}_i} \left[ (B_{FT}^T V_{FT} - B_*^T V_i^*)^T x x^T (B_{FT}^T V_{FT} - B_*^T V_i^*) \right]$  $= \frac{1}{2C} \sum_{i \in [C]} (B_{FT}^T V_{FT} - B_*^T V_i^*)^T \bigg[ \mathbb{E}_{x \sim \mathcal{D}_i} x x^T \bigg] (B_{FT}^T V_{FT} - B_*^T V_i^*)$  $= \frac{1}{2C} \sum_{i \in [C]} (B_{FT}^T V_{FT} - B_*^T V_i^*)^T (e_i e_i^T + \epsilon^2 I_d) (B_{FT}^T V_{FT} - B_*^T V_i^*)$  $= \frac{1}{2C} \sum_{i \in [C]} (B_{FT}^T V_{FT} - B_*^T V_i^*)^T (e_i e_i^T) (B_{FT}^T V_{FT} - B_*^T V_i^*)$  $+ \frac{1}{2C} \sum_{i \in [C]} (B_{FT}^T V_{FT} - B_*^T V_i^*)^T (\epsilon^2 I_d) (B_{FT}^T V_{FT} - B_*^T V_i^*)$  $= \frac{1}{2C} \sum_{i \in [C]} (B_{FT}^T V_{FT} - B_*^T V_i^*)^T (e_i e_i^T) (B_{FT}^T V_{FT} - B_*^T V_i^*)$  $+\epsilon^2 \frac{1}{2C} \sum_{i \in [C]} \left\| (B_{FT}^T V_{FT} - B_*^T V_i^*) \right\|_2^2.$ 

We have:

 $B_{*}^{T}V_{i}^{*} = \sum_{i=1}^{m} (V_{com}^{*})_{j}b_{j}^{*} + \lambda b_{m+i}^{*}$ 

$$\begin{array}{ll} & 1453 \\ 1454 \\ 1455 \\ 1455 \\ 1456 \\ 1457 \end{array} + \sum_{j=1}^{m} \left( \gamma \lambda(B_{*})_{m+1,1}(B_{*})_{j,1} \right) b_{j}^{*} \\ & + \sum_{j=1}^{m} \left( (V_{com}^{*})_{j} + \eta \lambda(B_{*})_{m+1,1}(B_{*})_{j,1} \right) \left( b_{j}^{*} + \eta \lambda(V_{com}^{*})_{j} \epsilon^{2} \sigma^{2} b_{m+1}^{*} + \eta \lambda(V_{com}^{*})_{j} (B_{*})_{m+1,1} e_{1} \right) \\ \end{array}$$

(5)

1458 Therefore, we can obtain: 1459  $B_{FT}^{T}V_{FT} - B_{*}^{T}V_{i}^{*} = \sum_{j=1}^{m} (V_{com}^{*})_{j} \left(b_{j}^{*} + \eta\lambda(V_{com}^{*})_{j}\epsilon^{2}\sigma^{2}b_{m+1}^{*} + \eta\lambda(V_{com}^{*})_{j}(B_{*})_{m+1,1}e_{1}\right)$ 1462 1463  $h_{T}^{*}V_{FT} - B_{*}^{T}V_{i}^{*} = \sum_{j=1}^{m} (V_{com}^{*})_{j} \left(b_{j}^{*} + \eta\lambda(V_{com}^{*})_{j}\epsilon^{2}\sigma^{2}b_{m+1}^{*} + \eta\lambda(V_{com}^{*})_{j}(B_{*})_{m+1,1}e_{1}\right)$ 

$$+ \eta \lambda \epsilon^2 \sigma^2 b_{m+1}^* - \lambda b_{m+i}^* + \sum_{j=m+1}^{n} \Big( \eta \lambda(B_*)_{m+1,1}(B_*)_{j,1} \Big) b_j^*$$
(6)

$$+\sum_{j=1}^{m} \left( \eta \lambda(B_{*})_{m+1,1}(B_{*})_{j,1} \right) \left( b_{j}^{*} + \eta \lambda(V_{com}^{*})_{j} \epsilon^{2} \sigma^{2} b_{m+1}^{*} + \eta \lambda(V_{com}^{*})_{j}(B_{*})_{m+1,1} e_{1} \right).$$

$$(7)$$

From equation (4), we observe that  $\mathcal{L}_G(V_{LPFT}, B_{LPFT})$  is a monotonically decreasing function of and as  $\lambda$  approaches zero,  $\mathcal{L}_G(V_{LPFT}, B_{LPFT})$  also converges to zero. In contrast, combining equations (5) and (6), we find that  $\mathcal{L}_G(V_{FT}, B_{FT})$  does not converge to zero as  $\lambda$  approaches zero due to the presence of constant terms independent of  $\lambda$ . Therefore, it follows that there always exists a threshold  $\lambda^*$  such that for all  $\lambda \leq \lambda^*$ :

$$\mathcal{L}_G(V_{LPFT}, B_{LPFT}) \le \mathcal{L}_G(V_{FT}, B_{FT}).$$

G EMPIRICAL PERFORMANCE OF LP-FT AND FT UNDER THEOREM 4.5 CONDITIONS

To give a better understanding of Theorem 4.5, we give a simple visualization of two randomly generated data-generating functions for different clients and compute the global loss of LP-FT and FT based on equations (4) and (5).



Figure 6: (a) Global loss of LP-FT and FT as a function of the heterogeneity parameter  $\lambda$ , with  $\eta = 0.1$ ,  $\epsilon = 0.1$ , matrix  $B_*$  as a  $10 \times 20$  random matrix, and number of clients C = 5. (b) Global loss of LP-FT and FT as a function of the heterogeneity parameter  $\lambda$ , with  $\eta = 0.1$ ,  $\epsilon = 1$ , matrix  $B_*$  as a  $10 \times 20$  random matrix, and number of clients C = 5.

1503 1504

1475

1476

1477 1478

1479

1480 1481

These examples illustrate, within the theoretical setting of Sec. 4.2, the behavior of the loss functions for LP-FT and FT with a randomly generated labeling function  $y = V_i^{*T} B_* x$ , a fixed learning rate  $\eta$ , noise parameter  $\epsilon$ , and a fixed number of clients C. To compute this, we generated 1000 random matrices  $B_*$  and 1000 randomly chosen linear heads  $V_i^*$  as ground-truth labeling functions, ensuring they adhere to the theoretical assumptions. Using equations (4) and (5), we calculated the average loss of LP-FT and FT across these random trials.

As shown in Fig. 6, there exists a threshold  $\lambda^*$  such that when  $\lambda \leq \lambda^*$ , LP-FT consistently outperforms FT. While this is a simplified example with a fixed number of clients, learning rate, noise parameter,

and dimensionality of the ground-truth parameters  $B_*$  and  $V_i^*$ , the observed trend remains similar across different parameter settings. The purpose of this figure is to provide an intuitive understanding of Theorem 4.5 in a controlled, simplified context. More comprehensive experiments in Sec. 5 demonstrate that LP-FT globally outperforms FT across a broader range of heterogeneity levels in real-world settings.

- 1517 1518
- 1519 1520

1521

### H COMPUTATIONAL COST OF LP-FT COMPARED TO FT

1522 In this section, we want to study the computation cost of adding one step of linear probing (LP) to 1523 the full-fine tuning (FT) to see how this additional LP step affects computational cost. Suppose the 1524 dimension of the output of the feature extractor layer (input of the linear head) is d and the dimension 1525 of the output of the linear head is m. In fact, the linear head will be a  $d \times m$  linear layer. We assume 1526 having n samples. We want to see what is the computational cost of fine-tuning this linear head.

To estimate the computational cost of training a linear neural network layer with d inputs, m outputs, and n samples, we analyze the steps involved:

| 1529                         | 1. Forward Pass: A linear neural network computes outputs as:   |
|------------------------------|---|
| 1530                         | V VII/  |
| 1531                         | Y = XW,   |
| 1532                         | where:  |
| 1533                         | • $X \in \mathbb{R}^{n \times d}$ is the input matrix (with <i>n</i> samples each of dimension <i>d</i> )   |
| 1534                         | • $W \in \mathbb{D}^{d \times m}$ is the weight matrix  |
| 1535                         | • $V \in \mathbb{R}^{n \times m}$ is the output matrix.   |
| 1536                         | • $I \in \mathbb{R}$ is the output matrix.  |
| 1537                         | The cost of this matrix multiplication is $O(ndm)$ .  |
| 1538<br>1539                 | 2. <b>Backward Pass (Gradient Computation):</b> To update $W$ , the gradient of the loss $\mathcal{L}$ with respect to $W$ is computed. We know that:   |
| 1540                         |   |
| 1541                         | $ abla_W \mathcal{L} = X^T \left(  abla_Y \mathcal{L}  ight)$   |
| 1542                         | Therefore computing $(\nabla_{uv} f)$ involves:   |
| 1543                         | $(V_W \omega) \text{ involves.}$  |
| 1544                         | • Computing the gradient of the loss with respect to the outputs Y, which has a cost of $O(nk)$   |
| 1545                         | $O(n\kappa),$   |
| 1546<br>1547                 | • matrix multiplication $\nabla_W \mathcal{L} = X^+ (\nabla_Y \mathcal{L})$ which involves a matrix multiplication with a cost of $O(ndm)$ .  |
| 1548                         | 3. Weight Update: If using gradient descent, the cost of updating the weights is $O(dm)$ .  |
| 1549<br>1550<br>1551<br>1552 | <b>Total Computational Cost:</b> The total cost for one forward and backward pass through the data is dominated by $O(ndm)$ , which accounts for both forward propagation and gradient computation. If the training involves multiple epochs, the total cost scales as: |
| 1553                         | O(e, ndm)   |
| 1554                         | $O(c \cdot nam),$   |
| 1555                         | where $e$ is the number of epochs.  |
| 1556                         |   |
| 1557                         |   |
| 1558                         |   |
| 1559                         |   |
| 1560                         |   |
| 1561                         |   |
| 1562                         |   |
| 1563                         |   |
| 1564                         |   |
| 1565                         |   |
|                              |   |