

# FedReFT: Federated Representation Fine-Tuning with All-But-Me Aggregation

Fatema Siddika<sup>1\*</sup>, Md Anwar Hossen<sup>1\*</sup>, Juan Pablo Munoz<sup>2</sup>,  
Tanya Roosta<sup>3,4</sup>, Anuj Sharma<sup>1</sup>, Ali Jannesari<sup>1</sup>

<sup>1</sup>Iowa State University, Ames, USA

<sup>2</sup>Maro Systems, USA

<sup>3</sup>University of California, Berkeley, <sup>4</sup>Amazon

{fatemask, manwar, anuj, jannesari}@iastate.edu

pablo.munoz@maro-systems.com, troosta@ischool.berkeley.edu

## Abstract

Parameter-efficient fine-tuning (PEFT) adapts large pre-trained models by updating only a small subset of parameters. Recently, Representation Fine-Tuning (ReFT) has emerged as an effective alternative. ReFT shifts the fine-tuning paradigm from updating model weights to directly manipulating hidden representations that capture rich semantic information, and outperform state-of-the-art PEFTs in standalone settings. However, its application in Federated Learning (FL) remains challenging due to heterogeneity in clients' data distributions, model capacities, and computational resources. To address these challenges, we introduce **Federated Representation Fine-Tuning (FedReFT)**, a novel approach to fine-tune clients' hidden representations. FedReFT applies sparse intervention layers to steer hidden representations directly, offering a lightweight and semantically rich fine-tuning alternative ideal for edge devices. However, representation-level updates are especially vulnerable to aggregation mismatch under different task heterogeneity, where naive averaging can corrupt semantic alignment. To mitigate this issue, we propose All-But-Me (ABM) aggregation, where each client receives the aggregated updates of others and partially incorporates them, enabling stable and personalized learning by balancing local focus with global knowledge. We further design an adaptive update strategy inspired by Test-Time Computing (TTC) to balance local and global contributions under heterogeneous conditions. FedReFT achieves state-of-the-art performance on commonsense reasoning, arithmetic reasoning, and GLUE benchmarks, while delivering  $1\times\text{--}49\times$  higher parameter efficiency compared to leading LoRA-based methods. The paper code is available at [Anonymous Repository](#)

## 1 Introduction

Fine-tuning has emerged as a core strategy for adapting large language models (LLMs) to various downstream tasks, allowing for a broad generalization from minimal task-specific data (Ding et al., 2023; Ziegler et al., 2019). However, traditional full fine-tuning is computationally expensive and memory-intensive, which poses scalability challenges. This is further amplified in resource-constrained environments, such as smartphones, where full model updates are often infeasible due to limited resources. To address these challenges, parameter-efficient fine-tuning (PEFT) methods, such as Adapter Tuning (Houlsby et al., 2019), BitFit (Zaken et al., 2022), Prefix Tuning (Li and Liang, 2021), Prompt Tuning (Lester et al., 2021), and Low-Rank Adaptation (LoRA) (Hu et al., 2021a), have been proposed. These methods significantly reduce the cost of adaptation by updating only a small subset of model weights.

PEFT has emerged as the preferred method for efficiently adapting large language models (LLMs) without sacrificing performance. However, most PEFT approaches assume centralized data access, which is unrealistic in many real-world scenarios where data is distributed across users or devices with varying tasks and privacy concerns. Federated Learning (FL) offers a solution by enabling collaborative model training without centralizing data, but prior FL work often emphasizes task-specific tuning rather than learning generalizable representations. In practice, clients frequently work on diverse or specialized tasks, making global representation learning both more difficult and more essential.

While PEFT typically modifies model weights, recent interpretability research highlights the potential of hidden representations, which encode rich semantic information. ReFT (Wu et al., 2024c) leverages this by training small interventions that act as

---

\*Equal contribution

lightweight transformations on the model’s internal representations, steering model behavior for downstream tasks without altering the original weights. This representation-level adaptation enables ReFT to achieve stronger performance than weight-based methods such as LoRA. Despite ReFT’s success in centralized settings, it has yet to be adapted to FL setting, where challenges such as data heterogeneity, varying model capacities, and limited computational resources complicate aggregation and reduce effectiveness. To investigate the challenges of representation-level fine-tuning in heterogeneous federated settings and assess the effectiveness of our proposed aggregation strategy, we formulate the following research questions:

**RQ1:** How can representation-level updates be aggregated in FL while preserving semantic alignment across task-heterogeneous clients?

**RQ2:** Is simple weighted averaging sufficient for aligning semantically rich, hidden representations, or is a more robust and personalized strategy required to maintain local semantics while leveraging global knowledge?

To address these challenges, we propose **Federated Representation Fine-Tuning (FedReFT)**, a framework for personalized and parameter-efficient federated fine-tuning. FedReFT extends ReFT by injecting lightweight intervention components (sparse low-rank matrices  $W$ ,  $R$ ,  $b$ ) directly into hidden representations, making it suitable for resource-limited edge devices. To avoid semantic misalignment from naive aggregation such as FedAvg (McMahan et al., 2017), we introduce *All-But-Me (ABM)* aggregation, which builds a global aggregated intervention by computing the geometric median over updates from all other clients. The key contributions of our work are as follows:

**Contribution 1:** We address a critical gap in adapting ReFT to federated learning by introducing FedReFT, the first framework that enables personalized and parameter-efficient fine-tuning through sparse representation-level interventions, making it effective for resource-constrained clients.

**Contribution 2:** To support this framework, we propose the *All-But-Me (ABM)* aggregation strategy, specifically designed for representation-level interventions. ABM mitigates semantic misalignment caused by naive averaging and enables stable, personalized collaboration under heterogeneous conditions. Furthermore, FedReFT incorporates an adaptive mixing mechanism inspired by *Test-Time Computing (TTC)* to dynamically learn how

to combine local and ABM intervention parameters for each client.

**Contribution 3:** We evaluate the framework by simulating task heterogeneity, i.e., assigning different tasks to clients, all derived from a common dataset. This setup mimics real-world scenarios where clients pursue distinct objectives over structurally similar data, allowing us to evaluate the effectiveness of FedReFT and ABM under realistic conditions. The rest of the paper is organized as follows. Section 2 defines the problem and challenges in heterogeneous FL. Section 3 presents our FedReFT framework and ABM aggregation. Section 4 reports experimental results, and Section 5 concludes with insights and future directions. Additional details are provided in the appendix.

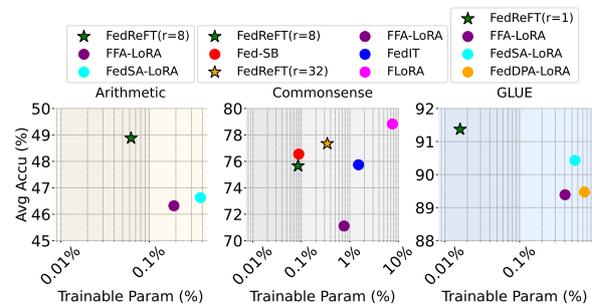


Figure 1: Average accuracy vs. trainable parameters (%) for federated PEFT methods on Arithmetic, Commonsense, and GLUE benchmarks using LLaMA-3 8B, LLaMA-3.2B, and RoBERTa-large models, respectively. FedReFT attains state-of-the-art accuracy while training far fewer parameters, improving communication efficiency and reducing transmission cost in FL.

## 2 Problem Formulation and Motivation

In this section, we motivate applying ReFT to FL and formalize the problem of personalized representation adaptation. While ReFT enables parameter-efficient updates in the representation space, its use in FL is hindered by task heterogeneity, semantic misalignment, and unstable aggregation. Therefore, we define the objective as achieving parameter-efficient and semantically aligned adaptation in FL with ReFT.

**Challenge 1: LoReFT in FL Settings (RQ1):** ReFT (Wu et al., 2024c) offers an attractive alternative by modifying hidden activations instead of model weights. By intervening directly in structured semantic subspaces, ReFT supports interpretable, modular, and task-aligned adaptation, particularly advantageous in task-heterogeneous FL

settings. However, deploying full ReFT in federated settings can be challenging, as transmitting high-dimensional representation updates increases communication overhead. Moreover, when clients employ models with varying capacities or architectures, aligning representation spaces becomes nontrivial, making integration across clients difficult.

To bridge this gap, we adopt Low-Rank Linear Subspace ReFT (LoReFT)(Wu et al., 2024c), which is a lightweight ReFT variant that constrains interventions to a learnable low-rank subspace. This design significantly reduces overhead while maintaining semantic control, making it a promising candidate for FL. We follow the LoReFT intervention formulation from (Wu et al., 2024c) on hidden representations  $h \in \mathbb{R}^d$  which is defined as:

$$\Phi_{\text{LoReFT}}(h) = h + R^\top (Wh + b - Rh), \quad (1)$$

where,  $W \in \mathbb{R}^{r \times d}$  is a low-rank projection matrix with  $d$  as the representation dimension and  $r$  as the subspace intervention dimension,  $R \in \mathbb{R}^{r \times d}$  is a low-rank projection matrix with orthonormal rows, and  $b \in \mathbb{R}^r$ , with  $r \ll d$ . This structure, inspired by Distributed Interchange Intervention (DII) (Geiger et al., 2024), enables semantically grounded, low-rank adaptation suitable for scalable and privacy-preserving FL. Despite its efficiency, applying LoReFT in FL raises several non-trivial challenges: (i) LoReFT modifies internal representations that are sensitive to client-specific data distributions. Aggregating these interventions naïvely using FedAvg can cause semantic interference or collapse. (ii) Without global synchronization, low-rank updates may evolve in divergent directions, especially when tasks are dissimilar. (iii) Applying shared LoReFT interventions across clients risks overfitting to shared patterns while ignoring local semantics. Considering all these challenges, the major research question is:

*Can representation-level adaptation via LoReFT achieve personalization and stability in federated environments without collapsing under task and data heterogeneity?*

FedReFT uses All-But-Me(ABM) aggregation to robustly combine intervention parameters while preserving personalization in heterogeneous FL. **Challenge 2: Federated Fine-Tuning under Task Heterogeneity (RQ1 & RQ2):** A central motivation of our work is to address task heterogeneity in real-world FL, where clients perform fundamentally different tasks rather than optimizing a

shared objective. For example, clients may work on distinct reasoning tasks within natural language question answering that demand different semantic skills. While centralized fine-tuning has proven effective for such tasks, it assumes access to all data, which is unrealistic in decentralized settings. In FL, each client sees only a local, task-specific subset of the broader reasoning space, leading to highly heterogeneous training distributions, a common challenge in multi-department or cross-domain deployments. This raises the question: *How can we learn a global representation that generalizes across tasks when each client trains only on a fragment of the broader task distribution?* Standard methods like FedAvg (McMahan et al., 2017) struggle in this regime, as they average semantically misaligned updates, often resulting in degraded performance or collapsed representations. Formally, let each client  $i$  have a dataset  $\mathcal{T}_i = \{X_i, Y_i\}$  and optimize a personalized model  $\theta_i$  by solving:

$$\min_{\Theta} \frac{1}{N} \sum_{i=1}^N \mathcal{L}(X_i, Y_i, \theta_i), \quad (2)$$

where  $\mathcal{L}$  is the task-specific loss and  $\Theta = \{\theta_i\}_{i=1}^N$  is the set of client-specific models. Our proposed method, FedReFT, can successfully address this research challenge. FedReFT enables scalable, personalized representation learning across heterogeneous tasks, allowing global reasoning capabilities to emerge from decentralized, task-specific updates. **Challenge 3: Learnable Parameter Sharing with the Server (RQ2):** When applying ReFT (Wu et al., 2024c) in a FL setting from the perspective of learnable parameter sharing, a fundamental question is:

*Which of these parameters should be communicated to the server for collaborative aggregation?* In FedReFT, each client fine-tunes hidden representations by introducing learnable low-rank intervention parameters  $W$ ,  $R$ , and a bias  $b$  into a frozen backbone model. Sharing only some part of the intervention parameters leads to incomplete information transfer and breaks the low-rank structure critical for generalization.  $W$  projects representations into a low-dimensional space, and  $R$  reconstructs them; omitting either disrupts compositionality and limits alignment across heterogeneous clients. In Table 1, empirical results show that partial sharing significantly degrades performance and representation alignment under task heterogeneity. Therefore, FedReFT shares the full set of learnable interven-

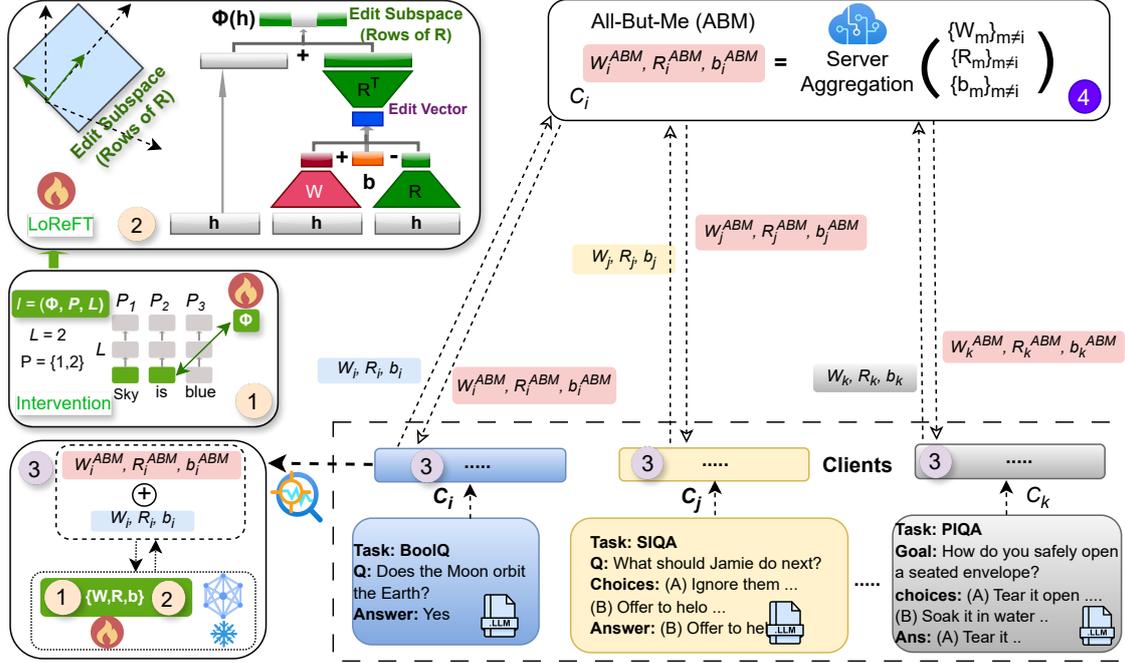


Figure 2: **FedReFT with ABM Aggregation.** Clients cross-task demonstrate personalization while maintaining alignment with the global representation. (1)-(2): Each client applies LoReFT (Wu et al., 2024c) interventions to train learnable parameter  $\{W, R, b\}$  to modify hidden representations  $h$  in a low-rank edit subspace. (3): Clients fine-tune  $\{W, R, b\}$  locally and partially fuse received *All-But-Me* aggregated updates with their own. (4): The server performs ABM aggregation using the geometric median over other clients’ intervention parameters to generate  $W_k^{ABM}, R_k^{ABM}, b_k^{ABM}$ .

tion parameters  $(W, R, B)$  from each client with the server.

### 3 Methodology

In this section, we introduce FedReFT, designed to address the challenges we discussed in the previous section. An illustrative overview of FedReFT is shown in Figure 2.

#### 3.1 Intervention Parameter Sharing Strategies

An intervention is a small, learnable operation inserted into a model’s hidden layers that slightly modifies its internal representations, known as activations. Instead of updating the model’s weights, it transforms these hidden features using a few parameters in Equation 1; a low-rank projection matrix  $R \in \mathbb{R}^{r \times d}$  that captures compact feature subspaces, a transformation matrix  $W \in \mathbb{R}^{r \times d}$  that maps these subspaces back into the model’s space, and a bias vector  $b \in \mathbb{R}^r$  that shifts the representation to adjust how information flows through the network. To balance communication overhead with personalization, we evaluate three sharing strate-

gies. The **Full Intervention** approach transmits the complete parameter set  $\{W \in \mathbb{R}^{r \times d}, R \in \mathbb{R}^{r \times d}, b \in \mathbb{R}^r\}$ , capturing all transformation aspects for optimal global performance. The **No Bias** variant  $\{W \in \mathbb{R}^{r \times d}, R \in \mathbb{R}^{r \times d}\}$  retains subspace alignment but misses fine-grained shifts. Finally, the **No  $W$**  configuration  $\{R \in \mathbb{R}^{r \times d}, b \in \mathbb{R}^r\}$  maximizes efficiency but severely degrades performance by omitting the crucial encoding matrix  $W$ . After observing results in Table 1, FedReFT shared the full Intervention parameters.

#### 3.2 Intervention Design for Federated Classification Tasks

Following the formulation in ReFT (Wu et al., 2024c), for a given client, we define the classification head  $H_\psi$  with parameters  $\psi = \{W_o, b_o, W_d, b_d\}$  which operates on the CLS token representation  $z \in \mathbb{R}^d$  from the final layer:

$$H_\psi(z) = \text{softmax}(W_o \cdot \tanh(W_d z + b_d) + b_o). \quad (3)$$

We jointly optimize the intervention parameters  $\phi$  and the classifier  $\psi$  using cross-entropy loss over

input  $x$  and label  $y$ :

$$\min_{\phi, \psi} \{-\log H_{\psi}(y | z_{\phi}(x))\}. \quad (4)$$

Table 1: Performance vs. parameter efficiency for different LoReFT sharing strategies (Uplink) for  $C$  clients on commonsense reasoning task following the second experiment design.

Task	Strategy	TP(% ↓)	Accu↑
GLUE	W, R, b	0.01384	94.91
ROBERTa	W, R	0.01383	64.03
	R, b	0.00693	74.12
Arithmetic	W, R, b	0.03114	33.31
LLaMA-2 7B	W, R	0.03114	26.13
	R, b	0.01557	25.77
Commonsense	W, R, b	0.03114	76.55
LLaMA-2 7B	W, R	0.03114	70.63
	R, b	0.01557	68.02

### 3.3 All-But-Me (ABM) Aggregation on Server

In heterogeneous FL, the integration of shared knowledge without compromising local task-specific adaptation remains a core challenge. Standard aggregation methods, such as FedAvg (McMahan et al., 2017), which averages client models into a single global model, are often suboptimal in non i.i.d. scenarios. They risk overwriting valuable client-specific representations and rely on fixed mixing weights that may further reduce personalization. To overcome these limitations, we propose the *All-But-Me* (ABM) aggregation strategy. Instead of initializing clients with a global model, each client continues to update its local parameters while partially incorporating knowledge aggregated from other clients. Specifically, each client  $k$  receives a robustly aggregated set of intervention parameters  $\{W_k^{\text{ABM}}, R_k^{\text{ABM}}, B_k^{\text{ABM}}\}$ , calculated from the updates of all other clients using a geometric median:

$$X = \{W, R, B\}, \quad X_k^{\text{ABM}} = \text{ABM}(\{X_m^c\}_{m \neq k}). \quad (5)$$

**ABM via Geometric Median.** The geometric median (also known as the spatial or  $L_1$  median) offers a robust alternative to the arithmetic mean, particularly under client heterogeneity and adversarial conditions (Maronna and Martin, 2006; Weiszfeld,

1937). Given a set of vectors  $\mathcal{S} = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$ , it is defined as:

$$x^* = \arg \min_{x \in \mathbb{R}^d} \sum_{i=1}^n \|x - x_i\|_2, \quad (6)$$

which minimizes the sum of Euclidean distances to all elements in the set. This estimator is robust to outliers and misaligned updates, making it well-suited for federated settings. We instantiate the ABM function using the geometric median, where each client  $k$  receives an aggregated intervention vector computed from  $\mathcal{S}_k = \{x_m\}_{m \neq k}$ :

$$\text{ABM}(\mathcal{S}_k) = \arg \min_{x \in \mathbb{R}^d} \sum_{x_m \in \mathcal{S}_k} \|x - x_m\|_2. \quad (7)$$

To solve this optimization efficiently, we employ

Table 2: Comparison of FedAvg, Mean\_ABM, and GeoMedian\_ABM across commonsense, arithmetic, and GLUE tasks. Geometric median-based ABM aggregation achieves better accuracy over all other aggregation strategies.

Task, Model	Agg. Method	Acc	$\Delta \text{Acc} \uparrow$
Commonsense	FedAvg	74.89	
LLaMA-2 7B	Mean_ABM	75.15	+0.25
	<b>GeoMed_ABM</b>	<b>76.55</b>	<b>+1.66</b>
Arithmetic	FedAvg	24.87	
LLaMA-2 7B	Mean_ABM	25.12	+0.25
	<b>GeoMed_ABM</b>	<b>26.09</b>	<b>+1.22</b>
GLUE	FedAvg	93.57	
RoBERTa	Mean_ABM	94.10	+0.53
	<b>GeoMed_ABM</b>	<b>94.91</b>	<b>+1.34</b>

Weiszfeld’s algorithm, an iterative method known to converge under mild conditions. Details of the algorithm are provided in Table 2, Figure 5, and Appendix D.3. By avoiding direct averaging and incorporating semantically meaningful low-rank intervention updates through robust aggregation, ABM enables each client to benefit from the knowledge of others without sacrificing local personalization. This approach enhances stability and generalization across non-i.i.d. and task-heterogeneous FL environments. Geometric median has complexity of  $O(T \cdot d)$ , with  $T$  being the number of iterations and  $d$  the number of parameters.

### 3.4 Local Model Update with ABM

#### Aggregation using Test-Time Computing

In FedReFT, local interventions capture client-specific semantics, while ABM aggregation con-

veys shared global knowledge. However, naively averaging representation-level updates can risk disrupting semantic consistency across clients, potentially leading to misaligned feature spaces. Moreover, uniform aggregation overlooks client heterogeneity, limiting adaptability in diverse edge environments.

These challenges highlight the need for a mechanism that can dynamically balance personalization and global alignment. Without such adaptation, clients may either overfit to their own data or lose semantic fidelity when forced into uniform global updates. FedReFT introduces an adaptive mixing strategy inspired by Test-Time Computing (TTC) (Wang et al., 2021), which learns client-specific coefficients to combine local and ABM aggregated intervention parameters. Unlike conventional fine-tuning that updates model weights, TTC performs lightweight optimization during inference, using a small, controlled compute budget to iteratively refine predictions and learn adaptive mixing parameters. Here, unconstrained logits  $\beta \in \mathbb{R}$  define mixing weights  $\alpha = \sigma(\beta) \in (0, 1)$  via the sigmoid  $\sigma(\cdot)$ . For any intervention tensor  $X \in W, R, B$ , the mixed parameter is:

$$X_{\text{mixed}}(\alpha) = \alpha X^{\text{local}} + (1 - \alpha) X^{\text{ABM}}. \quad (8)$$

TTC runs for  $S$  steps over  $B$  calibration batches of clients’ test data. At each step, it injects  $X_{\text{mixed}}(\alpha)$  into the model interventions, computes a forward pass to evaluate the test loss  $\mathcal{L}_{\text{test}}$ , backpropagates gradients  $\nabla_{\beta} \mathcal{L}$  while keeping  $X^{\text{local}}$  and  $X^{\text{ABM}}$  frozen, and then updates  $\beta \leftarrow \beta - \eta \nabla_{\beta} \mathcal{L}$ . The gradient  $\nabla_{\beta} \mathcal{L} = \nabla_{\alpha} \mathcal{L} \sigma(\beta) (1 - \sigma(\beta))$  directly indicates whether increasing  $\alpha$  reduces the test loss: negative gradients push  $\beta$  higher (increasing  $\alpha$ ), while positive gradients push it lower. Over iterations,  $\beta$  converges to:

$$\beta^* = \arg \min_{\beta} \mathcal{L}_{\text{test}}(X_{\text{mixed}}(\sigma(\beta))) + \lambda_1 \mathcal{H}(\alpha) + \lambda_2 \mathcal{C}(\alpha) + \lambda_3 \mathcal{D}(\alpha), \quad (9)$$

where  $\mathcal{H}$  (entropy) discourages collapse,  $\mathcal{C}$  (consistency) promotes stable mixing across keys, and  $\mathcal{D}$  (diversity) avoids extreme  $\alpha$  saturation and the details of hyper parameter  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  described in Appendix D.2.

**Final Mixed Interventions.** The optimal  $\alpha^* = \sigma(\beta^*)$  is used to mix the local and ABM intervention parameters:

$$X_k^{\text{new}} = \alpha^* X_k^{\text{local}} + (1 - \alpha^*) X_k^{\text{ABM}}, \quad (10)$$

$$X \in \{W, R, B\}.$$

Before the next local training with  $R_k^{\text{new}}$ , we apply an orthogonal transformation to  $R_k^{\text{new}}$  to preserve the original property of  $R$ . We continue the discussion in Appendix D.2 and Table 12, 13, 14, and 15 show that the performance of TTC-based adaptive mixing is over the balanced mixing.

**Computational Overhead.** The TTC process costs  $S \times B$  forward/backward passes per client. With  $N$  clients,  $K$  intervention parameters, batch size  $B$ ,  $L$  layers, and hidden width  $d_h$ , an approximate FLOPs is  $\mathcal{O}(S \cdot M \cdot N \cdot K \cdot B \cdot L \cdot d_h^2)$ . Amortized over  $T$  communication rounds and  $|\mathcal{D}_{\text{test}}|$  inference batches, the effective per-inference overhead scales as  $\frac{S \cdot B}{T \cdot |\mathcal{D}_{\text{test}}|}$ . Here, small  $S=50$  and  $B=8$  keep the overhead modest while still improving robustness under shift.

## 4 Experimental Validation

To evaluate FedReFT, we conduct extensive experiments on three different NLP benchmarks covering over 12 datasets. Our objective is to present a comprehensive assessment of how this approach performs in various NLP tasks. We experiment with both masked and autoregressive language models, including RoBERTa-large, TinyLLaMA-1B, LLaMA 7B, LLaMA-2 7B and 13B, LLaMA-3.2B and LLaMA-3 8B, across multiple settings and scales. Our comparisons include state-of-the-art baselines, such as LoRA (Hu et al., 2021b), FedIT (Zhang et al., 2024), FFA-LoRA (Sun et al., 2024), FedDPA-LoRA (Long et al., 2024), FedSA-LoRA (Guo et al., 2024), Fed-SB (Singhal et al., 2025) and FLoRA (Wang et al., 2024) focusing on both parameter efficiency and performance trade-offs. We align the experimental setup configurations with the baseline papers to ensure fair comparisons. To optimize memory usage, we load all base language models with torch.bfloat16 precision. The results are averaged over two runs to report the mean performance.

**Hyperparameter Configuration.** In the experiments, we determine the number of interventions to learn and the specific layers and input positions where they are applied. Interventions are inserted at a fixed number of layers ( $L$ ) and at the prefix ( $p$ ) and suffix ( $s$ ) positions of the input prompt. We narrow the hyperparameter search space for FL in Appendix B by adopting the configuration used in the centralized ReFT (Wu et al., 2024c) paper.

**Task Distribution Rationale.** We design two experimental setups for commonsense and arithmetic

Table 3: Performance of LLaMA-3.2 3B across five commonsense reasoning tasks with Mixed-Task (MT) setup, where clients train on heterogeneous task mixtures to promote generalizable representations. **Trainable Parameter TP Efficiency Rank 8** quantifies the relative parameter efficiency of FedReFT(R8), indicating how many times fewer trainable parameters it requires compared to baseline methods.

Method	R	TP(M) ↓	TP Effi. (R8) ↓	BoolQ	PIQA	SIQA	HellaS.	WinoG	Avg Acc↑
FedIT*	32	48.63	17.68×	62.99	81.50	73.13	76.83	71.51	75.74
FFA-LoRA*	32	24.31	8.84×	62.87	80.03	68.53	70.02	65.56	71.11
Fed-SB*	120	2.83	1.03×	64.86	81.66	74.87	81.67	75.22	75.66
Fed-SB*	200	7.8	2.84×	66.66	83.79	77.22	85.42	79.56	78.53
	4	1.38	0.5×	63.35	82.72	72.96	91.37	69.70	76.02
<b>FedReFT (ours)</b>	<b>8</b>	<b>2.75</b>	<b>1.0×</b>	<b>65.50</b>	<b>82.32</b>	<b>73.28</b>	<b>91.43</b>	<b>70.24</b>	<b>76.55</b>
	<b>16</b>	<b>5.5</b>	<b>1.0×</b>	<b>64.56</b>	<b>82.20</b>	<b>75.95</b>	<b>89.80</b>	<b>80.59</b>	<b>78.62</b>

Table 4: Performance of Federated fine-tuning of Llama-3.2 3B across eight commonsense reasoning datasets in a highly data-heterogeneous setting, which is denoted as Distinct Task (DT). Trainable Parameter (TP) Efficiency Rank quantifies the relative parameter efficiency of FedReFT(R8). Best results are in **bold**.

Method	R	TP(M) ↓	BoolQ	PIQA	SIQA	HellaS	WinoG	ARC-e	ARC-c	OBQA	Avg↑
FedIT*	32	48.63	60.89	78.22	69.92	73.18	67.78	81.21	67.04	66.91	70.80
FFA-LoRA*	32	24.31	60.73	76.91	65.37	68.61	61.89	79.41	62.92	63.12	67.17
FedEx-LoRA*	32	243.15	62.55	79.36	71.41	71.78	72.45	82.69	67.80	70.25	73.13
FLoRA*	32	243.15	62.55	79.36	71.41	71.78	72.45	82.69	67.80	70.25	73.13
Fed-SB*	200	7.85	63.28	80.34	73.56	82.07	76.01	84.01	69.02	72.46	75.21
<b>FedReFT(ours)</b>	<b>8</b>	<b>2.75</b>	<b>63.84</b>	<b>78.01</b>	<b>77.53</b>	<b>89.74</b>	<b>72.30</b>	<b>84.57</b>	<b>69.31</b>	<b>77.00</b>	<b>76.41</b>
	<b>16</b>	<b>5.50</b>	<b>64.67</b>	<b>81.88</b>	<b>78.06</b>	<b>89.88</b>	<b>77.82</b>	<b>86.11</b>	<b>69.37</b>	<b>76.20</b>	<b>78.00</b>

reasoning tasks to study how global representations converge under diverse task distributions. In the Mixed-Task (MT) setup, each client trains on a subset of a combined reasoning dataset but is evaluated on a single task, encouraging generalized, transferable representations through ABM aggregation. This reflects collaborative learning across varied yet related tasks. In the Distinct Task (DT) setup, each client trains on a unique reasoning task, enabling personalized fine-tuning while still leveraging global updates. Despite higher task heterogeneity, this setup maintains stable performance as model capacity increases. Both setups show that FedReFT supports effective generalization in MT and robustness in DT.

#### 4.1 Commonsense Reasoning

We evaluate global representation generation on eight commonsense reasoning tasks using the Commonsense170K dataset inspired by (Singhal et al., 2025; Wu et al., 2024c). We use the same hyperparameter of (Singhal et al., 2025) and tune the

intervention parameter in the Appendix B.1. This helps us tune important hyperparameters efficiently and also test their robustness across multiple commonsense reasoning tasks.

**Datasets.** For the first setup, as a Mixed-Task (MT) design, we split the commonsense reasoning tasks dataset Commonsense170K (Hu et al., 2023) among clients and use them for fine-tuning. Each client evaluates one of the commonsense reasoning tasks. BoolQ, PIQA, SIQA, HellaSwag, and WinoGrande. For the second setup as Distinct Task (DT) design, each client fine-tunes on only one of these five commonsense reasoning tasks and is evaluated using the same task. We adopt the prompt template from Hu et al. (Hu et al., 2023). The Distinct Task (DT) experiment is described in the Appendix 16, as no baseline was found in this setting.

**Results.** In Table 3, our proposed FedReFT method demonstrates strong parameter efficiency while maintaining competitive accuracy across five commonsense reasoning tasks. Notably, FedReFT with rank 8 uses only 2.75M (0.0857%) trainable pa-

Table 5: Performance comparison across arithmetic reasoning tasks with the Distinct Task (DT) and Mixed Task (MT) setup using different model sizes. We report results under adaptive mixing with Test-Time Computing, which dynamically balances local and global interventions at inference. Avg represents the average of the clients’ accuracy.

FedReFT Models	Distinct Task (DT)				Mixed Task (MT)			
	AQuA	SVAMP	MAWPS	Avg Acc $\uparrow$	AQuA	SVAMP	MAWPS	Avg Acc $\uparrow$
LLaMA 7B	26.12	26.71	49.94	34.26	20.86	15.60	28.22	22.23
LLaMA-2 7B	30.96	33.31	59.51	41.93	22.05	23.50	32.74	26.09
LLaMA-3 8B	35.36	49.41	75.25	53.34	33.45	51.28	73.51	52.75

Table 6: Performance comparison across GLUE Tasks on RoBERTa model for  $C = 3$ , FedReFT use rank 1. All baseline methods use LoRA rank 8. FedReFT achieves higher accuracy over all the baselines in FL settings. **TP Efficiency** quantifies the relative trainable parameter efficiency of FedReFT, indicating how many times fewer parameters it requires compared to baselines.

Setup	Method	TP(M) $\downarrow$	TP Effi. $\downarrow$	MNLI-m	MNLI-mm	SST-2	QNLI	QQP	Avg $\uparrow$
Standalone	Full Tuning	355	6698.11 $\times$	88.8	88.56	96.0	93.8	91.5	91.73
	LoRA*	1.83	34.53 $\times$	88.71	88.21	95.16	91.16	85.33	89.71
	LoReFT	0.053	1.0 $\times$	89.2	89.26	96.20	94.10	88.5	91.45
FL	FFA-LoRA*	1.44	27.17 $\times$	88.83	88.27	94.95	91.52	86.71	89.39
	FedDPA-LoRA*	2.62	49.44 $\times$	88.99	88.43	95.50	90.74	85.73	89.47
	FedSA-LoRA*	1.83	10.40 $\times$	90.18	88.88	96.00	92.13	87.48	90.43
	<b>FedReFT</b>	<b>0.053</b>	<b>1.0<math>\times</math></b>	<b>89.75</b>	<b>89.31</b>	<b>95.75</b>	<b>94.91</b>	<b>87.15</b>	<b>91.37</b>

rameters, achieving accuracy close to or better than several baselines. Compared to existing methods our approach reduces the trainable parameter count by factors of 1.03 $\times$  to 13.68 $\times$ , with minimal to no compromise in performance, also shown in Figure 1. The experiments results on the Mixed-Task (MT) setup are shown in Appendix Table 16.

## 4.2 Arithmetic Reasoning

For the arithmetic reasoning tasks, we design three experimental settings to fine-tune models on various arithmetic reasoning tasks. We follow the same hyperparameter tuning strategy as used in Commonsense170K in Appendix B.1, which uses a development set to select the best-performing configuration. Evaluation is based solely on the final numeric or multiple-choice answer, disregarding intermediate reasoning steps.

**Datasets.** In the first setting following Mixed-Task (MT), we split the arithmetic reasoning dataset, MATH10K (Hu et al., 2023), which includes four tasks with chain-of-thought solutions generated by a language model. Each client reports performance using test set one of three tasks: AQuA, SVAMP, and MAWPS. In the second setting fol-

lowing Distinct-Task (DT), each client is assigned one arithmetic reasoning task for both fine-tuning and evaluation. In the third setting, following (Guo et al., 2024; Kuang et al., 2024), we split the dataset GSM8K into 3 clients under an IID distribution.

**Results.** In Table 7, FedReFT achieves the highest accuracy among all methods while requiring substantially fewer trainable parameters, underscoring its efficiency and superior performance. In Table 5, the DT setup represents task heterogeneity, enabling clients to learn personalized, task-specific representations while benefiting from global aggregation. Consequently, the DT setup achieves higher performance. In contrast, the Mixed-Task (MT) setup trains clients on heterogeneous task mixtures to promote globally generalizable representations, but this blending can reduce performance on specific evaluation tasks due to representation misalignment and conflicting objectives.

## 4.3 Natural Language Understanding

We evaluate the effectiveness of FedReFT in learning generalizable representations for Natural Language Understanding (NLU) using the GLUE benchmark (Wang et al., 2018). The objective is to

fine-tune NLU to learn global representations that capture task-level semantics. By aligning intermediate representations for downstream classification performance. This setup allows us to test whether lightweight intervention tuning can align representations across clients within a single NLU task.

**Results.** Table 6 depicts that our approach performs strongly across GLUE tasks while using very few trainable parameters. It performs competitively or surpasses other methods, demonstrating its ability to learn strong representations even under federated conditions. Despite utilizing  $10\times-49\times$  fewer trainable parameters than several baselines, it achieves comparable accuracy, underscoring its efficiency and effectiveness.

Table 7: Performance comparison on arithmetic reasoning task for GSM8K on LLaMA-3 8B with rank 8, where clients enable consistent evaluation of representation generalization. FedReFT was also evaluated with Mistral-7B and Gemma-2 9B models. **TP Effi** measures how many times fewer trainable parameters it uses compared to the baselines, while  $\Delta\text{Acc}$  indicates the accuracy gain achieved by FedReFT over baseline.

Method	TP(M) $\downarrow$	TP Effi. $\downarrow$	Acc	$\Delta\text{Acc}\uparrow$
<b>LLaMA-3 8B</b>				
LoReFT	4.19	1.0 $\times$	48.33	+0.55
LoRA*	30.40	7.26 $\times$	46.23	+2.65
FedSA-LoRA*	30.40	7.26 $\times$	46.63	+2.25
FFA-LoRA*	15.2	3.63 $\times$	46.32	+2.56
<b>FedReFT(ours)</b>	<b>4.19</b>	<b>1.0<math>\times</math></b>	<b>48.88</b>	
<b>Mistral-7B</b>				
<b>FedReFT(ours)</b>	<b>4.19</b>	<b>1.0<math>\times</math></b>	<b>47.63</b>	-1.25
<b>Gemma-2 9B</b>				
<b>FedReFT(ours)</b>	<b>4.81</b>	<b>1.14<math>\times</math></b>	<b>67.34</b>	+18.46

#### 4.4 Partial Client Participation (PCP)

The ability of a federated learning system to maintain performance when only a fraction of clients participate in each round, known as Partial Client Participation (PCP), is crucial for real-world deployment. This robustness to PCP demonstrates the system’s ability to handle practical constraints like communication bandwidth limitations, client availability, and device battery life. To evaluate partial client participation (PCP), the experiment’s results in Table 8 use a total of  $C = 5$  clients. When all five clients participated in every global round, FedReFT obtained an average accuracy of 91.45%.

Under PCP, where only  $C = 3$  clients were randomly selected per round, FedReFT achieved an average accuracy of 90.62%. The small difference between the two settings indicates that FedReFT maintains stable performance even with reduced client participation.

Table 8: Impact of Partial Client Participation (PCP) on FedReFT performance for the GLUE task. The table compares global accuracy when all clients participate ( $C = 5$ ) versus a random subset ( $C = 3$ ) per round, demonstrating the model’s robustness to limited client availability and communication constraints.

Method	MNLI-m	SST-2	QNLI	QQP	Avg $\uparrow$
PCP	89.17	92.83	93.81	86.65	90.62
All	89.61	95.25	94.02	86.93	91.45

#### 4.5 Ablation Studies

We conduct ablation studies in Appendix D to further evaluate the effectiveness of FedReFT, focusing on the impact of geometric-median-based All-But-Me aggregation, intervention parameter sharing strategies, and local model update approaches using balanced and TTC-based adaptive mixing.

#### 4.6 Computational Resources

All experiments are executed on a single NVIDIA A100-SXM4-80GB GPU, except for LLaMA-2 13B, which is run on a GPUH200x8 141GB system to accommodate the computational demands of large-scale federated fine-tuning.

### 5 Conclusion

In this work, we bridge the gap between Representation Fine-Tuning (ReFT) and Federated Learning by introducing FedReFT, a unified framework that enables personalized and parameter-efficient federated representation learning. FedReFT employs the *All-But-Me* aggregation strategy to mitigate semantic misalignment caused by naive averaging, enabling clients to adapt using a robust average of others’ interventions. Additionally, an adaptive mixing mechanism inspired by *Test-Time Computing* dynamically balances local and global representations, enhancing robustness under heterogeneous conditions. Extensive experiments demonstrate that FedReFT consistently improves convergence, generalization, and parameter efficiency across diverse FL settings.

## Limitations

Due to computational constraints, our current study focuses primarily on LoReFT-based interventions within language models under a fixed set of hyperparameters. In future work, we aim to automate the parameter search space using a multi-agent coordination framework to better explore optimal low-rank configurations for each client. Although our current set-up does not explicitly address privacy, we are actively investigating how to integrate differential privacy mechanisms, such as DP-SGD, into the FedReFT framework without sacrificing personalization. Initial experiments in this direction are ongoing. Additionally, we are exploring the theoretical properties of ABM aggregation under adversarial or noisy clients, and whether it can be extended to other modalities beyond language, such as vision-language models in federated systems.

## Data and Model Usage

We use publicly available models including LLaMA-1.1B, LLaMA-2 (7B, 13B), LLaMA-3 8B, LLaMA-3.2 3B, Gemma-2 9B, Mistral-7B and RoBERTa-large. LLaMA-2 and LLaMA-3 models are licensed under Meta’s community license permitting commercial use. RoBERTa-large is under the MIT License, and TinyLLaMA use Apache 2.0, while the original LLaMA-1 7B is for non-commercial research only. We will release code and configurations under an open-source license with usage documentation to support reproducibility and responsible use.

We employ publicly available datasets across commonsense and arithmetic reasoning tasks, each released under open-source licenses. For commonsense reasoning, BoolQ is under CC BY-SA 3.0, PIQA under Apache 2.0, SIQA and WinoGrande under CC BY 4.0, HellaSwag under MIT, ARC under CC BY-SA 4.0, and OBQA under CC BY 4.0. For arithmetic reasoning, AddSub, AQuA, MAWPS, and MultiArith are under Apache 2.0, GSM8K and SVAMP under MIT, and SingleEq under CC BY 4.0. For natural language understanding, GLUE consists of multiple datasets, each with its own license, allowing for research use and redistribution.

## Environmental Impact

Our approach FedReFT achieves  $1\times$ – $49\times$  higher parameter efficiency than existing PEFT methods, using fewer trainable parameters. This reduces en-

ergy consumption and training time, making our method more resource-efficient and environmentally friendly.

## Societal Impacts

Our method FedReFT adapts ReFT for Federated Learning, enabling efficient model personalization with minimal computational overhead. This promotes broader accessibility of large language models on edge devices, including in low-resource or privacy-sensitive environments. While improving inclusivity and deployment scalability, care must be taken to mitigate potential misuse or bias propagation across decentralized systems.

## Bias and Fairness

Our approach FedReFT considers the potential for bias introduced by non-IID client data in Federated Learning. While we do not explicitly optimize for fairness, we acknowledge that imbalanced participation or data diversity may lead to uneven model performance. Future work should explore fairness-aware objectives to mitigate such disparities across clients and demographic groups.

## Responsible Deployment

To support responsible use, we include clear documentation outlining the intended use cases of our framework and advise against applying it in safety-critical settings without thorough validation. We encourage users to follow ethical standards, such as the ACL Code of Ethics, when deploying our method. Our released code comes with usage instructions to promote safe adoption and reduce the risk of misuse. This work is licensed under CC BY 4.0, allowing reuse and adaptation, even commercially, with proper attribution.

## AI Assistants in Research Writing

We used AI assistants to support writing and code refinement during the preparation of this paper. All AI-generated content was reviewed and verified by the authors.

## References

- Matan Avitan, Ryan Cotterell, Yoav Goldberg, and Shauli Ravfogel. 2024. What changed? converting representational interventions to natural language. *arXiv preprint arXiv:2402.11355*.
- Jiamu Bai, Daoyuan Chen, Bingchen Qian, Liuyi Yao, and Yaliang Li. 2024. Federated fine-tuning of

- large language models under heterogeneous tasks and client resources. *Advances in Neural Information Processing Systems*, 37:14457–14483.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, and Yejin Choi. 2020. **Piqa: Reasoning about physical commonsense in natural language**. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 7432–7439.
- Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. 2017. Machine learning with adversaries: Byzantine tolerant gradient descent. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Dongqi Cai, Yaozong Wu, Shangguang Wang, Felix Xiaozhu Lin, and Mengwei Xu. 2023. Efficient federated learning for modern nlp. In *Proceedings of the Annual International Conference on Mobile Computing and Networking*, pages 1–16.
- Jinyu Chen, Wenchao Xu, Song Guo, Junxiao Wang, Jie Zhang, and Haozhao Wang. 2022. Fedtune: A deep dive into efficient federated fine-tuning with pre-trained transformers. *arXiv preprint arXiv:2211.08025*.
- Lili Chen, Lijun Su, and Jinhui Xu. 2017. Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. *arXiv preprint arXiv:1705.10301*.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. **Boolq: Exploring the surprising difficulty of natural yes/no questions**. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, pages 2924–2936.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, and 1 others. 2023. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235.
- Wenzhi Fang, Dong-Jun Han, Liangqi Yuan, Seyyedali Hosseinalipour, and Christopher G Brinton. 2025. Federated sketching lora: On-device collaborative fine-tuning of large language models. *arXiv preprint arXiv:2501.19389*.
- Shangqian Gao, Ting Hua, Yen-Chang Hsu, Yilin Shen, and Hongxia Jin. 2024. Adaptive rank selections for low-rank approximation of language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 227–241.
- Atticus Geiger, Zhengxuan Wu, Christopher Potts, Thomas Icard, and Noah Goodman. 2024. Finding alignments between interpretable causal variables and distributed neural representations. In *Causal Learning and Reasoning*, pages 160–187. PMLR.
- Pengxin Guo, Shuang Zeng, Yanran Wang, Huijie Fan, Feifei Wang, and Liangqiong Qu. 2024. Selective aggregation for low-rank adaptation in federated learning. *arXiv preprint arXiv:2410.01463*.
- Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2021. Towards a unified view of parameter-efficient transfer learning. *arXiv preprint arXiv:2110.04366*.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *ICML*.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021a. **Lora: Low-rank adaptation of large language models**. *arXiv preprint ArXiv:2106.09685*.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021b. **Lora: Low-rank adaptation of large language models**. *arXiv preprint arXiv:2106.09685*.
- Zhiqiang Hu, Lei Wang, Yihuai Lan, Wanyu Xu, Ee-Peng Lim, Lidong Bing, Xing Xu, Soujanya Poria, and Roy Ka-Wei Lee. 2023. Llm-adapters: An adapter family for parameter-efficient fine-tuning of large language models. *arXiv preprint arXiv:2304.01933*.
- Rik Koncel-Kedziorski, Subhro Roy, Tao Zhang, and Hannaneh Hajishirzi. 2016. Mawps: A math word problem repository. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1152–1157.
- Weirui Kuang, Bingchen Qian, Zitao Li, Daoyuan Chen, Dawei Gao, Xuchen Pan, Yuexiang Xie, Yaliang Li, Bolin Ding, and Jingren Zhou. 2024. Federatedscope-llm: A comprehensive package for fine-tuning large language models in federated learning. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 5260–5271.
- Kenneth Lange. 2016. *MM optimization algorithms*. SIAM.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In *EMNLP*.
- Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In *KR*.

- Dengchun Li, Yingzi Ma, Naizheng Wang, Zhiyuan Cheng, Lei Duan, Jie Zuo, Cal Yang, and Mingjie Tang. 2024a. Mixlor: Enhancing large language models fine-tuning with lora based mixture of experts. *arXiv preprint arXiv:2404.15159*.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2024b. Inference-time intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36.
- Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. In *ACL*.
- Vladislav Lialin, Sherin Muckatira, Namrata Shivagunde, and Anna Rumshisky. 2023. Relora: High-rank training through low-rank updates. In *The Twelfth International Conference on Learning Representations*.
- Wang Ling, Dani Yogatama, Chris Dyer, and Philip Blunsom. 2017. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, page 158–167.
- Sheng Liu, Haotian Ye, Lei Xing, and James Zou. 2023. In-context vectors: Making in context learning more effective and controllable through latent space steering. *arXiv preprint arXiv:2311.06668*.
- Zefang Liu and Jiahua Luo. 2024. Adamole: Fine-tuning large language models with adaptive mixture of low-rank adaptation experts. *arXiv preprint arXiv:2405.00361*.
- Guodong Long, Tao Shen, Jing Jiang, Michael Blumstein, and 1 others. 2024. Dual-personalizing adapter for federated foundation models. *Advances in Neural Information Processing Systems*, 37:39409–39433.
- Qikai Lu, Di Niu, Mohammadamin Samadi Khoshkho, and Baochun Li. 2024. Hyperflora: Federated learning with instantaneous personalization. In *Proceedings of the 2024 SIAM International Conference on Data Mining (SDM)*, pages 824–832.
- Ricardo A Maronna and Douglas Martin. 2006. Yohai robust statistics. *Wiley Series in Probability and Statistics. John Wiley and Sons*, 2:3.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguerre y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR.
- J. Pablo Muñoz, Jinjie Yuan, and Nilesh Jain. 2025. [Low-rank adapters meet neural architecture search for llm compression](#). In *AAAI'25 workshop on CoLoRAI - Connecting Low-Rank Representations in AI*.
- Yahao Pang, Xingyuan Wu, Xiaojin Zhang, Wei Chen, and Hai Jin. 2025. Fedeat: A robustness optimization framework for federated llms. *arXiv preprint arXiv:2502.11863*.
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are nlp models really able to solve simple math word problems? *arXiv preprint arXiv:2103.07191*.
- Krishna Pillutla, Sham M Kakade, and Zaid Harchaoui. 2022. Robust aggregation for federated learning. *IEEE Transactions on Signal Processing*, 70:1142–1154.
- Hossein Rajabzadeh, Mojtaba Valipour, Tianshu Zhu, Marzieh Tahaei, Hyock Ju Kwon, Ali Ghodsi, Boxing Chen, and Mehdi Rezagholizadeh. 2024. Qdy-lora: Quantized dynamic low-rank adaptation for efficient large language model tuning. *arXiv preprint arXiv:2402.10462*.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavathula, and Yejin Choi. 2021. [Winogrande: An adversarial winograd schema challenge at scale](#). *Communications of the ACM*, 64(9):99–106.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. 2019. Socialiqa: Commonsense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*.
- Shashwat Singh, Shauli Ravfogel, Jonathan Herzig, Roei Aharoni, Ryan Cotterell, and Ponnurangam Kumaraguru. 2024. Mimic: Minimally modified counterfactuals in the representation space. *arXiv preprint arXiv:2402.09631*.
- Raghav Singhal, Kaustubh Ponske, Rohit Vartak, Lav R Varshney, and Praneeth Vepakomma. 2025. Fed-sb: A silver bullet for extreme communication efficiency and performance in (private) federated lora fine-tuning. *arXiv preprint arXiv:2502.15436*.
- Guangyu Sun, Matias Mendieta, Taojiannan Yang, and Chen Chen. 2022. Exploring parameter-efficient fine-tuning for improving communication efficiency in federated learning. In *International Conference on Learning Representations (ICLR)*.
- Youbang Sun, Zitao Li, Yaliang Li, and Bolin Ding. 2024. Improving lora in privacy-preserving federated learning. *arXiv preprint arXiv:2403.12313*.
- Van-Tuan Tran, Quoc-Viet Pham, and 1 others. 2025. Revisiting sparse mixture of experts for resource-adaptive federated fine-tuning foundation models. In *ICLR 2025 Workshop on Modularity for Collaborative, Decentralized, and Continual Deep Learning*.
- Mojtaba Valipour, Mehdi Rezagholizadeh, Ivan Kobyzev, and Ali Ghodsi. 2022. Dylora: Parameter efficient tuning of pre-trained models using dynamic search-free low-rank adaptation. *arXiv preprint arXiv:2210.07558*.

- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355.
- Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. 2021. **TENT: Fully test-time adaptation by entropy minimization**. In *International Conference on Learning Representations (ICLR)*.
- Yaqing Wang, Sahaj Agarwal, Subhabrata Mukherjee, Xiaodong Liu, Jing Gao, Ahmed Hassan Awadallah, and Jianfeng Gao. 2022. Adamix: Mixture-of-adaptations for parameter-efficient model tuning. *arXiv preprint arXiv:2205.12410*.
- Ziyao Wang, Zheyu Shen, Yexiao He, Guoheng Sun, Hongyi Wang, Lingjuan Lyu, and Ang Li. 2024. Flora: Federated fine-tuning large language models with heterogeneous low-rank adaptations. *arXiv preprint arXiv:2409.05976*.
- Endre Weiszfeld. 1937. Sur le point pour lequel la somme des distances de  $n$  points donnés est minimum. *Tohoku Mathematical Journal, First Series*, 43:355–386.
- Feijie Wu, Zitao Li, Yaliang Li, Bolin Ding, and Jing Gao. 2024a. Fedbiot: Llm local fine-tuning in federated learning without full model. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 3345–3355.
- Xun Wu, Shaohan Huang, and Furu Wei. 2024b. Mixture of lora experts. *arXiv preprint arXiv:2404.13628*.
- Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D. Manning, and Christopher Potts. 2024c. **Refit: Representation fine-tuning for language models**. *arXiv preprint ArXiv:2404.03592*.
- Yunlu Yan, Chun-Mei Feng, Wangmeng Zuo, Rick Siow Mong Goh, Yong Liu, and Lei Zhu. Federated residual low-rank adaptation of large language models. In *The Thirteenth International Conference on Learning Representations*.
- Yifan Yang, Kai Zhen, Ershad Banijamal, Athanasios Mouchtaris, and Zheng Zhang. 2024. Adazeta: Adaptive zeroth-order tensor-train adaption for memory-efficient large language models fine-tuning. *arXiv preprint arXiv:2406.18060*.
- Liping Yi, Han Yu, Gang Wang, and Xiaoguang Liu. 2023. Fedlora: Model-heterogeneous personalized federated learning with lora tuning. *arXiv preprint arXiv:2310.13283*.
- Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. 2018. Byzantine-robust distributed learning: Towards optimal statistical rates. In *International conference on machine learning*, pages 5650–5659. Pmlr.
- Ben Zaken, Yoav Goldberg, and Amir Globerson. 2022. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. In *ACL Findings*.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. **Hellaswag: Can a machine really finish your sentence?** *arXiv preprint arXiv:1905.07830*.
- Jianyi Zhang, Saeed Vahidian, Martin Kuo, Chunyuan Li, Ruiyi Zhang, Tong Yu, Guoyin Wang, and Yiran Chen. 2024. Towards building the federatedgpt: Federated instruction tuning. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6915–6919. IEEE.
- Zhuo Zhang, Yuanhang Yang, Yong Dai, Qifan Wang, Yue Yu, Lizhen Qu, and Zenglin Xu. 2023. Fedpetuning: When federated learning meets the parameter-efficient tuning methods of pre-trained language models. In *Annual Meeting of the Association of Computational Linguistics (ACL)*, pages 9963–9977.
- Changhai Zhou, Shijie Han, Shiyang Zhang, Shichao Weng, Zekai Liu, and Cheng Jin. 2024. Rankadaptor: Hierarchical dynamic low-rank adaptation for structural pruned llms. *arXiv preprint arXiv:2406.15734*.
- Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*.
- Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, and 1 others. 2023. Representation engineering: A top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*.

## A Related Works

### A.1 Parameter-Efficient Fine-Tuning (PEFT)

Fine-tuning LLMs is resource-intensive due to their large parameter counts. Parameter-efficient fine-tuning (PEFT) methods mitigate this by updating only a small subset of parameters while keeping pre-trained weights frozen (Li and Liang, 2021; He et al., 2021; Wang et al., 2022). Several PEFT approaches have been proposed, Adapter Tuning (Houlsby et al., 2019), BitFit (Zaken et al., 2022), Prefix Tuning (Li and Liang, 2021), Prompt Tuning (Lester et al., 2021), and Low-Rank Adaptation (LoRA) (Hu et al., 2021a). Among them,

LoRA is widely adopted for its efficiency in approximating weight updates via low-rank matrices. Extensions such as ReLoRA (Lialin et al., 2023) and RankAdapter (Zhou et al., 2024) improve memory use and adapt ranks dynamically, though they lack theoretical guarantees. AdaZeta (Yang et al., 2024) introduces zeroth-order optimization with convergence guarantees, while others (Gao et al., 2024; Rajabzadeh et al., 2024; Valipour et al., 2022) explore adaptive ranks without formal proofs. LoRA has been integrated with Mixture-of-Experts models (Li et al., 2024a; Wu et al., 2024b), as in AdaMoLE (Liu and Luo, 2024), to enable dynamic expert selection, and also with Neural Architecture Search for LLM compression (Muñoz et al., 2025). These approaches primarily target weight updates, overlooking direct interventions in hidden representations, which are discussed next.

## A.2 Representation Fine-Tuning (ReFT)

ReFT shifts fine-tuning from model weights to hidden representations, leveraging their semantic structure for efficient adaptation (Wu et al., 2024c). Inspired by activation steering and representation engineering (Avitan et al., 2024; Li et al., 2024b; Liu et al., 2023; Singh et al., 2024), ReFT enables task-specific control through fixed or learned interventions without updating the full model. Notably, Inference-Time Intervention (ITI) (Li et al., 2024b) improves LLM truthfulness by modifying activations, while representation engineering (Zou et al., 2023) combines representation reading and control for interpretable model behavior. Minimally Modified Counterfactuals (MMC) (Singh et al., 2024) unify erasure and steering to reduce bias, and can be mapped to natural language edits (Avitan et al., 2024), enhancing interpretability. These findings support direct representation manipulation as a lightweight and effective alternative to weight-based PEFT methods like LoRA.

## A.3 Federated Fine-Tuning

Federated Learning (FL) (McMahan et al., 2017) poses challenges for fine-tuning LLMs, including data heterogeneity, communication constraints, and model diversity. PEFT methods have emerged to address these issues efficiently (Sun et al., 2022; Chen et al., 2022; Zhang et al., 2023). LoRA-based approaches such as FedLoRA (Yi et al., 2023), Hyper-FloRA (Lu et al., 2024), and Efficient FL Adapter (Cai et al., 2023) offer modular and personalized adaptation across clients. Recent advances

further incorporate privacy (FFA-LoRA (Sun et al., 2024)), heterogeneous adaptation (FloRA (Wang et al., 2024)), instruction tuning (FedIT (Zhang et al., 2024)), residual learning in FRLoRA (Yan et al.) which tackles client drift by directly adding residual low-rank weight products to the global model parameters in each round, and heterogeneous resources in FlexLoRA (Bai et al., 2024) which mitigates the bucket effect by leveraging SVD to aggregate and redistribute LoRA weights with varying ranks among clients, and model compression in FedBiOT (Wu et al., 2024a) which enables LLM fine-tuning without requiring clients to access the full model by using a bi-level optimization scheme to align a compressed emulator with a lightweight adapter, and adaptive sketching in FSLoRA (Fang et al., 2025) which reduces communication and computation overhead by allowing clients to selectively update submatrices of the global LoRA modules based on sketching ratios, and expert routing (DualFed (Long et al., 2024), Sparse-FedMoE (Tran et al., 2025)). In contrast, our proposed FEDREFT shifts from weight updates to direct representation-level tuning via sparse intervention layers and introduces an All-But-Me (ABM) aggregation strategy to preserve semantic alignment while enabling robust knowledge sharing across non-IID clients.

## A.4 Aggregation Methods in FL

To address the inherent heterogeneity and robustness challenges in federated learning, median-based aggregation strategies have been extensively studied as alternatives to simple averaging. Unlike the arithmetic mean, the geometric and coordinate-wise medians are significantly more resilient to outliers and adversarial updates, making them suitable for secure and personalized FL scenarios. For instance, coordinate-wise median aggregation has been proposed to defend against Byzantine clients in distributed optimization (Blanchard et al., 2017). This was extended with geometric median-based gradient descent to improve statistical guarantees across diverse loss landscapes (Yin et al., 2018). Further work demonstrated that coordinate-wise median and trimmed-mean-based methods achieve order-optimal convergence not only for strongly convex losses but also under non-strongly convex and even non-convex population losses (Chen et al., 2017). Additionally, a one-round median-based algorithm was shown to maintain statistical optimality under quadratic convexity, offering

a communication-efficient solution (Chen et al., 2017). RFA (Pillutla et al., 2022) maintains privacy and demonstrates improved robustness over standard averaging techniques, particularly in environments with high levels of data corruption. FedEAT (Pang et al., 2025) integrates adversarial training in the embedding space with geometric median-based aggregation to enhance robustness while preserving performance. This work demonstrates that LoRA-based FL systems can effectively leverage geometric median aggregation. Inspired by these findings, we adopt geometric median aggregation in our FL framework to aggregate the All-But-Me (ABM) intervention parameter, weight  $\mathbf{W}$ , rotation  $\mathbf{R}$ , and bias  $\mathbf{b}$ . This provides stability across diverse client behaviors and loss geometries, improving personalization performance under data and objective heterogeneity.

## B Hyperparameter Search Space

### B.1 Hyperparameter Search Space for Commonsense and Arithmetic Reasoning

Following the ReFT framework (Wu et al., 2024c), we construct a development set using the GSM8K dataset and consider only the last 300 samples. We trained the clients using LLaMA 7B model with the remaining training data and determined the best-performing hyperparameters based on the model’s performance on the development set. We further use this hyperparameter in another model directly. We set the maximum input sequence length to 512 tokens during training and tuning, and limit inference to 32 generated tokens. We use the same setup for commonsense reasoning with Commonsense170K dataset. The hyperparameter search space is summarized in Tables 9 and 10.

During inference, we use greedy decoding (without sampling) for the commonsense reasoning benchmark, as it is a multi-token classification task. For arithmetic reasoning, we follow the decoding setup from (Hu et al., 2023), using a higher temperature of 0.3. This change helps avoid errors in HuggingFace’s decoding caused by unstable probabilities

### B.2 Hyperparameter Search Space for GLUE Benchmark

We perform hyperparameter tuning following common practice for PEFT methods (Hu et al., 2023) in FL on RoBERTa separately for each GLUE task in Table 11, selecting the optimal settings based

Table 9: Narrow down the hyperparameter(HP) search space of LLaMA 7B models with FedReFT on the GSM8K development set, inspired from (Wu et al., 2024c). The best-performing settings are underlined. We apply greedy decoding without sampling during hyperparameter tuning.

HP	FedReFT
prefix+suffix, $p + s$	{p5+s5, <u>p7+s7</u> , p9+s9}
Tied weight $\phi$	{True, <u>False</u> }
Rank $r$	{ <u>8</u> , 16, 32, 64}
Layer $L$	{ <u>all</u> }
Dropout	{ <u>0.00</u> , 0.05}
Optimizer	AdamW
LR	{6, <u>9</u> } $\times 10^{-4}$
Weight decay	{ <u>0</u> , $1 \times 10^{-3}$ , $2 \times 10^{-3}$ }
LR scheduler	Linear
Batch size	{ <u>16</u> , 32}
Warmup ratio	{0.06, <u>0.10</u> }
Clients	{3, 5}
Epochs	{3, 4, <u>5</u> , 6}
Rounds	10

Table 10: Narrow down the hyperparameter (HP) search space of LLaMA 7B models with FedReFT on the Commonsense170K development set, following the Appendix B.1. The best-performing settings are underlined. We apply greedy decoding without sampling during hyperparameter tuning.

HP	FedReFT
prefix+suffix, $p + s$	{p5+s5, <u>p7+s7</u> }
Tied weight $p, s$	{True, <u>False</u> }
Rank $r$	{ <u>8</u> , 16, 32, 64}
Layer $L$	{ <u>all</u> }
Dropout	{ <u>0.00</u> , 0.05}
Optimizer	AdamW
LR	{4, <u>6</u> , 9} $\times 10^{-4}$
Weight decay	{ <u>0</u> }
LR scheduler	Linear
Batch size	{ <u>16</u> , 32}
Warmup ratio	{ <u>0.1</u> }
Clients	{3, 5}
Epochs	{2, <u>3</u> , 4}
Rounds	10

on validation performance using a fixed random seed of 42. Final evaluations are conducted using two additional unseen seeds, {43, 44}, to ensure robustness.

Table 11: Hyperparameter(HP) settings of RoBERTa-large models on selected GLUE tasks for FedReFT, inspired from (Wu et al., 2024c)

HP	MNLI	SST-2	QNLI	QQP
position $p$	$p1$	$p3$	$p11$	$p11$
Tied weight		False		
Rank $r$		1		
Layer $L$		all		
Dropout		0.05		
Optimizer		AdamW		
LR		$6 \times 10^{-4}$		
Weight decay		0.00		
LR scheduler		Linear		
Batch size		32		
Warmup ratio	0.00	0.10	0.10	0.06
Epochs		5		
Rounds		50		

## C Theoretical Foundation: Geometric Median via Weiszfeld’s Algorithm

The geometric median offers a robust alternative to the arithmetic mean, particularly suitable for federated settings with heterogeneous or noisy client updates. For a given set of vectors  $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subset \mathbb{R}^d$ , the geometric median  $\mathbf{y}^*$  is defined as:

$$\mathbf{y}^* = \arg \min_{\mathbf{y} \in \mathbb{R}^d} \sum_{i=1}^n \|\mathbf{y} - \mathbf{x}_i\|_2. \quad (11)$$

This optimization is non-smooth and convex, and generally lacks a closed-form solution. However, Weiszfeld’s algorithm (Weiszfeld, 1937) provides an efficient iterative method to approximate  $\mathbf{y}^*$ . We now derive and justify this algorithm via the Majorization-Minimization (MM) framework.

We define the cost function to be minimized: This function is convex but non-differentiable at points where  $\mathbf{y} = \mathbf{x}_i$ . Weiszfeld’s algorithm avoids such points during updates by construction.

The MM algorithm minimizes a difficult objective  $f(\mathbf{y})$  by iteratively minimizing a surrogate function  $Q(\mathbf{y}|\mathbf{y}^{(k)})$  that: Majorizes  $f$ :

$Q(\mathbf{y}|\mathbf{y}^{(k)}) \geq f(\mathbf{y})$  for all  $\mathbf{y}$ , Touches  $f$  at the current iterate:  $Q(\mathbf{y}^{(k)}|\mathbf{y}^{(k)}) = f(\mathbf{y}^{(k)})$ .

We define the surrogate using Jensen’s inequality and the convexity of the norm:

$$Q(\mathbf{y}|\mathbf{y}^{(k)}) = \sum_{i=1}^n \frac{\|\mathbf{y} - \mathbf{x}_i\|_2^2}{2\|\mathbf{y}^{(k)} - \mathbf{x}_i\|_2} + C(\mathbf{y}^{(k)}), \quad (12)$$

where  $C(\mathbf{y}^{(k)})$  is a constant that does not depend on  $\mathbf{y}$ . This function is differentiable and strictly convex in  $\mathbf{y}$ .

To find the minimizer of  $Q(\mathbf{y}|\mathbf{y}^{(k)})$ , we take the gradient and set it to zero:

$$\nabla Q(\mathbf{y}) = \sum_{i=1}^n \frac{\mathbf{y} - \mathbf{x}_i}{\|\mathbf{y}^{(k)} - \mathbf{x}_i\|_2} = 0. \quad (13)$$

Solving the above yields the Weiszfeld update rule:

$$\mathbf{y}^{(k+1)} = \frac{\sum_{i=1}^n \frac{\mathbf{x}_i}{\|\mathbf{y}^{(k)} - \mathbf{x}_i\|_2}}{\sum_{i=1}^n \frac{1}{\|\mathbf{y}^{(k)} - \mathbf{x}_i\|_2}}. \quad (14)$$

The update is only valid when  $\mathbf{y}^{(k)} \neq \mathbf{x}_i$  for all  $i$ , a condition that can be enforced by initialization and step-size dampening if needed. From MM theory (Lange, 2016), each iteration satisfies:

$$\begin{aligned} f(\mathbf{y}^{(k+1)}) &\leq Q(\mathbf{y}^{(k+1)}|\mathbf{y}^{(k)}) \\ &\leq Q(\mathbf{y}^{(k)}|\mathbf{y}^{(k)}) = f(\mathbf{y}^{(k)}), \end{aligned} \quad (15)$$

ensuring that  $f(\mathbf{y}^{(k)})$  is non-increasing. Under mild conditions (excluding cases where  $\mathbf{y}^{(k)} = \mathbf{x}_i$ ), Weiszfeld’s algorithm converges to the geometric median  $\mathbf{y}^*$ .

### C.1 All-But-Me (ABM) Aggregation Strategy

In our FedReFT framework, each client receives an All-But-Me (ABM) aggregated update for intervention parameters computed as the geometric median of the corresponding parameters from all other clients. For client  $k$ , the ABM aggregated parameter is:

$$\mathbf{W}_k^{\text{ABM}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{m \neq k} \|\mathbf{w} - \mathbf{W}_m^{\text{local}}\|_2. \quad (16)$$

We compute this using Weiszfeld’s algorithm for each parameter type independently, ensuring robustness to outlier clients and misaligned updates. This enables stable and personalized aggregation without sacrificing task-specific semantics. Weiszfeld’s algorithm provides a theoretically grounded and computationally efficient way

to compute the geometric median, making it ideal for ABM aggregation in heterogeneous FL. By leveraging this algorithm in FedReFT, we ensure robustness in aggregation and improve both convergence and personalization in non-i.i.d. federated environments.

## C.2 Geometric Median over mean

The Geometric Median used in All-but-Me (ABM) aggregation is computationally more expensive, with a time complexity of  $O(T \cdot d)$  due to its iterative Weiszfeld optimization and the need to accumulate all  $N$  client updates. In this paper, the dimension is only  $d = 20$ , which makes the overhead relatively small. In contrast, the arithmetic mean used in FedAvg has a lower complexity of  $O(d)$ . However, the slightly higher cost of the Geometric Median is justified by its stronger robustness: by minimizing the  $L_1$  error, it becomes significantly more resilient to client heterogeneity and outlier updates. This leads to consistently higher accuracy, with clear performance improvements observed over FedAvg on the Commonsense task. While outlier client updates can render the aggregated FedAvg model unstable or even unusable, the Geometric Median ensures a stable and reliable update, producing a more accurate final model.

For a heterogeneous task setting of a commonsense task using 8 clients, the Geometric Median provides superior outlier detection in Figure 3. The Distance from GeoMed plot visually proves Client 7 is an extreme outlier, clearly separating it from the benign majority’s distance. This maximum separation occurs because GeoMed anchors near the majority by  $L_1$  linear error minimization, resisting the outlier’s pull. Conversely, the Mean plot shows the Mean center is contaminated, minimizing the outlier’s distance because the Mean minimizes  $L_2$  squared error in figure 4. This distortion makes the aggregated model unusable, but GeoMed’s robust calculation ensures a stable, accurate final update.

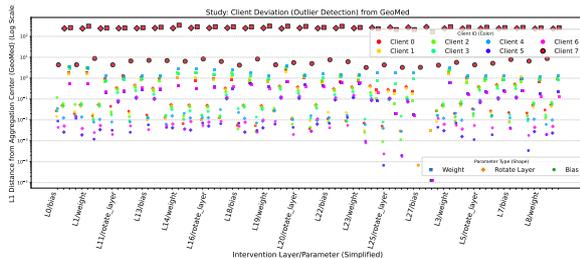


Figure 3: Study\_Distance\_from\_GeoMed\_ShapeColor

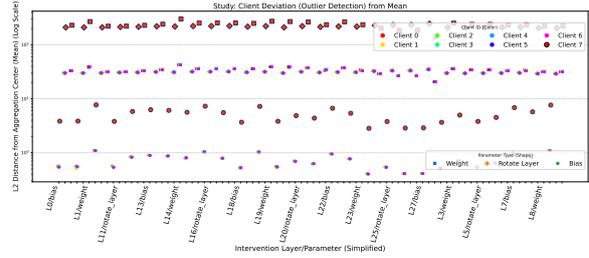


Figure 4: Study\_Distance\_from\_Mean\_ShapeColor

## D Ablation Study

### D.1 Intervention Parameter Sharing Strategy

To reduce communication overhead while maintaining personalization, we explore three strategies for sharing local intervention parameters with the server. These strategies represent different trade-offs between expressiveness and communication efficiency:

- **Full Intervention Sharing:**  $\{W \in \mathbb{R}^{r \times d}, R \in \mathbb{R}^{r \times d}, b \in \mathbb{R}^r\}$  This strategy shares the complete set of intervention parameters, capturing client-specific compression ( $W$ ), transformation ( $R$ ), and translation ( $b$ ). It enables the most accurate reconstruction of local updates and yields the best global performance, especially under high heterogeneity.
- **No Bias Sharing:**  $\{W \in \mathbb{R}^{r \times d}, R \in \mathbb{R}^{r \times d}\}$  This variant omits the bias term  $b$  but retains the directional transformation via  $W$  and  $R$ . While it allows the server to align low-rank subspace transformations across clients, it lacks the ability to model per-dimension translation shifts, which can hinder fine-grained personalization.
- **No  $W$  Sharing:**  $\{R \in \mathbb{R}^{r \times d}, b \in \mathbb{R}^r\}$ . This configuration excludes  $W$ , giving the server access only to the reconstruction and shift parameters. Without knowledge of how the local signals were encoded, the server’s ability to interpret or align updates is severely limited.

The  $\{W, R, b\}$  strategy provides the highest fidelity for aggregation,  $\{W, R\}$  offers a balanced compromise, and  $\{R, b\}$  prioritizes communication efficiency at the cost of semantic alignment and global performance.

## D.2 Integrating Test-Time Computing (TTC) with Adaptive Mixing in Local Model Updates

In FedReFT, local interventions encode client-specific semantics, while ABM aggregation conveys shared global knowledge. However, naively averaging representation-level updates can disrupt semantic consistency across clients, leading to misaligned feature spaces. Furthermore, uniform aggregation disregards client heterogeneity, reducing adaptability in diverse edge environments.

To address these challenges, we evaluate the effectiveness of FedReFT under two mixing strategies: (i) adaptive mixing via Test-Time Computing (TTC), where local and global interventions are dynamically balanced, and (ii) balanced mixing, where a fixed coefficient  $\alpha$  governs the trade-off between local ( $\alpha$ ) and global ( $1 - \alpha$ ) contributions. While TTC provides a dynamic, task-adaptive mechanism expected to outperform balanced mixing ( $\alpha = 0.5$ ), we also include results for the fixed  $\alpha = 0.5$  setting to illustrate this trade-off. Table 12, 13, 14, and 15 show that the performance of TTC-based adaptive mixing is over the balanced mixing.

In equation 9, the Test-Time Computing (TTC) optimization uses three small regularization weights,  $\lambda_1, \lambda_2, \lambda_3$ , all set to 0.001. This low value means the primary goal is minimizing the test loss  $\mathcal{L}_{\text{test}}$ , not strict following of the regularization rules.  $\lambda_1$  Entropy and  $\lambda_3$  Diversity gently stop the mixing coefficient  $\alpha$  from completely favoring local (1) or global (0) updates.  $\lambda_2$  Consistency mildly allows  $\alpha$  to change across different intervention layers, which helps with adaptive personalization. Because the  $\lambda$  values are small, the resulting mean  $\alpha \approx 0.57$  strongly favors local updates based on the client’s specific data, but not the fully Local-Only assumption of  $\alpha = 1.0$ . If these  $\lambda$  values were increased significantly, the regularization would become more important than the test loss. This would push  $\alpha$  closer to 0.5 for balanced mixing across all layers, reducing the model’s ability to personalize.

## D.3 Comparison of Aggregation Methods on Different Tasks

We use training data from the COMMONSENSE170K dataset, split among three clients, and evaluate the models using the SIQA task. Figure 5 shows that Geometric Median ABM aggregation outperforms all other approaches. Similarly, we

split the MATH10K dataset among three clients, train each client for only five local epochs, and evaluate the results using the GSM8K evaluation set. Additionally, we use the QNLI GLUE dataset for the natural language understanding (NLU) task. **Why small gains matter.** We emphasize that even modest accuracy improvements are meaningful in the federated learning setting, particularly under highly heterogeneous task distributions. As shown in Figure 5 the arithmetic mean has time complexity  $O(d)$ , while the geometric median (Weiszfeld’s algorithm) has  $O(T \cdot d)$ , with  $T$  as iterations and  $d$  as parameters. Both methods have memory complexity  $O(d)$ . Despite the higher computational cost, the geometric median provides greater robustness to heterogeneity and consistently yields higher accuracy in FL settings.

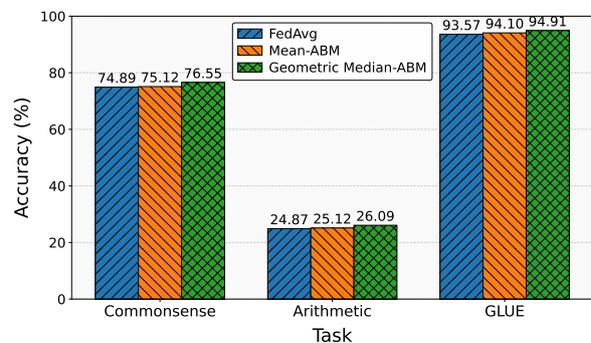


Figure 5: Comparison of aggregation strategies across tasks. Results for FedAvg, Mean-ABM, and Geometric Median-ABM on Commonsense Reasoning, Arithmetic Reasoning, and GLUE show that Geometric Median-ABM consistently outperforms others, demonstrating greater robustness in heterogeneous federated settings.

These improvements hold consistently across three diverse tasks. We chose Geometric Median-ABM not only for its peak accuracy but also for its robustness to outliers and task drift. The benefits are especially pronounced in arithmetic reasoning, where client diversity is highest. While geometric median incurs additional server-side computation, it is applied once per round on low-dimensional sparse parameters, making the overhead negligible. Client-side operations remain lightweight.

$$f(\mathbf{y}) = \sum_{i=1}^n \|\mathbf{y} - \mathbf{x}_i\|_2. \quad (17)$$

## E Communication Efficiency

As shown in Table 17, FedReFT is communication and computationally efficient as it uses only a very

Table 12: Federated fine-tuning performance of LLaMA-3.2 3B across five commonsense reasoning tasks with the Mixed-Task (MT) setup, where clients train on heterogeneous task mixtures. We report the effectiveness of **FedReFT** under two mixing strategies: (i) **adaptive mixing** using Test-Time Computing (TTC), where local and global interventions are dynamically balanced, and (ii) **balanced mixing**, where a coefficient  $\alpha$  controls the trade-off ( $\alpha$  for local and  $1 - \alpha$  for global). While TTC is intended to provide a dynamic, task-adaptive solution expected to surpass balanced mixing ( $\alpha = 0.5$ ), we also include results for  $\alpha = 0.5$  to illustrate the trade-off. **R** = Rank, **TP(M)** = Trainable Parameters in Millions, and **TP Effi. (R8)** = parameter efficiency of FedReFT with rank 8.

Method	R	TP(M) ↓	TP Effi. (R8) ↓	BoolQ	PIQA	SIQA	HellaS.	WinoG	Avg Acc ↑
<b>FedReFT (<math>\alpha=0.5</math>)</b>	4	1.38	0.5×	63.09	82.10	72.36	90.27	69.22	75.41
	8	2.75	1.0×	64.01	81.18	72.11	90.71	71.01	75.66
	16	5.5	0.5×	63.42	81.61	73.64	91.23	71.35	76.05
	32	11.01	0.25×	64.53	81.34	73.39	91.51	71.32	76.22
<b>FedReFT (tie <math>\phi</math>)</b>	4	0.688	0.0214	49.94	81.23	72.72	89.84	68.43	72.43
	8	1.38	0.0428	57.15	81.22	72.77	90.56	68.50	74.04
<b>FedReFT (TTC)</b>	4	1.38	0.5×	63.35	82.72	72.96	91.37	69.70	76.02
	8	2.75	1.0×	65.50	82.32	73.28	91.43	70.24	76.55
	<b>16</b>	<b>5.5</b>	<b>1.0×</b>	<b>64.56</b>	<b>82.20</b>	<b>75.95</b>	<b>89.80</b>	<b>80.59</b>	<b>78.62</b>

Table 13: Performance comparison across GLUE tasks on RoBERTa with  $C = 3$ . We report results under two mixing strategies: (i) **adaptive mixing** with Test-Time Computing (TTC), which dynamically balances local and global interventions at inference, and (ii) **balanced mixing**, where  $\alpha$  denotes the proportion of local contribution and  $1 - \alpha$  the global contribution. The balanced  $\alpha = 0.5$  cases illustrate trade-offs between local personalization and global knowledge, while TTC is intended as a dynamic alternative expected to generalize better across heterogeneous conditions. **TP(M)** denotes the number of trainable parameters (in millions).

Method	TP(M) ↓	MNLI-m	MNLI-mm	SST-2	QNLI	QQP	Avg ↑
FedReFT ( $\alpha=0.5$ )	0.053	88.86	88.76	95.17	94.52	86.57	90.93
<b>FedReFT (TTC)</b>	<b>0.053</b>	<b>89.75</b>	<b>89.31</b>	<b>95.75</b>	<b>94.91</b>	<b>87.15</b>	<b>91.37</b>

Table 14: Performance comparison across arithmetic reasoning tasks with the Distinct Task (DT) and Mixed Task (MT) setup using different model sizes. We report results under two mixing strategies: (i) **adaptive mixing** with Test-Time Computing (TTC), which dynamically balances local and global interventions at inference, and (ii) **balanced mixing**, where  $\alpha = 0.5$  denotes the proportion of local contribution and  $1 - \alpha$  the global contribution.

<b>FedReFT (<math>\alpha = 0.5</math>)</b>	<b>Distinct Task (DT)</b>				<b>Mixed Task (MT)</b>			
	AQuA	SVAMP	MAWPS	Avg ↑	AQuA	SVAMP	MAWPS	Avg ↑
LLaMA 7B	25.59	25.47	49.80	33.62	22.83	14.33	27.10	21.42
LLaMA-2 7B	29.53	32.45	57.30	39.76	21.65	20.39	31.50	24.51
LLaMA-3 8B	34.64	48.98	73.60	52.41	31.89	48.90	70.04	50.48
<b>FedReFT (TTC)</b>	<b>Distinct Task (DT)</b>				<b>Mixed Task (MT)</b>			
LLaMA 7B	26.12	26.71	49.94	34.26	20.86	15.60	28.22	22.23
LLaMA-2 7B	30.96	33.31	59.51	41.93	22.05	23.50	32.74	26.09
LLaMA-3 8B	35.36	49.41	75.25	53.34	33.45	51.28	73.51	52.75

Table 15: Performance comparison on arithmetic reasoning tasks for GSM8K on LLaMA-3 8B model with LoRA rank 8. **Trainable Parameter (TP) Efficiency** indicates the efficiency of FedReFT compared to the baselines.  $\Delta\text{Acc}$  shows the accuracy improvement of FedReFT.

Method	TP(M)↓	TP Effi.↓	Acc	$\Delta\text{Acc}$ ↑
FedReFT (tie $\phi$ )	2.09	0.5×	47.27	+1.61
FedReFT $\alpha = 0.5$	4.19	1.0×	48.39	+0.49
<b>FedReFT (TTC)</b>	<b>4.19</b>	<b>1.0×</b>	<b>48.88</b>	

small percentage of trainable parameters (TP) compared to the total model parameters. For example, in LLaMA-7B and LLaMA-2 7B, only 0.0311% of the total parameters are trained. In RoBERTa Large, this number is even smaller, at just 0.0138%. Even for large models like LLaMA-2-13B, the trainable portion remains as low as 0.0503%. This shows that FedReFT is highly parameter-efficient. Despite using such a small fraction of parameters, FedReFT still achieves strong performance, as discussed in the experimental analysis section 4. This highlights the benefit of using FedReFT in resource-constrained or communication-limited federated learning settings.

### E.1 Intervention Parameter sharing Across Tokens in Same Layer

In this section, we also conducted some additional experiments to show the robustness of FedReFT in different setups. We experiment with whether to share (tie) the intervention parameters  $\phi$  across different input positions within the same layer. Given the positions  $P = \{1, \dots, p\} \cup \{n - s + 1, \dots, n\}$ , we define the untied and tied variants (Wu et al., 2024c):

$$\begin{aligned} \mathbf{I}_{\text{untied}} &= \{\langle \Phi, \{p\}, l \rangle \mid p \in P, l \in L\}, \\ \mathbf{I}_{\text{tied}} &= \{\langle \Phi, P, l \rangle \mid l \in L\}. \end{aligned} \quad (18)$$

while FedReFT (tie  $\phi$ ) offers a compelling trade-off between performance and efficiency. These results highlight the scalability and efficiency of our representation-tuning approach. Appendix Table 12 and 14 depicts these.

### E.2 Additional Experimental Validation

We experiment with different LLaMA model sizes across  $C = 3$  clients following the Distinct Task (Dt) framework for the commonsense reasoning

Table 16: We vary LLaMA model sizes with  $C = 3$  clients following the Distinct Task (DT) design for the commonsense reasoning task, alongside a centralized LoReFT baseline. As model capacity increases, we observe notable performance gains, with the largest model approaching the accuracy of the centralized setting.

Method	BoolQ	PIQA	HellaS.	Avg ↑
Tiny LLaMA 1B	63.92	50.36	47.21	53.83
LLaMA 7B	66.64	78.34	67.92	70.97
LLaMA-2 7B	69.71	75.45	78.26	74.47
LLaMA-3.2 3B	65.93	78.37	82.36	75.55

task in Table 16, along with a centralized LoReFT baseline for comparison. As the model size grows, we observe steady performance improvements, with the largest variant achieving results close to the centralized model. The first four experiments are conducted under the standalone (centralized) setup, and the following four are performed in the federated learning (FL) environment.

Table 17: Trainable Intervention Parameters across Models (in Millions) in FedReFT

Model	Total P(M)	TP(M)	TP(%)↓
LLaMA-1.1B	1100.05	0.72	0.0655
LLaMA 7B	6,738.42	2.10	0.0311
LLaMA-2 7B	6,738.42	2.10	0.0311
LLaMA-3 8B	8,030.27	2.10	0.0261
LLaMA-2-13B	13,015.86	6.55	0.0503
RoBERTa Large	355.36	0.050	0.0138

## F Dataset Description

### F.1 Commonsense Reasoning

We train and evaluate our models on eight commonsense reasoning datasets spanning different types of open-ended QA tasks, following (Hu et al., 2021a), we construct all examples. Table 18 shows the dataset samples.

- **BoolQ** (Clark et al., 2019): A yes/no question answering dataset consisting of naturally occurring questions. We remove the associated passages to ensure a fair comparison.
- **PIQA** (Bisk et al., 2020): A dataset for physical commonsense reasoning. The model must

Table 18: Examples from commonsense reasoning tasks: BoolQ (Clark et al., 2019), PIQA (Bisk et al., 2020), HellaSwag (Zellers et al., 2019), and SIQA (Sap et al., 2019). Each instruction is followed by the selected answer during evaluation.

Dataset	Instruction / Question	Answer
BoolQ	<i>Please answer the following question with true or false:</i> Question: Do Iran and Afghanistan speak the same language?	True
PIQA	<i>Please choose the correct solution to the question:</i> Question: When boiling butter, when it's ready, you can Solution1: Pour it onto a plate Solution2: Pour it into a jar	Solution2
HellaSwag	<i>Please choose the correct ending to complete the given sentence:</i> Removing ice from car: Then, the man writes over the snow covering the window of a car, and a woman wearing winter clothes smiles. then Ending1: , the man adds wax to the windshield and cuts it. Ending2: , a person boards a ski lift... Ending3: , the man puts on a christmas coat... Ending4: , the man continues removing the snow on his car.	Ending4
SIQA	<i>Please choose the correct answer to the question:</i> Cameron decided to have a barbecue and gathered her friends together. How would others feel as a result? Answer1: like attending Answer2: like staying home Answer3: a good friend to have	Answer1

Table 19: Examples from math reasoning tasks: AQuA (Ling et al., 2017), GSM8K (Cobbe et al., 2021), SVAMP (Patel et al., 2021), and MAWPS (Koncel-Kedziorski et al., 2016). Each instruction is followed by the correct answer derived through step-by-step reasoning.

Dataset	Instruction / Question	Answer
AQuA	<i>Solve the following word problem:</i> A car is driven in a straight line toward the base of a vertical tower. It takes 10 minutes for the angle of elevation to change from $45^\circ$ to $60^\circ$ . After how much more time will the car reach the base of the tower? Answer Choices: (A) $5(\sqrt{3} + 1)$ , (B) $6(\sqrt{3} + \sqrt{2})$ , (C) $7(\sqrt{3} - 1)$ , (D) $8(\sqrt{3} - 2)$ , (E) None of these.	(A)
GSM8K	<i>Solve the following question:</i> Janet's ducks lay 16 eggs per day. She eats 3 eggs and uses 4 for baking. She sells the rest at \$2 per egg. How much money does she make daily?	\$18
SVAMP	<i>Solve the following arithmetic question:</i> Each pack of DVDs costs \$76. A discount of \$25 is applied. What is the final price per pack?	\$51
MAWPS	<i>Solve the following math word problem:</i> Tom has 3 boxes of pencils, each containing 12 pencils. He gives 7 pencils to his friend. How many pencils does Tom have left in total?	29

Table 20: Examples from GLUE benchmark (Wang et al., 2018) tasks: MNLI, SST-2, QNLI, and QQP. Each instruction is followed by the corresponding ground truth label.

Dataset	Instruction / Question	Answer
MNLI	Premise: The dog is running through the field. Hypothesis: An animal is moving. Label: entailment	Entailment
SST-2	Sentence: A touching and thought-provoking piece of cinema. Label: positive	Positive
QNLI	Question: What is the capital of France? Sentence: Paris is the capital and most populous city of France. Label: entailment	Entailment
QQP	Question1: How do I learn to play guitar? Question2: What is the best way to learn guitar? Label: duplicate	Duplicate

select the more plausible solution to everyday physical tasks.

- **SIQA** (Sap et al., 2019): Focuses on social interaction reasoning by asking the model to choose responses based on human intent and consequences.
- **HellaSwag** (Zellers et al., 2019): Requires choosing the most coherent sentence completion given a context, often involving physical or temporal common sense.
- **WinoGrande** (Sakaguchi et al., 2021): Inspired by the Winograd Schema Challenge (Levesque et al., 2012), this dataset contains fill-in-the-blank problems with binary choices requiring commonsense coreference reasoning.

We follow the experimental setup in (Hu et al., 2021a) by fine-tuning our models on a combined training corpus referred to as Commonsense170K, which merges all of the above datasets. Evaluation is conducted individually on each dataset’s test split.

## F.2 Arithmetic Reasoning

We evaluate arithmetic reasoning using seven benchmark datasets that cover a range of math word problem types. As in (Hu et al., 2021a), we construct all examples without using golden or retrieved passages. Data samples are shown in Table 19.

- **AQuA** (Ling et al., 2017): Presents algebraic word problems in a multiple-choice format.
- **GSM8K** (Cobbe et al., 2021): A widely used benchmark of grade-school math problems requiring multi-step reasoning.
- **SVAMP** (Patel et al., 2021): A more challenging dataset that tests robustness to paraphrased and structurally altered word problems.
- **MAWPS** (Koncel-Kedziorski et al., 2016): A large repository of math word problems aggregated from multiple sources, covering diverse arithmetic and algebraic reasoning tasks expressed in natural language.

Following (Hu et al., 2021a), we train our models on a combined training set named **MATH10K**.

## F.3 Natural Language Understanding

For NLU, we evaluate on the GLUE benchmark following the evaluation protocol in (Wu et al., 2024c). Data samples are shown in Table 20.

- The validation set is split into two subsets one for in-training evaluation and the other for final testing.
- For large datasets (QQP, MNLI, QNLI), 1,000 samples are used for in-training validation.
- For smaller datasets, half of the validation set is used during training.