

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 MITIGATING NON-IID DRIFT IN ZEROTH-ORDER FED- ERATED LLM FINE-TUNING WITH TRANSFERABLE SPARSITY

Anonymous authors

Paper under double-blind review

## ABSTRACT

Federated Learning enables collaborative fine-tuning of Large Language Models (LLMs) across decentralized Non-Independent and Identically Distributed (Non-IID) clients, but such models’ massive parameter sizes lead to significant memory and communication challenges. This work introduces MEERKAT, a sparse zeroth-order optimization (ZO) method designed for federated LLM fine-tuning. By limiting fine-tuning to a transferable, static, extremely sparse subset of parameters, MEERKAT achieves remarkable communication efficiency, enabling cost-effective high-frequency synchronization. With theoretical analysis and experiments, we show that this high-frequency communication effectively mitigates Non-IID data challenges and leads to superior performance compared to full-parameter ZO. Furthermore, experiment results show that MEERKAT outperforms existing sparsity baselines with better performance at the same communication frequency. To further handle Non-IID drift, MEERKAT leverages traceable local updates and forms a *virtual path* for each client. This virtual path mechanism reveals the GradIP phenomenon: the inner products between LLM pre-training gradients maintained by server and client gradients estimated via ZO converges for extreme Non-IID clients but oscillates for IID ones. This distinct behavior provides a signal for identifying clients with extreme data heterogeneity. Using this signal, MEERKAT-VP is proposed to analyze GradIP trajectories to identify extreme Non-IID clients and applies early stopping to enhance aggregated model quality. Experiments confirm that MEERKAT and MEERKAT-VP significantly improve the efficiency and effectiveness of ZO federated LLM fine-tuning.

## 1 INTRODUCTION

Federated Learning (FL) McMahan et al. (2017) has emerged as a powerful paradigm for enabling decentralized collaboration, particularly relevant for fine-tuning Large Language Models (LLMs) across numerous client devices Dubey et al. (2024); Brown et al. (2020). Unlike centralized training, FL allows clients to train models locally and share only model updates with a central server. However, fine-tuning LLMs in a FL setting faces two major challenges: the massive model parameter size and the Non-Independent and Identically Distributed (Non-IID) data distribution across clients. The former leads to high computation demands on clients and significant communication overhead, while the latter causes client drift and hinder global convergence. These challenges make LLM fine-tuning impractical on resource-constrained clients and hinder the effective use of decentralized data.

Zerth-order Optimization (ZO) provides a promising avenue for addressing some of these challenges in federated LLM fine-tuning. By estimating gradients through model perturbations and forward passes, ZO bypasses the need for backpropagation and the storage of intermediate activations, leading to more memory-efficient learning on client devices Zhang et al. (2021); Fang et al. (2022); Ling et al. (2024); Liu et al. (2024); Malladi et al. (2023). However, applying standard ZO directly to the massive parameter space of LLMs can still be computationally inefficient and the optimization process unstable Malladi et al. (2023). Moreover, adapting ZO for federated LLM fine-tuning remains challenging, particularly in balancing computational efficiency, communication overhead, and model performance under Non-IID data heterogeneity.

In order to address the above challenges, we propose MEERKAT, a sparse ZO method designed for efficient federated LLM fine-tuning. MEERKAT addresses the computational and communication burdens by focusing ZO updates on a static, extremely sparse (less than 0.1%), and transferable subset of LLM parameters. This subset is strategically identified using gradients derived from pre-training data, ensuring that updates target parameters most sensitive to the loss function. This selective approach dramatically reduces communication overhead and supports cost-effective high-frequency synchronization. As we will demonstrate through theoretical analysis and extensive experiments, the combination of high communication frequency and sparsity in MEERKAT enables frequent yet lightweight synchronization. This effectively reduces the convergence error floor in theory and practice, leading to consistently superior performance compared to full-parameter ZO fine-tuning and other sparsity methods under the same communication frequency.

Leveraging MEERKAT’s efficient high-frequency synchronization to effectively mitigate Non-IID data challenges, we further enhance its adaptability to weak network conditions. By employing a virtual path mechanism to track client updates, we enable the server to analyze client training dynamics without accessing raw data, thus facilitating robust operation even when frequent direct communication is constrained. Within this virtual path, we observe the **GradIP phenomenon**, a pattern revealed by the GradIP score, which computes the inner product between local client gradients estimated via ZO and server pre-training gradients. GradIP scores converge for Non-IID clients while oscillating for IID clients, serving as a clear indicator of data heterogeneity. Leveraging this insight, we propose MEERKAT-VP that introduces a virtual path client selection method to identify clients with significant Non-IID characteristics and apply early stopping, thereby reducing their adverse impact on the aggregated model and enhancing its quality.

In summary, this paper makes the following contributions:

- **Performance Improvement with Sparsity.** Meerkat consistently outperforms full-parameter ZO optimization in both IID and Non-IID settings, demonstrating the effectiveness of our sparse update strategy. Extensive experiments show that Meerkat surpasses not only full-parameter ZO but also other sparse methods, such as LoRA and weight-magnitude, achieving superior performance.
- **High Frequency Communication with Sparsity Can Lower the Error Floor.** MEERKAT leverages extreme model sparsity to reduce local computational memory. Exchanging scalar gradients drastically decreases communication costs, enabling high-frequency communication.
- **Traceable Local Updates and GradIP Phenomenon:** MEERKAT leverages traceable sparse local updates and forms a *virtual path*. The virtual paths reveals the GradIP phenomenon: the inner product between LLM pre-training gradients maintained by server and client gradients estimated via ZO converges for extreme Non-IID clients but oscillates for IID ones. This distinct behavior serves as a signal for detecting clients with extreme data heterogeneity.
- **MEERKAT-VP: Early Stopping for Extreme Non-IID Clients.** Leveraging the GradIP phenomenon via virtual path client selection, MEERKAT-VP effectively manages extreme Non-IID clients, by early stopping these clients to improve global model quality.
- **Theoretical and Experimental Validation.** We present theoretical analysis and extensive experiments across diverse FL settings, validating the scalability and performance benefits of both MEERKAT and MEERKAT-VP.

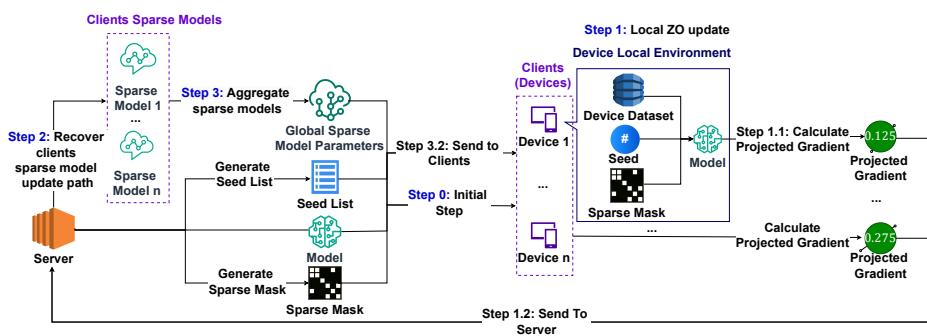


Figure 1: MEERKAT: Sparse zeroth-order optimization for federated LLM fine-tuning workflow.

108 

## 2 SPARSE ZEROTH-ORDER OPTIMIZATION FOR FEDERATED LLM 109 FINE-TUNING 110

111 This section introduces MEERKAT, a sparse ZO method for federated LLM fine-tuning, and its  
112 upgraded version, MEERKAT-VP, which incorporates Virtual Path Client Selection(VPCS) strategy.  
113 This strategy leverages the traceable virtual path of client local updates to identify clients with  
114 extremely Non-IID data and applies early stopping to mitigate their adverse impact on global model  
115 convergence. We first introduce the technical details of MEERKAT, as illustrated in Figure 1, and  
116 subsequently describe MEERKAT-VP, shown in Figure 5. We then present theoretical convergence  
117 analysis for both methods and discuss their strengths in terms of cost-effectiveness, traceability, and  
118 the use of early stopping to mitigate client drift caused by Non-IID data.  
119

120 

### 2.1 MEERKAT: EXTREME SPARSE ZEROTH-ORDER FEDERATED LLM FINE-TUNING

121

122 **Sparse ZO On-Device LLM Fine-Tuning.** MEERKAT performs sparse ZO for LLM fine-tuning  
123 on the client device. Let  $\mathcal{D}$  denote the client dataset we would like an LLM to fine-tune with loss  
124 function  $f$ . Given the LLM weight  $\mathbf{w} \in \mathbb{R}^d$ , we perform an iterative optimization by randomly  
125 sampling a batch  $\mathcal{B} \subset \mathcal{D}$  for each step and performing the local update step as  
126

$$g = \frac{f(\mathbf{w} + \epsilon(\mathbf{z} \odot \mathbf{m}); \mathcal{B}) - f(\mathbf{w} - \epsilon(\mathbf{z} \odot \mathbf{m}); \mathcal{B})}{2\epsilon}, \quad \hat{\nabla} f = g(\mathbf{z} \odot \mathbf{m}). \quad (1)$$
127

128 where  $\mathbf{z} \in \mathbb{R}^d$  is a random vector sampled from a Gaussian distribution  $\mathcal{N}(0, I_d)$ ,  $\epsilon \in \mathbb{R}$  is the  
129 perturbation magnitude, and  $\mathbf{m} \in \{0, 1\}^d$  is a binary sparse mask with density ratio  $u$  that selects a  
130 subset of parameters for updates.  
131

132 **Extremely Sparse Parameters Obtained from Pre-Training.** According to the formulation in  
133 Eq equation 1, we focus the perturbation of the LLM on a subset of parameters determined by a  
134 binary mask  $\mathbf{m}$ . The mask  $\mathbf{m}$  is derived from the pre-training process of the LLM. We compute the  
135 average squared gradients of each parameter over a subset of the C4 dataset Raffel et al. (2020). Then,  
136 we select the top  $u$  parameters with the highest average squared gradient values and mark them as 1  
137 in  $\mathbf{m}$ . In practice, we set  $u$  to 0.1%, resulting in extremely sparse updates.  
138

139 **FL with MEERKAT.** The workflow of MEERKAT is illustrated in Figure 1 and Algorithm 2.  
140 MEERKAT first loads each client with the pre-trained weight  $\mathbf{w}_0$  and the sparse mask  $\mathbf{m}$ . Next,  
141 MEERKAT initializes a random seed list  $\{s_1^1, \dots, s_1^T\}$  at the server to generate the random Gaussian  
142 vector  $\mathbf{z}$  for each local step in the first round. Next, MEERKAT performs an iterative federated  
143 optimization with  $R$  rounds of client-server synchronization with each round as follows.  
144

145 (1) *Local ZO update at each client.* Upon receiving global model weights  $\mathbf{w}_{r-1}$  and seed list  
146  $\{s_r^1, \dots, s_r^T\}$  from the server, each client performs  $T$  local iteration steps. In each local step  $t$ , the  
147 client perturbs the model parameters selected by  $\mathbf{m}$  with the random vector  $\mathbf{z}_k^t$  generated by the  
148 random seed  $s_r^t$ . Each client then computes projected gradient  $g_k^t$  (a scalar) according to Eq. equation 1.  
149 Using  $g_k^t$ , each client calculates the local gradient  $\hat{\nabla} f_k^t$  and updates the local model  $w_k$  with learning  
150 rate  $\eta$ . After  $T$  local steps, each client uploads a list of projected gradients  $\{g_k^1, g_k^2, \dots, g_k^T\}$  to the  
151 server. (2) *Server reconstructs client update with virtual path.* Since the server shares the same  
152 random seed list with clients for the round, it can reconstruct each client’s local model update path  
153 upon receiving their projected gradients. We term this server-side reconstruction process the *virtual  
154 path*, as it allows the server to follow the client’s local steps without accessing raw data. As shown in  
155 Step 2 of Algorithm 2, the server uses the preserved random seed and receives project gradients of  
156 each local step from each client to recover the local model update path for each client. (3) *Server  
157 aggregates and initiate the next round:* After virtual path reconstruction, the server aggregates the  
158 reconstructed client model weights  $\mathbf{w}_k^T$  to sparsely update the global model to  $\mathbf{w}_r$ . Subsequently, the  
159 server sends  $\mathbf{w}_r$  and a new seed list  $\{s_{r+1}^1, \dots, s_{r+1}^T\}$  to clients and initializes next round.  
160

161 **MEERKAT-VP: Virtual Path Client Selection and Early Stopping.** MEERKAT-VP extends  
162 MEERKAT by incorporating a VPCS strategy designed for heterogeneous environments. Lever-  
163 aging the virtual path reconstruction capability, the server analyzes client update trajectories to  
164 identify those with extremely Non-IID data distributions. MEERKAT-VP then applies an early stop-  
165 ping mechanism to these identified clients, restricting them to a single local step to mitigate the  
166 negative impact of their skewed updates on global model convergence and performance.  
167

162 2.2 THEORETICAL CONVERGENCE ANALYSIS  
163

164 We theoretically analyze the convergence of MEERKAT and MEERKAT-VP under the  
165 Polyak–Łojasiewicz (PL)-type non-convex condition. All technical assumptions and the corre-  
166 sponding proof are presented in Appendix C.

167 **Theorem 2.1** (Convergence rate of MEERKAT). *Under Assumptions C.1–C.6, if the learning rate  
168 satisfy  $\eta = \min \left\{ \frac{1}{L(u+2)}, \frac{\mu \sqrt{c} (1+\sqrt{c_h})}{2 L^2 (2+u)^2} \right\}$ , then the global model  $\{\mathbf{w}^r\}$  generated by the MEERKAT  
169 algorithm satisfies the following convergence bound:*

$$172 \quad \frac{1}{R} \sum_{r=0}^{R-1} (f(\mathbf{w}^r) - f^*) \leq \mathcal{O} \left( \frac{(2+u)^2}{TR} \cdot \mathbb{E}[f(\mathbf{w}^0) - f(\mathbf{w}^R)] \right) + \mathcal{O} \left( \frac{T}{2+u} \right) + O(1) \quad (2)$$

173 **Theorem 2.2** (Convergence rate of MEERKAT-VP). *Under Assumptions C.1–C.6, if the learning rate  
174 satisfies  $\eta = \min \left\{ \frac{1}{L(u+2)}, \frac{\mu \sqrt{c} (K_g T + K_b)}{2 K (2+u)^2 L^2 T \gamma} \right\}$  and each client  $k \in K_b$  performs  $T = 1$  local step  
175 while the remaining  $K_g$  clients perform  $T$  local steps, then the global model  $\{\mathbf{w}^r\}$  generated by the  
176 MEERKAT-VP algorithm satisfies the following convergence bound:*

$$177 \quad \frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E}_{\tilde{z}} [f(\mathbf{w}^r) - f^*] \leq \mathcal{O} \left( \frac{(K_g + K_b)^2 (2+u)^2 \gamma T}{c (K_g T + K_b)^2 R} \right) + \mathcal{O} \left( \frac{1+u}{K_g + K_b} \sum_{k=1}^{K_g + K_b} \Delta_k \right) \\ 178 \quad + \mathcal{O} \left( \frac{c T K_g}{(K_g + K_b)(1+u) \gamma} \right) + \mathcal{O} \left( \frac{c K_b \sigma_h^2}{(K_g + K_b)(1+u) T \gamma} \right) + O(1). \quad (3)$$

179 The detailed theoretical analysis and proofs for Theorem 2.1 (MEERKAT) can be found in Ap-  
180 pendix C.4, and for Theorem 2.2 (MEERKAT-VP) in Appendix C.5.

181 **Insights of MEERKAT.** MEERKAT’s convergence reveals the intricate interplay of local steps  $T$   
182 and density  $u$  on performance. (1) MEERKAT’s sparsity can theoretically improve performance.  
183 Lower  $u$  (higher sparsity) quadratically benefits the rate-dependent term ( $\propto (2+u)^2$ ), favoring  
184 faster initial convergence. However, it inflates the steady-state error ( $\propto \frac{1}{2+u}$ ). Comparing to the  
185 full-parameter case ( $u = 1$ ), sparsity ( $u < 1$ ) can reduce the overall bound by decreasing the  
186 rate-dependent term, offering communication and computational benefits. Yet, excessive sparsity can  
187 increase the steady-state error, suggesting an optimal density level  $u \in (0, 1]$ . (2) High frequency  
188 communication with sparsity can lower the error floor. Increasing  $T$  improves the transient term  
189 scaling with  $\mathcal{O}(\frac{(2+u)^2}{TR})$ , potentially accelerating convergence towards the steady state; however, it  
190 expands the steady-state term  $\mathcal{O}(\frac{T}{2+u})$ , thereby increasing the error floor. Conversely, decreasing  $T$   
191 reduces the steady-state term, leading to a tighter final accuracy. Although smaller  $T$  can lead to  
192 larger rate-dependent term. Its impact diminishes as the number of rounds  $R$  increases. This analysis  
193 suggests that operating with frequent communication can theoretically reduce the steady-state error.

194 **Advantages of MEERKAT-VP.** We compare each component of the error bound under the same  
195  $T$  and  $R$ . First, the transient term ratio between MEERKAT-VP and MEERKAT is approximately  
196  $\gamma(1+\sqrt{c_h})^2 < 1$ , and as  $c_h \rightarrow 1$  so  $\gamma \rightarrow 0$ , the product  $\gamma(1+\sqrt{c_h})^2 \rightarrow 0$ , causing the transient error  
197 to vanish. Second, the noise term ratio is given by  $\frac{\sigma_h^2/2}{\sigma_h^2/(\mu(1+\sqrt{c_h})^2)} = \frac{\mu(1+\sqrt{c_h})^2}{2}$ , which remains  
198 below 1 whenever  $\mu(1+\sqrt{c_h})^2 < 2$ . Since  $\mu < 1$  empirically, this condition typically holds.  
199 Moreover, MEERKAT-VP introduces an additional variance term  $\frac{c K_b \sigma_h^2}{2 K (2+u) L T \gamma}$  that decays as  $\mathcal{O}(1/T)$ ,  
200 making it negligible for large local steps. Lastly, in terms of heterogeneity, the coefficient of the  
201 heterogeneity term  $\sum_k \Delta_k$  in MEERKAT-VP is smaller:  $\frac{(2+u)L}{4K} < \frac{L}{K}$ , and the extra variance term  
202 scales inversely with  $K$ , thus diminishing in larger systems. Therefore,  $E_{\text{MEERKAT-VP}} < E_{\text{MEERKAT}}$   
203 and this gap widens as data heterogeneity  $c_h$  increases. The detailed mathematical derivations and  
204 analysis, please refer to the Appendix C.5.

216 2.3 CLAIM 1: MEERKAT CAN OUTPERFORMS FULL-PARAMETER FEDERATED ZO UNDER  
217 SAME SYNCHRONIZATION FREQUENCY  
218219 We claim that with fixed and extreme sparsity, MEERKAT outperforms full-parameter ZO in federated  
220 LLM fine-tuning under the same synchronization frequency and effectively mitigates the Non-IID  
221 client data problem through frequent synchronization and sparsity.222 **Advantages of Sparsity in Federated ZO.** ZO has an intrinsic need for sparsity due to its reliance  
223 on nearly uniform perturbations across dimensions. Research on ZO shows that selecting sensitive  
224 parameters using gradient-based methods consistently outperforms alternative strategies such as  
225 weight magnitude or random parameter selection Guo et al. (2024). Following this idea, MEERKAT  
226 produces LLM-sensitive parameters with gradient-based sparsification on pre-training data such as  
227 C4 Raffel et al. (2020). Moreover, MEERKAT fine-tunes LLMs by estimating gradients through  
228 forward passes, completely bypassing backpropagation. This approach minimizes the need to cache  
229 gradients and activations, leading to significant memory savings. Focusing on sensitive parameters,  
230 MEERKAT ensures efficient and effective fine-tuning even under extreme sparsity levels (e.g., updating  
231 only 0.1% of the parameters). Furthermore, these sensitive parameters exhibit transferability across  
232 downstream tasks. Theoretical analysis (Appendix C.4) also confirms that lower density  $u$  leads to  
233 faster convergence via improved rate-dependent terms  $\mathcal{O}((2+u)^2/(TR))$ , while excessive sparsity  
234 increases the steady-state error  $\mathcal{O}(T/(2+u))$ , suggesting an optimal sparsity trade-off.235 **Performance Under High Synchronization Frequency.** The lightweight communication of  
236 MEERKAT enables frequent client-server synchronization at a low cost, which is crucial for ad-  
237 dressing data heterogeneity Yang et al. (2024); Mendieta et al. (2022) in FL. In high-frequency  
238 communication scenarios, both the clients and the server only exchange a list of scalars (projected  
239 gradients) whereas in lower-frequency synchronization, clients have to upload projected gradients  
240 but still download sparse model parameters. By eliminating the need to download sparse model  
241 parameters in high-frequency synchronization, this approach is significantly more bandwidth-efficient,  
242 further minimizing communication overhead. We present the high-frequency synchronization algo-  
243 rithm of MEERKAT in Appendix C Algorithm 3. By facilitating frequent synchronization, training  
244 can better prevent clients from drifting. Our previous theoretical analysis also demonstrates that a  
245 smaller  $T$  might influence the rate-dependent term, its beneficial impact on reducing the steady-state  
246 error is significant for achieving a tighter final accuracy over many rounds  $R$ .

## 247 2.4 CLAIM 2: EMPIRICAL GRADIP PHENOMENON REVEALS DATA HETEROGENEITY

248 MEERKAT’s traceable virtual path allows us to analyze client local training dynamics, revealing an  
249 empirical phenomenon related to data heterogeneity via a metric we call GradIP.  
250251 **Definition 2.3.** Gradient Inner Product (GradIP) score: Let  $\hat{\nabla}f_k^t$  (see Algorithm 2) denote the ZO  
252 gradient of LLM with Eq equation 1 on client  $k$  at local step  $t$ . Let  $\nabla f_p$  denote the gradient of LLM  
253 computed by backpropagation on pre-training data. We define the GradIP score as  $\langle \nabla f_p, \hat{\nabla}f_k^t \rangle$ .254 **GradIP As Indicator for Data Heterogeneity.** Leveraging the virtual path reconstruction capability  
255 of MEERKAT, the server can trace each client’s local training trajectory. This process uses the  
256 uploaded projected gradients  $g_k^t$  along with the shared random seeds (which regenerate  $\mathbf{z}_k^t$ ) and the  
257 sparse mask  $\mathbf{m}$  to reconstruct the local gradient  $\hat{\nabla}f_k^t$ . To understand the impact of a client’s local  
258 data distribution on its training process, we introduce the *GradIP* metric. Inspired by the use of  
259 pre-training data gradients to identify sensitive parameters, GradIP quantifies the cosine similarity  
260 between the local gradient computed during client training and the LLM pre-training gradient.261 **Empirical GradIP Phenomenon.** Through the traceable virtual path provided by MEERKAT, we  
262 empirically investigated the behavior of the GradIP score among clients with different data distri-  
263 butions (IID and Non-IID) over their local training steps. Our analysis, presented in Appendix C.6,  
264 demonstrates distinct patterns in the dynamics of gradient norms based on data heterogeneity. While  
265 IID client gradient norms exhibit fluctuations, those of extremely Non-IID clients decay and converge  
266 towards zero. The GradIP definition depends on the fixed pre-training gradient norm, local client gra-  
267 dient norm, and the angle  $\theta$  between them. We hypothesize that  $\theta$  between these two gradient vectors  
268 is nearly orthogonal. This leads us to expect a different manifestation of the GradIP Phenomenon  
269 when comparing IID and extremely Non-IID clients, primarily influenced by their differing local  
gradient norm trajectories.

## 2.5 CLAIM 3: VIRTUAL PATH CLIENT SELECTION VIA GRADIP ANALYSIS

Building upon the traceable virtual path capability introduced in MEERKAT, we claim that VPCS, by leveraging GradIP analysis, effectively identifies and manages clients with extremely Non-IID data distribution, thereby improving global model performance and convergence. As established in Section 2.4, the GradIP score, computable by the server through virtual path reconstruction, provides an effective signal to identify such clients. VPCS utilizes this GradIP signal to detect extremely Non-IID clients. By analyzing the GradIP score trajectory and its behavior over local steps during a calibration phase, using metrics defined in Appendix table 3, the server empirically identifies clients exhibiting the characteristic diminishing GradIP behavior associated with extremely Non-IID data distribution. Upon identification via GradIP analysis, VPCS applies early stopping: these clients perform only one local training step per communication round. To ensure full data utilization over training, a data pointer tracks the batch processed, allowing clients to resume from that point in subsequent rounds. This strategy mitigates client drift from skewed data while ensuring their entire dataset is eventually processed. Algorithm 1 outlines the detailed procedure, and Figure 5 illustrates the workflow. Our previous theoretical analysis of MEERKAT-VP suggests that early stopping on extremely Non-IID clients can lead to improved global model performance.

### 3 EXPERIMENT

In this section, we aim to validate the effectiveness of MEERKAT and MEERKAT-VP. We aim to address the following research questions in response to claims in Section 2: (1) **RQ 1 for Claim 1 (2.3):** Is MEERKAT more effective than full parameter federated ZO under the same synchronization frequency, especially in heterogeneous environments? (2) **RQ 2 for Claim 2 (2.4):** Can the empirical GradIP phenomenon, observed via the virtual path, effectively reveal data heterogeneity by showing distinct behaviors for IID and Non-IID data distribution clients? (3) **RQ 3 for Claim 3 (2.5):** Can MEERKAT-VP, leveraging GradIP analysis, mitigate the impact of extreme Non-IID data compared to MEERKAT?

We focus on models Gemma-2-2b Team (2024), Qwen2-1.5B qwe (2024), Llama-3.2-1B Dubey et al. (2024). We conduct experiments on SST2 Socher et al. (2013), AG’s News Zhang et al. (2015), Yelp polarity (yelp) Zhang et al. (2015), RTE Wang (2018), BoolQ Clark et al. (2019), WSC Levesque et al. (2012), WiC Pilehvar & Camacho-Collados (2018) datasets. The datasets are partitioned across clients following a Dirichlet distribution to simulate clients with Non-IID data. For more experimental settings, we refer the readers to Appendix D.1.

### 3.1 ANSWER TO RQ1: SUPERIORITY OF MEERKAT COMPARED TO FULL-FEDZO IN FL

This section experimentally validates Claim 1 (Section 2.3), demonstrating MEERKAT’s superiority over full-parameter Federated ZO under the same synchronization frequency and its effectiveness in mitigating Non-IID challenges via high-frequency synchronization.

**Algorithm 1** MEERKAT-VP

- Input:** calibration step  $T_{\text{cali}}$ , pre-training gradients  $\nabla f_{c4}$ , projected gradients  $\{g_k^1, \dots, g_k^{T_{\text{cali}}}\}$ , seed  $s_r^t$ , sparse mask  $\mathbf{m}$ , initial phase steps  $T_{\text{init}}$ , later phase steps  $T_{\text{later}}$ , convergence threshold  $\sigma$ , Initial to later ratio  $\rho_{\text{later}}$ , quiescent step ratio  $\rho_{\text{quie}}$
- Step 1: Virtual Path Reconstruction & GradIP Calculation**
  - Generate  $\mathbf{z}_k^t$  using  $s_r^t$ .
  - Compute  $\hat{\nabla} f_k^t = g_k^t \cdot (\mathbf{z}_k^t \odot \mathbf{m})$
  - Compute  $\text{Gradip} = \hat{\nabla} f_k^t \cdot \nabla f_{c4}$  (Definition 2.3).
- Step 2: Identify Extremely Non-IID Clients**
  - Compute the average value of Gradip over the initial-phase steps.
$$\text{Gradip}_{\text{init\_avg}} = \frac{1}{T_{\text{init}}} \sum_{t=1}^{T_{\text{init}}} \text{Gradip}_t$$
  - Compute the average value of Gradip over the later-phase steps.
$$\text{Gradip}_{\text{later\_avg}} = \frac{1}{T_{\text{later}}} \sum_{t=1}^{T_{\text{later}}} \text{Gradip}_t$$
  - Compute the client's Initial to later ratio  $\rho_{\text{later\_client}}$  and quiescent step ratio  $\rho_{\text{quie\_client}}$
$$\rho_{\text{quie\_client}} = \frac{\{s \in \{1, 2, \dots, T_{\text{later}}\} \mid \text{Gradip}_s < \sigma\}}{T_{\text{later}}}$$

$$\rho_{\text{later\_client}} = \frac{\text{Gradip}_{\text{init\_avg}}}{\text{Gradip}_{\text{later\_avg}}}$$
- Record client IDs whose  $\rho_{\text{later\_client}}$  or  $\rho_{\text{quie\_client}}$  exceed  $\rho_{\text{later}}$  or  $\rho_{\text{quie}}$ .
- Step 3: Early Stopping**
- Require these identified clients to only perform one local training step.

324  
 325      Table 1: Performance comparison of MEERKAT and Full-FedZO on multiple non-IID data  
 326      distribution settings. “Acc” is the average test accuracy across tasks. Bold numbers indicate the  
 327      highest value in each row.

|                         | Methods          | Local Step | SST-2        | AgNews       | Yelp         | BoolQ        | RTE          | WSC          | WIC          | Acc          |
|-------------------------|------------------|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| 329 <b>LLaMA-3.2-1B</b> | Full-FedZO       | 10         | 0.909        | 0.705        | 0.940        | 0.641        | 0.542        | 0.634        | 0.523        | 0.699        |
|                         | Weight Magnitude | 10         | 0.902        | 0.857        | 0.951        | 0.696        | 0.551        | 0.519        | 0.546        | 0.717        |
|                         | Lora-FedZO       | 10         | 0.901        | 0.749        | 0.96         | 0.649        | 0.524        | 0.634        | 0.59         | 0.715        |
|                         | MEERKAT          | 10         | <b>0.916</b> | <b>0.872</b> | <b>0.964</b> | 0.695        | <b>0.600</b> | <b>0.653</b> | <b>0.614</b> | <b>0.759</b> |
|                         | Full-FedZO       | 30         | 0.904        | 0.706        | 0.935        | 0.636        | 0.533        | 0.634        | 0.539        | 0.698        |
|                         | Weight Magnitude | 30         | 0.902        | 0.84         | 0.946        | 0.674        | 0.542        | 0.556        | 0.550        | 0.716        |
|                         | Lora-FedZO       | 30         | 0.904        | 0.556        | 0.964        | 0.652        | 0.533        | 0.634        | 0.545        | 0.684        |
|                         | MEERKAT          | 30         | 0.897        | <b>0.862</b> | <b>0.965</b> | 0.646        | <b>0.577</b> | <b>0.644</b> | <b>0.583</b> | <b>0.739</b> |
|                         | Full-FedZO       | 50         | 0.889        | 0.696        | 0.935        | 0.633        | 0.542        | 0.634        | 0.529        | 0.694        |
|                         | Weight Magnitude | 50         | 0.897        | 0.838        | 0.948        | 0.662        | 0.551        | 0.562        | 0.554        | 0.716        |
| 335 <b>Qwen2-1.5b</b>   | Lora-FedZO       | 50         | 0.876        | 0.447        | 0.967        | 0.639        | 0.541        | 0.634        | 0.562        | 0.667        |
|                         | MEERKAT          | 50         | <b>0.909</b> | 0.827        | <b>0.965</b> | 0.647        | <b>0.595</b> | 0.634        | <b>0.567</b> | <b>0.734</b> |
|                         | Full-FedZO       | 100        | 0.901        | 0.705        | 0.939        | 0.632        | 0.533        | 0.634        | 0.525        | 0.695        |
|                         | Weight Magnitude | 100        | 0.885        | 0.83         | 0.946        | 0.66         | 0.56         | 0.534        | 0.548        | 0.709        |
|                         | Lora-FedZO       | 100        | 0.868        | 0.247        | 0.953        | 0.642        | 0.521        | 0.634        | 0.529        | 0.628        |
|                         | MEERKAT          | 100        | 0.896        | 0.777        | <b>0.961</b> | 0.658        | <b>0.577</b> | <b>0.644</b> | <b>0.573</b> | <b>0.726</b> |
|                         | Full-FedZO       | 10         | 0.888        | 0.700        | 0.928        | 0.694        | 0.808        | 0.673        | 0.639        | 0.761        |
|                         | Weight Magnitude | 10         | 0.881        | 0.84         | 0.939        | 0.681        | 0.795        | 0.672        | 0.623        | 0.776        |
|                         | Lora-FedZO       | 10         | 0.939        | 0.847        | 0.944        | 0.667        | 0.795        | 0.663        | 0.521        | 0.768        |
|                         | MEERKAT          | 10         | <b>0.949</b> | <b>0.881</b> | 0.934        | <b>0.752</b> | <b>0.813</b> | <b>0.682</b> | 0.628        | <b>0.805</b> |
| 343 <b>Gemma2-2b</b>    | Full-FedZO       | 30         | 0.892        | 0.699        | 0.926        | 0.708        | 0.791        | 0.663        | 0.594        | 0.753        |
|                         | Weight Magnitude | 30         | 0.88         | 0.843        | 0.939        | 0.681        | 0.786        | 0.673        | 0.594        | 0.771        |
|                         | Lora-FedZO       | 30         | 0.923        | 0.843        | 0.948        | 0.666        | 0.777        | 0.673        | 0.519        | 0.764        |
|                         | MEERKAT          | 30         | <b>0.944</b> | <b>0.878</b> | 0.928        | <b>0.734</b> | <b>0.800</b> | 0.663        | <b>0.624</b> | <b>0.795</b> |
|                         | Full-FedZO       | 50         | 0.868        | 0.696        | 0.922        | 0.707        | 0.773        | 0.663        | 0.594        | 0.746        |
|                         | Weight Magnitude | 50         | 0.883        | 0.855        | 0.938        | 0.703        | 0.768        | 0.673        | 0.595        | 0.774        |
|                         | Lora-FedZO       | 50         | 0.934        | 0.834        | 0.941        | 0.679        | 0.76         | 0.653        | 0.510        | 0.759        |
|                         | MEERKAT          | 50         | <b>0.948</b> | <b>0.872</b> | 0.926        | <b>0.746</b> | <b>0.795</b> | 0.663        | 0.594        | <b>0.792</b> |
|                         | Full-FedZO       | 100        | 0.864        | 0.691        | 0.917        | 0.675        | 0.777        | 0.653        | <b>0.620</b> | 0.742        |
|                         | Weight Magnitude | 100        | 0.888        | 0.842        | 0.934        | 0.695        | 0.768        | 0.656        | 0.579        | 0.766        |
| 352 <b>Qwen2-1.5b</b>   | Lora-FedZO       | 100        | 0.934        | 0.785        | 0.937        | 0.664        | 0.786        | 0.653        | 0.512        | 0.753        |
|                         | MEERKAT          | 100        | <b>0.936</b> | <b>0.878</b> | 0.925        | <b>0.741</b> | <b>0.795</b> | <b>0.663</b> | 0.610        | <b>0.792</b> |
|                         | Full-FedZO       | 10         | 0.928        | 0.721        | 0.943        | 0.731        | 0.564        | 0.644        | 0.595        | 0.732        |
|                         | Weight Magnitude | 10         | 0.931        | 0.849        | 0.955        | 0.778        | 0.711        | 0.634        | 0.595        | 0.779        |
|                         | Lora-FedZO       | 10         | 0.936        | 0.853        | 0.966        | 0.763        | 0.568        | 0.663        | 0.605        | 0.765        |
|                         | MEERKAT          | 10         | <b>0.939</b> | <b>0.869</b> | 0.96         | <b>0.804</b> | 0.591        | 0.634        | <b>0.609</b> | 0.772        |
|                         | Full-FedZO       | 30         | 0.927        | 0.802        | 0.932        | 0.725        | 0.568        | 0.634        | 0.581        | 0.738        |
|                         | Weight Magnitude | 30         | 0.935        | 0.851        | 0.951        | 0.771        | 0.653        | 0.634        | 0.598        | 0.770        |
|                         | Lora-FedZO       | 30         | 0.932        | 0.804        | 0.966        | 0.671        | 0.551        | 0.634        | 0.589        | 0.735        |
|                         | MEERKAT          | 30         | <b>0.94</b>  | <b>0.855</b> | 0.947        | 0.734        | 0.568        | <b>0.644</b> | <b>0.601</b> | 0.756        |
| 357 <b>Gemma2-2b</b>    | Full-FedZO       | 50         | 0.932        | 0.791        | 0.943        | 0.712        | 0.582        | <b>0.634</b> | 0.567        | 0.737        |
|                         | Weight Magnitude | 50         | 0.936        | 0.851        | 0.941        | 0.745        | 0.591        | 0.628        | 0.597        | 0.756        |
|                         | Lora-FedZO       | 50         | 0.91         | 0.779        | 0.942        | 0.664        | 0.557        | <b>0.634</b> | 0.597        | 0.726        |
|                         | MEERKAT          | 50         | <b>0.945</b> | <b>0.857</b> | <b>0.966</b> | <b>0.767</b> | <b>0.613</b> | <b>0.634</b> | <b>0.623</b> | <b>0.772</b> |
|                         | Full-FedZO       | 100        | 0.925        | 0.818        | 0.933        | 0.672        | 0.533        | 0.615        | 0.567        | 0.723        |
|                         | Weight Magnitude | 100        | 0.922        | 0.839        | 0.942        | 0.723        | 0.568        | 0.644        | 0.592        | 0.747        |
|                         | Lora-FedZO       | 100        | 0.922        | 0.247        | 0.942        | 0.62         | 0.541        | 0.634        | 0.573        | 0.640        |
|                         | MEERKAT          | 100        | <b>0.94</b>  | <b>0.851</b> | <b>0.951</b> | <b>0.745</b> | 0.551        | 0.634        | 0.574        | <b>0.749</b> |
|                         | Full-FedZO       | 100        | 0.925        | 0.818        | 0.933        | 0.672        | 0.533        | 0.615        | 0.567        | 0.723        |
|                         | Weight Magnitude | 100        | 0.922        | 0.839        | 0.942        | 0.723        | 0.568        | 0.644        | 0.592        | 0.747        |

364      First, to assess sparsity’s benefits, we compare MEERKAT to Full-FedZO and other sparse methods  
 365      (Weight Magnitude, LoRA-FedZO, Random-Select) with equivalent synchronization frequencies  
 366      (local steps  $T \in \{10, 30, 50, 100\}$ ). With a fixed 0.1% mask, MEERKAT reduces communication  
 367      budget by over 1000 $\times$  compared to Full-FedZO and achieves a strong computational and communica-  
 368      tion efficiency (Table 24). Using C4 as a calibration dataset, our analysis shows that the sensitivity  
 369      of the gradient is highly concentrated: the top 0.1% of the parameters have 52 $\times$  larger average  
 370      square gradients than the next 0.1–1% bucket (Table 9), which motivates extreme sparsity. The  
 371      mask is transferred across domain-shifted calibration datasets, and a client-aggregated UnionMask  
 372      performs comparably (Table 11). Across IID and Non-IID data distributions, MEERKAT outperforms  
 373      Full-FedZO and other sparsity methods on many tasks (Tables 1, 10, 13). Under the same settings,  
 374      MEERKAT also outperforms DeComFL Li et al. (2024) (Table 19).

375      Next, we evaluate performance under an extreme communication regime with a single local step  
 376      ( $T=1$ ). We compare MEERKAT with Full-FedZO and LoRA-FedZO in the IID and Non-IID data  
 377      distributions (Dirichlet  $\alpha \in \{0.5, 0.3, 0.1\}$ ). Figure 2 presents the results for  $\alpha = 0.5$ , the results for  
 378       $\alpha = 0.3$  and 0.1 are available in Appendix D.2 figure 6. Specifically, Figure 2 reveals a remarkable

finding: on the Qwen2-1.5b model, MEERKAT’s average test accuracy over seven tasks under Non-IID data distribution matches that under IID data distribution. Beyond this exact match, results show that at a local step of  $T = 1$ , MEERKAT effectively bridges the performance gap between IID and Non-IID data distribution settings, achieving nearly comparable test accuracy across both data distributions, and consistently outperforms baselines. Varying sparsity under  $T = 1$  (Table 15) confirms strong accuracy even at  $10^{-3}$ – $10^{-4}$ , substantially reducing client memory demands and making it ideal for resource-constrained FL. These results support **Claim 1**: high-frequency communication combined with extreme sparsity mitigates Non-IID drift. We also explored sensitive parameter selection using downstream task data. Since performance remained comparable under identical communication frequencies and sparsity levels, we prioritized pre-training data to better preserve client privacy (Appendix D.2, Tables 21, 20, 22).



Figure 2: This figure compares three methods—Full-FedZO, LoRA-FedZO, and MEERKAT—on three LLMs: LLaMA-3.2-1B, Qwen2-1.5b, and Gemma2-2b. The x-axis shows the different methods, and each method has two bars indicating performance under IID and Non-IID settings. The Non-IID results are obtained under a Dirichlet distribution with  $\alpha = 0.5$ . The y-axis represents the average test accuracy across multiple downstream tasks—SST2, AgNews, Yelp, BoolQ, RTE, WSC, and WiC. All detailed results for these tasks are provided in Appendix D.2, Table 15.

### 3.2 ANSWER TO RQ2: GRADIP TRAJECTORIES AS EFFECTIVE INDICATORS OF DATA HETEROGENEITY

This section experimentally validates Claim 2 (Section 2.4), investigating GradIP trajectories as indicators of data heterogeneity. Based on our theoretical analysis assuming single-label Non-IID data (Section C.6), we study the dynamics of gradient-related metrics during local training. We first compare two extremes: IID clients vs. clients with single-label (extreme Non-IID) data. We track three metrics: GradIP score, local gradient norm, and cosine value between the local and pre-training gradients. As shown in Figures 3 and 7, GradIP for extreme Non-IID clients steadily decays to zero over 100 steps, while for IID clients it fluctuates persistently. To understand this, we analyze its components: Figure 8(a) shows cosine value stays near zero (i.e., gradients are nearly orthogonal) for both settings, suggesting the gradient norm is the key factor. Indeed, Figure 8(b) shows that the gradient norm mirrors GradIP’s behavior across the two settings. Moreover, in later stages, GradIP declines more sharply for Non-IID clients than for IID ones, making this stage-wise mean difference an additional criterion for identifying Non-IID clients. We further extend our analysis to more general Non-IID scenarios (Figure 9, Figure 10, Figure 11), where GradIP exhibits similar dynamics that correlate with the degree of heterogeneity.

### 3.3 ANSWER TO RQ3: VPCS EARLY STOPPING EXTREMELY NON-IID DATA DISTRIBUTION CLIENTS

This section experimentally validates Claim 3 (Section 2.5). As established in Section 3.2, GradIP trajectories provide an effective signal for identifying clients with extremely Non-IID data, exhibiting

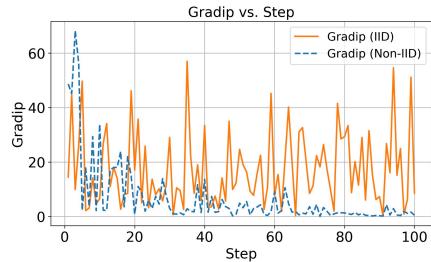
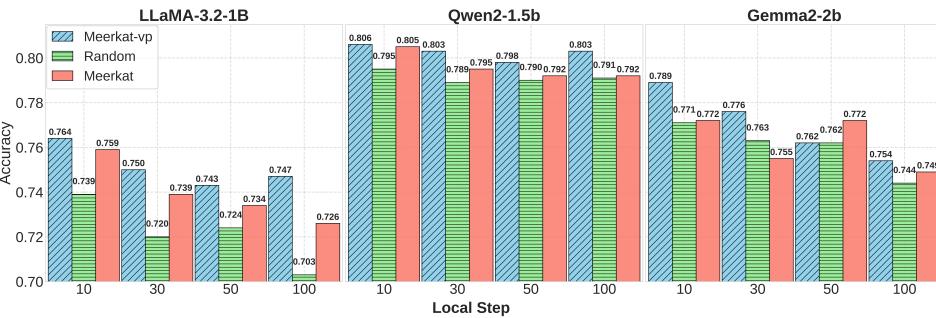


Figure 3: Under a density ratio of  $5 \times 10^{-3}$ , we track the GradIP (see Definition 2.3) over 100 local training steps on the SST-2 dataset using LLaMA-3.2-1B model, comparing a client with IID data to a client with Non-IID data.

432 distinct behaviors. Leveraging this signal, VPCS detects extremely Non-IID clients during a calibration  
 433 phase and applies early stopping, limiting them to one local training step per communication  
 434 round (Algorithm 1). To validate the effectiveness of this VPCS strategy in improving performance,  
 435 we compared MEERKAT-VP with MEERKAT and Random Client Selection, which randomly selects  
 436 the same number of clients for early stopping as VPCS, under Non-IID data distributions dirichlet  
 437  $\alpha = 0.5$  and the same communication frequencies. Crucially, for the same model, dataset, and  
 438 communication frequency, the three methods employed the same sparsity level. Figure 4 illustrates  
 439 the average test accuracy across multiple downstream tasks for MEERKAT-VP compared to  
 440 MEERKAT and RANDOM CLIENT SELECTION. Detailed results for individual tasks are presented  
 441 in Appendix D.2 Table 14. As shown in Figure 4, MEERKAT-VP consistently outperforms both  
 442 MEERKAT and RANDOM CLIENT SELECTION in different communication frequencies. Furthermore,  
 443 Table 25 shows that MEERKAT-VP achieves performance competitive with a back-propagation  
 444 upper bound and significantly outperforms an adapted FedDYN Acar et al. (2021) baseline. These  
 445 experimental results strongly validate Claim 3, confirming that VPCS effectively leverages GradIP  
 446 analysis to manage extremely Non-IID data distribution clients, leading to improved performance for  
 447 ZO federated LLM fine-tuning.



458 Figure 4: This figure compares two methods—MEERKAT-VP, MEERKAT and Random Client Selection—across three LLMs: LLaMA-3.2-1B, Qwen2-1.5b, and Gemma2-2b. The x-axis shows  
 459 the local step values (10, 30, 50, 100), while the y-axis indicates the average test accuracy over  
 460 multiple downstream tasks—SST-2, AgNews, Yelp, BoolQ, RTE, WSC, and WiC—in a Non-IID  
 461 data distribution setting. All detailed results for these tasks are presented in Appendix D.2 Table 14.  
 462

## 4 RELATED WORK

463 Our research leverages advances in ZO federated optimization, sparsity techniques for LLMs, and  
 464 communication frequency adjustments strategies for addressing data heterogeneity. ZO methods  
 465 significantly reduce computational and communication overhead. Integrating sparsity into LLM  
 466 fine-tuning amplifies these benefits, substantially decreasing resource demands during training and  
 467 inference. Concurrently, communication frequency adjustments mitigate performance degradation  
 468 induced by Non-IID data, emphasizing a crucial trade-off between communication budget and global  
 469 model performance. A detailed discussion is provided in Appendix B.

## 5 CONCLUSION

470 In this paper, we introduce MEERKAT, a sparse zeroth-order federated fine-tuning methodology.  
 471 Experiments show MEERKAT outperforms Full-FedZO and other sparsity methods on most tasks at  
 472 equivalent communication frequencies. MEERKAT’s efficiency enables high-frequency communication,  
 473 effectively mitigating Non-IID drift. Moreover, we propose MEERKAT-VP. This methodology  
 474 utilizes VPCS, which analyzes GradIP via virtual paths to enable the selective early stopping of  
 475 extreme Non-IID clients. This approach is shown to improve model performance. Our work thus  
 476 offers effective methods for efficient ZO federated LLM fine-tuning under varying network conditions  
 477 and data heterogeneity. Given the technical focus of this work on algorithm, there are no direct  
 478 negative societal consequences inherent to it that need to be emphasized; potential negative impacts  
 479 would arise from the specific applications where these methods are deployed.

486 ETHICS STATEMENT  
487488 This work follows the ICLR Code of Ethics. This paper aims to advance zeroth-order optimization for  
489 federated LLM fine-tuning by addressing key challenges related to efficiency and data heterogeneity.  
490 All datasets used are publicly available for academic evaluation. Our method is designed to protect  
491 user privacy by operating within the Federated Learning framework, where raw data remains on local  
492 devices. Respecting the broader research community, we acknowledge prior work appropriately and  
493 ensure our contributions are situated within ongoing academic efforts. We declare no conflicts of  
494 interest or external sponsorships associated with this work.  
495496 REPRODUCIBILITY STATEMENT  
497498 We are committed to ensuring the reproducibility of our work. The datasets and models used in this  
499 study are detailed in the experiments section (Section 3). The workflow for MEERKAT is presented  
500 in Figure 1, with further details in Section 2 and Algorithm 2. The workflow for MEERKAT-VP is  
501 demonstrated in Figure 5 and Algorithm 1. All experimental parameters are listed in Appendix 4  
502 and 5. The complete theoretical analysis for our methods can be found in Appendix C.  
503504 REFERENCES  
505

506 Qwen2 technical report. 2024.

508 Durmus Alp Emre Acar, Yue Zhao, Ramon Matas Navarro, Matthew Mattina, Paul N. Whatmough,  
509 and Venkatesh Saligrama. Federated learning based on dynamic regularization, 2021. URL  
510 <https://arxiv.org/abs/2111.04263>.511 Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine  
512 learning, 2018. URL <https://arxiv.org/abs/1606.04838>.514 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
515 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
516 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.517 Ilias Chalkidis, Manos Fergadiotis, and Ion Androutsopoulos. MultiEURLEX - a multi-lingual and  
518 multi-label legal document classification dataset for zero-shot cross-lingual transfer. In Marie-  
519 Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of  
520 the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 6974–6996,  
521 Online and Punta Cana, Dominican Republic, November 2021. Association for Computational  
522 Linguistics. doi: 10.18653/v1/2021.emnlp-main.559. URL <https://aclanthology.org/2021.emnlp-main.559/>.524 Jun Chen, Hong Chen, Bin Gu, and Hao Deng. Fine-grained theoretical analysis of federated  
525 zeroth-order optimization. *Advances in Neural Information Processing Systems*, 36, 2024.527 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina  
528 Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint  
529 arXiv:1905.10044*, 2019.531 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
532 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
533 *arXiv preprint arXiv:2407.21783*, 2024.534 Wenzhi Fang, Ziyi Yu, Yuning Jiang, Yuanming Shi, Colin N Jones, and Yong Zhou. Communication-  
535 efficient stochastic zeroth-order optimization for federated learning. *IEEE Transactions on Signal  
536 Processing*, 70:5058–5073, 2022.538 Pengxin Guo, Shuang Zeng, Yanran Wang, Huijie Fan, Feifei Wang, and Liangqiong Qu. Selective  
539 aggregation for low-rank adaptation in federated learning, 2025. URL <https://arxiv.org/abs/2410.01463>.

540 Wentao Guo, Jikai Long, Yimeng Zeng, Zirui Liu, Xinyu Yang, Yide Ran, Jacob R. Gardner, Osbert  
 541 Bastani, Christopher De Sa, Xiaodong Yu, Beidi Chen, and Zhaozhuo Xu. Zeroth-order fine-tuning  
 542 of llms with extreme sparsity, 2024. URL <https://arxiv.org/abs/2406.02913>.

543

544 Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data  
 545 distribution for federated visual classification, 2019. URL <https://arxiv.org/abs/1909.06335>.

546

547 Weiyu Huang, Yuezhou Hu, Guohao Jian, Jun Zhu, and Jianfei Chen. Pruning large language models  
 548 with semi-structural adaptive sparse training, 2024a. URL <https://arxiv.org/abs/2407.20584>.

549

550

551 Xinmeng Huang, Ping Li, and Xiaoyun Li. Stochastic controlled averaging for federated learning  
 552 with communication compression, 2024b. URL <https://arxiv.org/abs/2308.08165>.

553

554 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Sebastian U. Stich, and  
 555 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning, 2021.  
 556 URL <https://arxiv.org/abs/1910.06378>.

557

558 Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In  
 559 *Thirteenth international conference on the principles of knowledge representation and reasoning*,  
 560 2012.

561

562 Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. Federated learning on non-iid data silos: An  
 563 experimental study, 2021. URL <https://arxiv.org/abs/2102.02079>.

564

565 Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith.  
 566 Federated optimization in heterogeneous networks, 2020a. URL <https://arxiv.org/abs/1812.06127>.

567

568 Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of  
 569 fedavg on non-iid data, 2020b. URL <https://arxiv.org/abs/1907.02189>.

570

571 Zhe Li, Bicheng Ying, Zidong Liu, Chaosheng Dong, and Haibo Yang. Achieving dimension-  
 572 free communication in federated learning via zeroth-order optimization, 2024. URL <https://arxiv.org/abs/2405.15861>.

573

574 Zhenqing Ling, Daoyuan Chen, Liuyi Yao, Yaliang Li, and Ying Shen. On the convergence of  
 575 zeroth-order federated tuning for large language models. In *Proceedings of the 30th ACM SIGKDD  
 Conference on Knowledge Discovery and Data Mining*, pp. 1827–1838, 2024.

576

577 Yong Liu, Zirui Zhu, Chaoyu Gong, Minhao Cheng, Cho-Jui Hsieh, and Yang You. Sparse mezo:  
 578 Less parameters for better performance in zeroth-order llm fine-tuning, 2024. URL <https://arxiv.org/abs/2402.15751>.

579

580 Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhuo Xu, Anastasios  
 581 Kyriolidis, and Anshumali Shrivastava. Scissorhands: Exploiting the persistence of importance  
 582 hypothesis for LLM KV cache compression at test time. In Alice Oh, Tristan Naumann, Amir  
 583 Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural Information  
 584 Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023,  
 585 NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023a.

586

587 Zichang Liu, Jue Wang, Tri Dao, Tianyi Zhou, Binhang Yuan, Zhao Song, Anshumali Shrivastava,  
 588 Ce Zhang, Yuandong Tian, Christopher Re, et al. Deja vu: Contextual sparsity for efficient llms  
 589 at inference time. In *International Conference on Machine Learning*, pp. 22137–22176. PMLR,  
 590 2023b.

591

592 Xudong Lu, Aojun Zhou, Yuhui Xu, Renrui Zhang, Peng Gao, and Hongsheng Li. SPP: Sparsity-  
 593 preserved parameter-efficient fine-tuning for large language models. In *Forty-first International  
 Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=9Rroj9GIOQ>.

594 Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev  
 595 Arora. Fine-tuning language models with just forward passes. *Advances in Neural Information  
 596 Processing Systems*, 36:53038–53075, 2023.

597 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.  
 598 Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence  
 599 and statistics*, pp. 1273–1282. PMLR, 2017.

600 Matias Mendieta, Taojiannan Yang, Pu Wang, Minwoo Lee, Zhengming Ding, and Chen Chen. Local  
 601 learning matters: Rethinking data heterogeneity in federated learning. In *Proceedings of the  
 602 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8397–8406, 2022.

603 Mohammad Taher Pilehvar and Jose Camacho-Collados. Wic: the word-in-context dataset for  
 604 evaluating context-sensitive meaning representations. *arXiv preprint arXiv:1808.09121*, 2018.

605 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
 606 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text  
 607 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.

608 Hang Shao, Bei Liu, and Yanmin Qian. One-shot sensitivity-aware mixed sparsity pruning for large  
 609 language models. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and  
 610 Signal Processing (ICASSP)*, pp. 11296–11300. IEEE, 2024.

611 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and  
 612 Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank.  
 613 In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp.  
 1631–1642, 2013.

614 Junda Su, Zirui Liu, Zeju Qiu, Weiyang Liu, and Zhaozhuo Xu. In defense of structural sparse  
 615 adapters for concurrent llm serving. In *Findings of the Association for Computational Linguistics:  
 616 EMNLP 2024*, pp. 4948–4953, 2024.

617 Jingwei Sun, Ziyue Xu, Hongxu Yin, Dong Yang, Daguang Xu, Yiran Chen, and Holger R. Roth.  
 618 Fedbpt: Efficient federated black-box prompt tuning for large language models, 2023. URL  
 619 <https://arxiv.org/abs/2310.01467>.

620 Gemma Team. Gemma. 2024. doi: 10.34740/KAGGLE/M/3301. URL <https://www.kaggle.com/m/3301>.

621 Alex Wang. Glue: A multi-task benchmark and analysis platform for natural language understanding.  
 622 *arXiv preprint arXiv:1804.07461*, 2018.

623 Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H. Vincent Poor. A novel framework for the  
 624 analysis and design of heterogeneous federated learning. *IEEE Transactions on Signal Processing*,  
 625 69:5234–5249, 2021. doi: 10.1109/TSP.2021.3106104.

626 Ziheng Wang, Jeremy Wohlwend, and Tao Lei. Structured pruning of large language models. *arXiv  
 627 preprint arXiv:1910.04732*, 2019.

628 Zhaozhuo Xu, Zirui Liu, Beidi Chen, Shaochen Zhong, Yuxin Tang, Jue Wang, Kaixiong Zhou, Xia  
 629 Hu, and Anshumali Shrivastava. Soft prompt recovers compressed llms, transferably. In *Forty-first  
 630 International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*.  
 631 OpenReview.net, 2024. URL <https://openreview.net/forum?id=muBJPC1qZT>.

632 Zhiqin Yang, Yonggang Zhang, Yu Zheng, Xinmei Tian, Hao Peng, Tongliang Liu, and Bo Han.  
 633 Fedfed: Feature distillation against data heterogeneity in federated learning. *Advances in Neural  
 634 Information Processing Systems*, 36, 2024.

635 Qingsong Zhang, Bin Gu, Zhiyuan Dang, Cheng Deng, and Heng Huang. Desirable companion for  
 636 vertical federated learning: New zeroth-order gradient based algorithm. In *Proceedings of the 30th  
 637 ACM International Conference on Information & Knowledge Management*, pp. 2598–2607, 2021.

638 Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text  
 639 classification. *Advances in neural information processing systems*, 28, 2015.

648 Yihua Zhang, Pingzhi Li, Junyuan Hong, Jiaxiang Li, Yimeng Zhang, Wenqing Zheng, Pin-Yu Chen,  
 649 Jason D. Lee, Wotao Yin, Mingyi Hong, Zhangyang Wang, Sijia Liu, and Tianlong Chen. Revisiting  
 650 zeroth-order optimization for memory-efficient LLM fine-tuning: A benchmark. In *Forty-first*  
 651 *International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*.  
 652 OpenReview.net, 2024a. URL <https://openreview.net/forum?id=THPjMr2r0S>.

653 Yuxin Zhang, Lirui Zhao, Mingbao Lin, Yunyun Sun, Yiwu Yao, Xingjia Han, Jared Tanner, Shiwei  
 654 Liu, and Rongrong Ji. Dynamic sparse no training: Training-free fine-tuning for sparse llms, 2024b.  
 655 URL <https://arxiv.org/abs/2310.08915>.

656 Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated  
 657 learning with non-iid data. 2018. doi: 10.48550/ARXIV.1806.00582. URL <https://arxiv.org/abs/1806.00582>.

658 Haizhong Zheng, Xiaoyan Bai, Xueshen Liu, Z Morley Mao, Beidi Chen, Fan Lai, and Atul  
 659 Prakash. Learn to be efficient: Build structured sparsity in large language models. *arXiv preprint*  
 660 *arXiv:2402.06126*, 2024.

661 Yang Zhou, Zhuoming Chen, Zhaozhuo Xu, Victoria Lin, and Beidi Chen. Sirius: Contextual sparsity  
 662 with correction for efficient llms. *arXiv preprint arXiv:2409.03856*, 2024.

## 663 APPENDIX

664 In Section A, we discuss the usage of large language model usage in this work. In Section B, we  
 665 present the related work relevant to this study. In Section C, we present the theoretical convergence  
 666 analysis of MEERKAT, including its high-frequency communication method. Additionally, we analyze  
 667 the convergence of MEERKAT-VP and demonstrate its superior performance compared to MEERKAT.  
 668 We further prove that under extreme Non-IID settings, the gradient norm gradually vanishes during  
 669 convergence, whereas in IID settings, it tends to oscillate. In Section D, we provide details on  
 670 experimental hyperparameters and report supplementary results.

## 671 A LLM USAGE

672 We used an LLM-based writing assistant solely for grammar and typographical corrections to  
 673 improve the clarity of this paper. All outputs were carefully reviewed and revised by the authors  
 674 to ensure technical accuracy and consistency with the intended scientific meaning. The intellectual  
 675 contributions, methodological advances, and scientific insights are entirely original and author-driven.

## 676 B REVIEW OF RELATED WORKS

677 **Federated Zeroth-Order Optimization.** Zeroth-order optimization Malladi et al. (2023); Zhang  
 678 et al. (2024a) has gained increasing attention in federated learning Fang et al. (2022); Zhang et al.  
 679 (2021), particularly for addressing challenges in training costs, privacy, and communication overhead.  
 680 Fine-Grained Chen et al. (2024) demonstrates how clients can reduce upload overhead by sending  
 681 estimated gradients rather than full model parameters to the server, though download costs remain  
 682 significant due to complete model weight transfers. DeComFL Li et al. (2024) further advances  
 683 this approach by using gradient scalars for both uploads and downloads, substantially reducing  
 684 bidirectional communication costs. However, it does not address the challenges posed by data  
 685 heterogeneity (Non-IID) in federated learning. The integration of AirComp wireless technology  
 686 enables direct over-the-air aggregation of model updates Fang et al. (2022). In black-box settings  
 687 where pre-trained language model parameters are inaccessible, FedBPT Sun et al. (2023) employs ZO  
 688 to optimize prompt vectors, achieving efficient distributed optimization with reduced computational  
 689 and communication overhead. FedMeZO Li et al. (2020b) analyzes the convergence properties of ZO  
 690 for federated LLM fine-tuning.

691 **Sparsity in LLM.** Current research on sparsity in LLMs explores techniques such as pruning,  
 692 contextual sparsity prediction, and structured sparsity Zhang et al. (2024b); Liu et al. (2023b;a);  
 693 Lu et al. (2024); Zheng et al. (2024); Shao et al. (2024); Wang et al. (2019); Huang et al. (2024a);

Zhou et al. (2024); Su et al. (2024); Xu et al. (2024). These methods enhance both training and inference by improving computational efficiency, reducing memory usage, and enabling deployment in resource-constrained environments. Sparsity has also proven particularly effective in zeroth-order (ZO) optimization Guo et al. (2024); Liu et al. (2024), especially when combined with weight quantization for fine-tuning LLMs. Building on this, our work investigates the role of sparsity in resource-frugal federated fine-tuning of LLMs.

**High-Frequency Communication for Non-IID Federated Learning.** Data heterogeneity across clients is a major challenge in Federated Learning, significantly degrading performance compared to IID settings. Increasing communication frequency, by reducing local training steps per round, is explored as a strategy to mitigate this issue. Early work showed that merely reducing local steps had limited improvements in extreme non-IID scenarios Zhao et al. (2018). Theoretical analysis later confirmed that smaller local training steps can improve convergence speed under Non-IID conditions, but at the cost of increased communication budget, highlighting a critical trade-off Li et al. (2020b). To effectively handle challenges arising from non-IID data that often necessitate higher communication, various algorithms have been proposed: SCAFFOLD Karimireddy et al. (2021) highlights the 'client-drift' problem in FedAvg, noting it's exacerbated by increased local training steps (reduced communication frequency), and proposes using control variates to mitigate this drift, enabling improved convergence; FedDyn Acar et al. (2021) guarantees consistent convergence to the global optimum even with a larger number of local training steps (lower communication frequency). This overcomes the limitation of traditional methods where high communication frequency is needed to compensate for local-global optimum inconsistency. Empirical studies further demonstrate that performance is highly sensitive to the number of local training steps under different non-IID distributions, and the optimal communication frequency depends on the specific data heterogeneity Li et al. (2021). FedSA-LoRA Guo et al. (2025) tackles Non-IID heterogeneity in federated LoRA by showing that the two low-rank matrices play asymmetric roles—A learns global shared knowledge while B captures client-specific patterns—and then aggregating only the A matrices on the server while keeping the B matrices local, which reduces aggregation distortion and cross-client knowledge contamination under skewed client distributions. FedAvgM Hsu et al. (2019) directly studies FedAvg under synthetic label-skewed Dirichlet partitions and proposes adding server-side momentum to mitigate Non-IID client drift: by accumulating past aggregated updates, the server smooths noisy, biased client gradients and recovers much higher test accuracy and more stable training even when client label distributions are highly skewed and only a small fraction of clients participate in each round. Stochastic controlled averaging with compression Huang et al. (2024b) addresses Non-IID data in compressed FL by building SCALLION and SCAFCOM on a simplified SCAFFOLD backbone: control variates and local momentum keep local updates aligned with the global direction, and the authors prove (and verify empirically) that these algorithms remain robust to arbitrary client heterogeneity and partial participation even under aggressive communication compression. These works underscore the complex interplay between data heterogeneity, local computation, and communication frequency. This complexity motivates the development of algorithmic solutions to improve efficiency and robustness in FL under Non-IID settings.

## C THEORETICAL AND ALGORITHM ANALYSIS

### C.1 NOTATIONS AND DEFINITIONS

In this subsection, we formally define the assumptions, notations and concepts used in the convergence analysis of MEERKAT and MEERKAT-VP. Table 2 summarizes the key symbols.

### C.2 ASSUMPTIONS

We introduce the assumptions used in the convergence analysis of MEERKAT and MEERKAT-VP.

**Assumption C.1** (Lipschitz smoothness). We assume that each client  $k$ 's local objective function  $f_k(\mathbf{w})$  is differentiable and has  $L$ -Lipschitz continuous gradients:

$$\|\nabla f_k(\mathbf{w}_1) - \nabla f_k(\mathbf{w}_2)\| \leq L\|\mathbf{w}_1 - \mathbf{w}_2\|, \quad \forall \mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d.$$

Consequently, the global loss  $f(\mathbf{w}) = \sum_{k=1}^K p_k f_k(\mathbf{w})$  is also  $L$ -smooth.

Table 2: Notations used in our theoretical analysis.

| Notation                | Meaning                                                                                      |
|-------------------------|----------------------------------------------------------------------------------------------|
| $\mathbf{w}$            | global model parameter                                                                       |
| $K$                     | total number of clients in the federated system                                              |
| $p_k$                   | probability or weight assigned to client $k$                                                 |
| $f_k(\mathbf{w})$       | total loss computed over all data samples of the client $k$ .                                |
| $f(\mathbf{w})$         | global loss function evaluated by the global model over all data                             |
| $T$                     | number of local update steps per communication round                                         |
| $r$                     | communication round                                                                          |
| $t$                     | local update time step                                                                       |
| $\eta$                  | local learning rate                                                                          |
| $\epsilon$              | perturbation magnitude in ZO estimation                                                      |
| $\mathbf{z}_k^t$        | standard Gaussian vector for client $k$ at local step $t$ from $\mathcal{N}(0, I_d)$         |
| $\mathbf{m}$            | binary sparse mask vector ( $\mathbf{m} \in \{0, 1\}^d$ )                                    |
| $d$                     | model dimension                                                                              |
| $R$                     | federated learning training round                                                            |
| $u$                     | sparsity ratio                                                                               |
| $c$                     | gradient coverage                                                                            |
| $g_k^t$                 | projected gradient estimate for client $k$ at local step $t$                                 |
| $\hat{\nabla} f_k^t$    | zeroth-order gradient of client $k$ at local step $t$                                        |
| $L$                     | Lipschitz smoothness (Assumption 1)                                                          |
| $\mu$                   | PL inequality (Assumption 2)                                                                 |
| $f^*$                   | minimal global loss achieved by optimizing the global model                                  |
| $f_k^*$                 | minimal client loss achieved by optimizing the local model on client $k$                     |
| $c_h$ and $\sigma_h^2$  | heterogeneity-induced variance (Assumption 4)                                                |
| $\ \cdot\ _{\text{op}}$ | operator norm of a matrix                                                                    |
| $\sigma^2$              | variance of the sparse ZO gradient estimator (Assumption 6)                                  |
| $\gamma$                | The clients with balanced data distributions contribute to the global model during training. |

**Assumption C.2** (PL inequality). We assume that  $f(\mathbf{w})$  satisfies the Polyak-Łojasiewicz (PL) condition:

$$f(\mathbf{w}) - f^* \leq \frac{1}{2\mu} \|\nabla f(\mathbf{w})\|^2, \quad \forall \mathbf{w} \in \mathbb{R}^d,$$

$\mu > 0$  is the PL constant. This condition holds for a broad class of non-convex objectives and is commonly used in analyzing convergence of gradient-based and zeroth-order methods.

**Assumption C.3** (Global–Local Disparities in Non-i.i.d. Setting). For any  $\theta \in \mathbb{R}^d$ , the discrepancy between the local and global gradient is bounded by

$$\|\nabla f(\theta) - \nabla f_i(\theta)\|^2 \leq c_h \|\nabla f(\theta)\|^2 + \sigma_h^2,$$

where  $c_h > 0$  and  $\sigma_h^2 \geq 0$  are constants, and  $\theta$  is the global model parameter broadcast to all clients at the start of each round. We further assume  $c_h \in (0, 1)$ . In particular,

- A smaller  $c_h$  corresponds to lower *data heterogeneity*: local gradient deviations from the global gradient are small, indicating that client data distributions are nearly i.i.d.
- A larger  $c_h$  signals stronger non-i.i.d data distribution effects, with greater variation between each client’s gradient and the global gradient.

**Assumption C.4** (Bounded stochastic gradient variance). For any sample  $(x, y) \sim \mathcal{D}$  and any  $\mathbf{w} \in \mathbb{R}^d$ , denote  $f(\mathbf{w}; (x, y))$  as the loss on that single data point, and let  $\bar{f}(\mathbf{w}) := \mathbb{E}_{(x, y) \sim \mathcal{D}} [f(\mathbf{w}; (x, y))]$  be the average full-batch loss. We assume

$$\|\nabla f(\mathbf{w}; (x, y)) - \nabla \bar{f}(\mathbf{w})\|_2^2 \leq \sigma^2.$$

**Assumption C.5** (Local–Global Optimality Gap). For each client  $k$ , define the local–global optimality gap as

$$\Delta_k = \|\mathbf{w}_k^* - \mathbf{w}^*\|_2^2,$$

810 where  $\mathbf{w}_k^*$  is the local optimal model on client  $k$  and  $\mathbf{w}^*$  is the global optimal model.  
 811

812 **Assumption C.6** (Sensitive parameters are sparse). At each local step  $t$  (and for every client  $k$ ), there  
 813 exists a binary mask  $\mathbf{m} \in \{0, 1\}^d$  with exactly  $u$  non-zero entries and a constant  $c \in [0, 1]$  such that

$$814 \quad 815 \quad \|\mathbf{m} \odot \nabla f_k(\mathbf{w}_k^t; (\mathbf{x}_t, \mathbf{y}_t))\|^2 = c \|\nabla f_k(\mathbf{w}_k^t; (\mathbf{x}_t, \mathbf{y}_t))\|^2.$$

816 We further assume  $c \gg \frac{u}{d}$ , meaning this small subset of “sensitive” parameters captures a dispropor-  
 817 tionately large fraction of the gradient norm.  
 818

819 These assumptions are standard and foundational in optimization and FL literature Bottou et al. (2018);  
 820 Li et al. (2020a;b); Wang et al. (2021); Guo et al. (2024)  
 821

822 We start by formulating the expectation of the sensitive sparse ZO surrogate gradient norm square in  
 823 terms of its corresponding stochastic gradient norm square.  
 824

**Lemma C.7** (Sensitive sparse ZO surrogate gradient norm square).

$$825 \quad 826 \quad \mathbb{E}_{\bar{z}} \left\| \hat{\nabla} f(w_t, (x_t, y_t), \bar{z}_t) \right\|^2 = (2 + u)c \left\| \nabla f(w_t; (x_t, y_t)) \right\|^2.$$

827 *Proof.* Our masked perturbation  $\bar{z}$  is sampled as  $\bar{z} \sim \mathcal{N}(0, \tilde{I}_{d, \mathbf{m}})$ , where  $\tilde{I}_{d, \mathbf{m}}$  equals the identity  
 828 matrix  $I_d$  with its main diagonal masked by  $\mathbf{m}$ .  
 829

830 We expand the sensitive sparse ZO surrogate–gradient covariance matrix:  
 831

$$832 \quad 833 \quad \mathbb{E}_{\bar{z}} \hat{\nabla} f(w, (x, y), \bar{z}) \hat{\nabla} f(w, (x, y), \bar{z})^\top \\ 834 \quad = \mathbb{E}_{\bar{z}} [\bar{z} \bar{z}^\top ((\mathbf{m} \odot \nabla f(w; (x, y))) (\mathbf{m} \odot \nabla f(w; (x, y)))^\top) \bar{z} \bar{z}^\top] \\ 835 \quad = 2((\mathbf{m} \odot \nabla f(w; (x, y))) (\mathbf{m} \odot \nabla f(w; (x, y)))^\top) + \|\mathbf{m} \odot \nabla f(w; (x, y))\|^2 \tilde{I}_{d, \mathbf{m}}$$

836 The above expected squared norm is obtained by summing the diagonal elements of this covariance  
 837 matrix:  
 838

$$839 \quad 840 \quad \mathbb{E}_{\bar{z}} \left\| \hat{\nabla} f(w_t, x_t, \bar{z}_t) \right\|^2 = (\text{diag}[\mathbb{E}_{\bar{z}} \hat{\nabla} f(w, (x, y), \bar{z}) \hat{\nabla} f(w, (x, y), \bar{z})^\top])^2 \\ 841 \quad = 2c \|\nabla f(w_t; (x_t, y_t))\|^2 + uc \|\nabla f(w_t; (x_t, y_t))\|^2 \\ 843 \quad = (2 + u)c \|\nabla f(w_t; (x_t, y_t))\|^2.$$

845  $\square$   
 846

847 **Lemma C.8** (Unbiasedness of Masked Sparse ZO Surrogate Gradient).  
 848

$$849 \quad \mathbb{E}_{\bar{z}} [\hat{\nabla} f_k(\mathbf{w}_k^t, \bar{z})] = \mathbf{m} \odot \nabla f_k(\mathbf{w}_k^t), \quad \text{where } \bar{z} = z \odot \mathbf{m}. \quad (4)$$

850 *Proof.* First, consider the estimator defined as:  
 851

$$852 \quad 853 \quad \hat{\nabla} f_k(\mathbf{w}_k^t, z) = \frac{f_k(\mathbf{w}_k^t + \epsilon(z \odot \mathbf{m})) - f_k(\mathbf{w}_k^t - \epsilon(z \odot \mathbf{m}))}{2\epsilon} \cdot (z \odot \mathbf{m}).$$

855 To proceed, we apply a first-order Taylor expansion of  $f_k$  around  $\mathbf{w}_k^t$  for small  $\epsilon$ :  
 856

$$857 \quad f_k(\mathbf{w}_k^t \pm \epsilon(z \odot \mathbf{m})) = f_k(\mathbf{w}_k^t) \pm \epsilon \langle \nabla f_k(\mathbf{w}_k^t), z \odot \mathbf{m} \rangle + \mathcal{O}(\epsilon^2).$$

858 Substitute these expansions into the numerator of the estimator:  
 859

$$860 \quad 861 \quad f_k(\mathbf{w}_k^t + \epsilon(z \odot \mathbf{m})) - f_k(\mathbf{w}_k^t - \epsilon(z \odot \mathbf{m})) \\ 862 \quad = [f_k(\mathbf{w}_k^t) + \epsilon \langle \nabla f_k(\mathbf{w}_k^t), z \odot \mathbf{m} \rangle] \\ 863 \quad - [f_k(\mathbf{w}_k^t) - \epsilon \langle \nabla f_k(\mathbf{w}_k^t), z \odot \mathbf{m} \rangle] + \mathcal{O}(\epsilon^2).$$

864  
865

Simplify the expression:

866  
867

$$f_k(\mathbf{w}_k^t + \epsilon(z \odot \mathbf{m})) - f_k(\mathbf{w}_k^t - \epsilon(z \odot \mathbf{m})) = 2\epsilon \langle \nabla f_k(\mathbf{w}_k^t), z \odot \mathbf{m} \rangle + \mathcal{O}(\epsilon^2).$$

868

Thus, the estimator becomes:

869  
870  
871

$$\hat{\nabla} f_k(\mathbf{w}_k^t, z) = \frac{2\epsilon \langle \nabla f_k(\mathbf{w}_k^t), z \odot \mathbf{m} \rangle + \mathcal{O}(\epsilon^2)}{2\epsilon} \cdot (z \odot \mathbf{m}) = [\langle \nabla f_k(\mathbf{w}_k^t), z \odot \mathbf{m} \rangle + \mathcal{O}(\epsilon)] (z \odot \mathbf{m}).$$

872  
873As  $\epsilon \rightarrow 0$ , the  $\mathcal{O}(\epsilon)$  term disappears, yielding the approximation:874  
875

$$\hat{\nabla} f_k(\mathbf{w}_k^t, z) \approx \langle \nabla f_k(\mathbf{w}_k^t), z \odot \mathbf{m} \rangle \cdot (z \odot \mathbf{m}).$$

876  
877  
878Next, compute the expectation  $\mathbb{E}_z [\hat{\nabla} f_k(\mathbf{w}_k^t, z)]$ . Since the estimator is a vector, consider its  $j$ -th component:879  
880

$$[\hat{\nabla} f_k(\mathbf{w}_k^t, z)]_j \approx \langle \nabla f_k(\mathbf{w}_k^t), z \odot \mathbf{m} \rangle \cdot (z_j m_j).$$

881  
882

Express the inner product explicitly:

883  
884  
885  
886

$$\langle \nabla f_k(\mathbf{w}_k^t), z \odot \mathbf{m} \rangle = \sum_{i=1}^d (\nabla f_k(\mathbf{w}_k^t))_i z_i m_i.$$

887

Thus, the  $j$ -th component is:888  
889  
890  
891

$$[\hat{\nabla} f_k(\mathbf{w}_k^t, z)]_j \approx \left( \sum_{i=1}^d (\nabla f_k(\mathbf{w}_k^t))_i z_i m_i \right) z_j m_j.$$

892  
893Now, take the expectation over  $z \sim \mathcal{N}(0, \mathbf{I}_d)$ , where  $z_i$  are independent standard normal variables:894  
895  
896

$$\mathbb{E}_z \left[ \left( \sum_{i=1}^d (\nabla f_k(\mathbf{w}_k^t))_i z_i m_i \right) z_j m_j \right] = \sum_{i=1}^d (\nabla f_k(\mathbf{w}_k^t))_i m_i m_j \mathbb{E}[z_i z_j].$$

897

Since  $\mathbb{E}[z_i z_j] = \delta_{ij}$  (1 if  $i = j$ , 0 otherwise), the sum reduces to:898  
899  
900

$$(\nabla f_k(\mathbf{w}_k^t))_j m_j^2 \mathbb{E}[z_j^2].$$

901  
902Given  $m_j^2 = m_j$  (as  $m_j = 0$  or 1) and  $\mathbb{E}[z_j^2] = 1$ , this becomes:903  
904

$$(\nabla f_k(\mathbf{w}_k^t))_j m_j.$$

905  
906Thus, for each component  $j$ :907  
908  
909

$$\mathbb{E}_z \left[ [\hat{\nabla} f_k(\mathbf{w}_k^t, z)]_j \right] \approx m_j (\nabla f_k(\mathbf{w}_k^t))_j.$$

910

This implies:

911  
912  
913

$$\mathbb{E}_z [\hat{\nabla} f_k(\mathbf{w}_k^t, z)] \approx \mathbf{m} \odot \nabla f_k(\mathbf{w}_k^t).$$

914  
915Finally, as  $\epsilon \rightarrow 0$ , the higher-order terms in the Taylor expansion vanish, making the approximation exact:916  
917

$$\mathbb{E}_{\bar{z}} [\hat{\nabla} f_k(\mathbf{w}_k^t, \bar{z})] = \mathbf{m} \odot \nabla f_k(\mathbf{w}_k^t).$$

□

918 C.3 MEERKAT CONVERGENCE ANALYSIS  
919920 We consider the **federated zeroth-order optimization problem**, where the objective is to minimize  
921 the global loss function Ling et al. (2024):  
922

923  
924 
$$\min_{\mathbf{w}} f(\mathbf{w}) = \sum_{k=1}^K p_k f_k(\mathbf{w})$$
  
925  
926

927 Each client performs  $T$  local steps:  
928

929  
930 
$$\mathbf{w}_k^{t+1} = \mathbf{w}_k^t - \eta \nabla f_k^t(\mathbf{w}), \quad t = 0, 1, \dots, T-1$$
  
931

932 starting from the global model  $\mathbf{w}_k^0 = \mathbf{w}^r$ . After clients finish local updates, the server performs  
933 weighted aggregation of their model updates.  
934

935  
936 
$$\mathbf{w}^{r+1} = \sum_{k=1}^K p_k \mathbf{w}_k^r.$$
  
937  
938

939 **Theorem C.9.** [Client Local ZO Update Convergence] Let  $f_k$  be  $L$ -smooth and  $\hat{\nabla} f_k^t$  be an unbiased  
940 sparse zeroth-order gradient estimator with variance bounded by  $\sigma^2$ . Then we have  
941942 If we set constant learning rate  $\eta = \frac{1}{L(u+2)}$  and  $T$  local steps, the output of client  $k$  satisfies:  
943

944 
$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f_k(\mathbf{w}_k^t)\|^2 \leq \mathcal{O}\left(\frac{1}{T}\right) + \mathcal{O}(\sigma^2). \quad (5)$$
  
945

946 *Proof.* We start by proving Theorem C.9 equation 5 that each client achieves local convergence  
947 during training with sparse zeroth-order finetuning. Next, we demonstrate that server-side aggregation  
948 also converge. Finally, by leveraging the PL inequality, we prove that MEERKAT exhibits linear  
949 convergence to global minimum.  
950951 **Part 1: Client Local ZO Update Convergence**  
952953 We analyze the effect of one local step of MEERKAT under sparse zeroth-order updates. Let client  $k$   
954 perform the local update:  
955

956 
$$\mathbf{w}_k^{t+1} = \mathbf{w}_k^t - \eta \hat{\nabla} f_k^t,$$
  
957

958 where the stochastic sparse zeroth-order gradient estimator is defined as:  
959

960 
$$g_k^t = \frac{f_k(\mathbf{w}_k^t + \epsilon(\mathbf{z}_k^t \odot \mathbf{m})) - f_k(\mathbf{w}_k^t - \epsilon(\mathbf{z}_k^t \odot \mathbf{m}))}{2\epsilon}.$$
  
961

962 
$$\hat{\nabla} f_k^t = g_k^t \cdot (\mathbf{z}_k^t \odot \mathbf{m})$$
  
963

964 **Descent via Lipschitz smoothness.** Since  $f_k(\mathbf{w})$  is Lipschitz smoothness:  
965

966 
$$f_k(\mathbf{w}_k^{t+1}) \leq f_k(\mathbf{w}_k^t) + \langle \nabla f_k(\mathbf{w}_k^t), \mathbf{w}_k^{t+1} - \mathbf{w}_k^t \rangle + \frac{L}{2} \|\mathbf{w}_k^{t+1} - \mathbf{w}_k^t\|^2.$$
  
967

968 Substituting the update  $\mathbf{w}_k^{t+1} - \mathbf{w}_k^t = -\eta \hat{\nabla} f_k^t$ , we obtain:  
969

970 
$$f_k(\mathbf{w}_k^{t+1}) \leq f_k(\mathbf{w}_k^t) - \eta \langle \nabla f_k(\mathbf{w}_k^t), \hat{\nabla} f_k^t(\mathbf{w}, \bar{\mathbf{z}}_t) \rangle + \frac{L\eta^2}{2} \|\hat{\nabla} f_k^t(\mathbf{w}, \bar{\mathbf{z}}_t)\|^2.$$
  
971

972 Taking expectation, we have:  
973

974 
$$\mathbb{E}_{\bar{\mathbf{z}}}[f_k(\mathbf{w}_k^{t+1})] \leq \mathbb{E}_{\bar{\mathbf{z}}}[f_k(\mathbf{w}_k^t)] - \eta \mathbb{E}_{\bar{\mathbf{z}}} \|\mathbf{m} \odot \nabla f_k(\mathbf{w}_k^t)\|^2 + \frac{L\eta^2}{2} \mathbb{E}_{\bar{\mathbf{z}}} \|\hat{\nabla} f_k(\mathbf{w}_k^t, \bar{\mathbf{z}}_t)\|^2.$$
  
975

972

$$\mathbb{E}_{\bar{\mathbf{z}}}[f_k(\mathbf{w}_k^{t+1})] \leq \mathbb{E}_{\bar{\mathbf{z}}}[f_k(\mathbf{w}_k^t)] - c\eta \mathbb{E}_{\bar{\mathbf{z}}} \|\nabla f_k(\mathbf{w}_k^t)\|^2 + \frac{L\eta^2}{2}(2+u)c\mathbb{E}_{\bar{\mathbf{z}}} \|\nabla f_k(\mathbf{w}_k^t)\|^2.$$

973

$$\mathbb{E}_{\bar{\mathbf{z}}} f_k(\mathbf{w}_k^{t+1}) \leq \mathbb{E}_{\bar{\mathbf{z}}} f_k(\mathbf{w}_k^t) - \left( c\eta_t - \frac{L\eta_t^2}{2}c(u+2) \right) \|\nabla_{\mathbf{w}} f_k(\mathbf{w}_k^t)\|^2 + \frac{L\eta_t^2}{2}c(u+2)\sigma^2.$$

974

975 Denote  $\alpha = Lc(u+2)$ , we can rewrite as:

976

$$\mathbb{E}_{\bar{\mathbf{z}}} f_k(\mathbf{w}_k^{t+1}) \leq \mathbb{E}_{\bar{\mathbf{z}}} \left\{ f_k(\mathbf{w}_k^t) - \eta_t \left( c - \frac{\alpha}{2}\eta_t \right) \|\nabla_{\mathbf{w}} f_k(\mathbf{w}_k^t)\|^2 \right\} + \frac{\alpha}{2}\sigma^2\eta_t^2.$$

977

978 From the above inequality, we get  $\eta < \frac{2c}{\alpha}$ . Suppose we use a constant learning rate  $\eta_t = \eta = \frac{c}{\alpha} = \frac{1}{L(u+2)}$ , we get:

979

$$\mathbb{E}_{\bar{\mathbf{z}}} f_k(\mathbf{w}_k^{t+1}) \leq \mathbb{E}_{\bar{\mathbf{z}}} \left\{ f_k(\mathbf{w}_k^t) - \frac{c\eta}{2} \|\nabla_{\mathbf{w}} f_k(\mathbf{w}_k^t)\|^2 \right\} + \frac{\alpha}{2}\sigma^2\eta^2. \quad (6)$$

980

981 **Accumulating over  $T$  steps.** Summing equation 6 over  $t = 0$  to  $T-1$ , we get:

982

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{\bar{\mathbf{z}}} \|\nabla f_k(\mathbf{w}_k^t)\|^2 &\leq \frac{2}{c\eta T} (f_k(\mathbf{w}_k^0) - f_k^*) + \frac{1}{T} \sum_{t=0}^{T-1} \frac{\alpha}{2c\eta} \sigma^2 \eta^2 \\ &= \frac{2L(u+2)}{cT} (f_k(\mathbf{w}_k^0) - f_k^*) + \sigma^2 \\ &= \mathcal{O}\left(\frac{u}{T}(f_k(\mathbf{w}_k^0) - f_k^*)\right) + \mathcal{O}(1). \end{aligned} \quad (7)$$

983

984

#### 1000 C.4 MEERKAT CONVERGENCE ANALYSIS

1001

1002 We now proceed to analyze the convergence of the global model in our federated learning framework.  
1003 Having established the convergence properties of local client updates, we demonstrate how these  
1004 results extend to guarantee the convergence of the server-aggregated global model.

1005

1006 *Proof.* We approach this proof systematically by analyzing how the local convergence properties of  
1007 clients extend to the global model through the aggregation process.

1008

1009 **Global Model Update Representation.** First, the global model update can be represented as:

1010

$$\mathbf{w}^{r+1} - \mathbf{w}^r = \sum_{k=1}^K p_k (\mathbf{w}_k^T - \mathbf{w}^r)$$

1011

1012

1013 where each client  $k$  starts from the global model  $\mathbf{w}^r$  and performs  $T$  local updates to reach  $\mathbf{w}_k^T$ .

1014

1015

1016 **Client Local Update Accumulation** For any client  $k$ , the accumulated local updates can be expressed  
1017 as:

1018

$$\mathbf{w}_k^{r,T} - \mathbf{w}^r = -\eta \sum_{t=0}^{T-1} \hat{\nabla} f_k^t$$

1019

1020

1021 **Global Loss Descent Analysis** By the  $L$ -smoothness property (Assumption C.1), we have:

1022

1023

$$f(\mathbf{w}^{r+1}) \leq f(\mathbf{w}^r) + \langle \nabla f(\mathbf{w}^r), \mathbf{w}^{r+1} - \mathbf{w}^r \rangle + \frac{L}{2} \|\mathbf{w}^{r+1} - \mathbf{w}^r\|^2 \quad (8)$$

1024

1025

1026 For the inner product we can get:

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1026 According to the client local update process, we have:  
 1027

$$\begin{aligned}
 1029 \quad \sum_{k=1}^K p_k \langle \nabla f_k(\mathbf{w}^r), \mathbf{w}_k^{r,T} - \mathbf{w}^r \rangle &= -\eta \sum_{k=1}^K p_k \langle \nabla f(\mathbf{w}^r), \sum_{t=0}^{T-1} \hat{\nabla} f_k(\mathbf{w}^{r,t}, \bar{\mathbf{z}}_t) \rangle \\
 1030 \quad &= -\eta \sum_{k=1}^K p_k \sum_{t=0}^{T-1} \langle \nabla f(\mathbf{w}^r), \hat{\nabla} f_k(\mathbf{w}^{r,t}, \bar{\mathbf{z}}_t) \rangle
 \end{aligned}$$

1035  
 1036 We assume that each client's weight is equal  $p_k = 1/K$ , by substituting it into the above inequality,  
 1037 we have:  
 1038

$$\sum_{k=1}^K p_k \langle \nabla f_k(\mathbf{w}^r), \mathbf{w}_k^{r,T} - \mathbf{w}^r \rangle = -\frac{\eta}{K} \sum_{k=1}^K \sum_{t=0}^{T-1} \langle \nabla f(\mathbf{w}^r), \hat{\nabla} f_k(\mathbf{w}^{r,t}, \bar{\mathbf{z}}_t) \rangle. \quad (9)$$

1043 Based on the equation 9 and  $\hat{\nabla} f_k^t$  is unbiased, we have:  
 1044

$$\sum_{k=1}^K \sum_{t=0}^{T-1} \langle \nabla f(\mathbf{w}^r), \hat{\nabla} f_k(\mathbf{w}^{r,t}, \bar{\mathbf{z}}_t) \rangle = \sum_{k=1}^K \sum_{t=0}^{T-1} \langle \nabla f(\mathbf{w}^r), \mathbb{E}_{\bar{\mathbf{z}}}[\hat{\nabla} f_k(\mathbf{w}^{r,t}, \bar{\mathbf{z}}_t)] \rangle$$

1045 We substitute the equation 4 and get:  
 1046

$$\sum_{k=1}^K \sum_{t=0}^{T-1} \langle \nabla f(\mathbf{w}^r), \mathbb{E}_{\bar{\mathbf{z}}}[\hat{\nabla} f_k(\mathbf{w}^{r,t}, \bar{\mathbf{z}}_t)] \rangle = \sum_{k=1}^K \sum_{t=0}^{T-1} \langle \nabla f(\mathbf{w}^r), \mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t}) \rangle.$$

1053 Under the Cauchy–Schwarz inequality, we have:  
 1054

$$\langle \nabla f(\mathbf{w}^r), \mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t}) \rangle \leq \|\nabla f(\mathbf{w}^r)\| \|\mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t})\|$$

1055 We substitute Assumption C.6 get:  
 1056

$$\|\nabla f(\mathbf{w}^r)\| \|\mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t})\| = \sqrt{c} \|\nabla f(\mathbf{w}^r)\| \|\nabla f_k(\mathbf{w}^{r,t})\|.$$

1057 Thus we get:  
 1058

$$\langle \nabla f(\mathbf{w}^r), \mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t}) \rangle \leq \sqrt{c} \|\nabla f(\mathbf{w}^r)\| \|\nabla f_k(\mathbf{w}^{r,t})\|.$$

1059 By the triangle inequality, we have  
 1060

$$\|\nabla f_k(\mathbf{w}^{r,t})\| \leq \|\nabla f(\mathbf{w}^r)\| + \|\nabla f_k(\mathbf{w}^{r,t}) - \nabla f(\mathbf{w}^r)\|$$

1061 We substitute Assumption C.3 and use the properties of square roots we get:  
 1062

$$\begin{aligned}
 1063 \quad &\|\nabla f(\mathbf{w}^r)\| + \|\nabla f_k(\mathbf{w}^{r,t}) - \nabla f(\mathbf{w}^r)\| \\
 1064 \quad &\leq \|\nabla f(\mathbf{w}^r)\| + \sqrt{c_h \|\nabla f(\mathbf{w}^r)\|^2 + \sigma_h^2} \\
 1065 \quad &\leq (1 + \sqrt{c_h}) \|\nabla f(\mathbf{w}^r)\| + \sigma_h.
 \end{aligned}$$

1066 Using the bound  $\langle \nabla f(\mathbf{w}^r), \mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t}) \rangle \leq \sqrt{c} \|\nabla f(\mathbf{w}^r)\| \|\nabla f_k(\mathbf{w}^{r,t})\|$  from Cauchy–Schwarz and  
 1067 Assumption C.6, and then plugging in the above, we obtain  
 1068

$$\begin{aligned}
 1069 \quad \langle \nabla f(\mathbf{w}^r), \mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t}) \rangle &\leq \sqrt{c} \|\nabla f(\mathbf{w}^r)\| [(1 + \sqrt{c_h}) \|\nabla f(\mathbf{w}^r)\| + \sigma_h] \\
 1070 \quad &\leq \sqrt{c} (1 + \sqrt{c_h}) \|\nabla f(\mathbf{w}^r)\|^2 + \sqrt{c} \sigma_h \|\nabla f(\mathbf{w}^r)\|.
 \end{aligned}$$

1071 Recall that the server update inner product is  
 1072

$$\langle \nabla f(\mathbf{w}^r), \mathbf{w}^{r+1} - \mathbf{w}^r \rangle = -\frac{\eta}{K} \sum_{k=1}^K \sum_{t=0}^{T-1} \langle \nabla f, \mathbf{m} \odot \nabla f_k \rangle.$$

1080

Substituting the bound to equation 9. We have:

1081

1082

1083

$$\langle \nabla f(w^r), w^{r+1} - w^r \rangle \geq -\eta T \sqrt{c} (1 + \sqrt{c_h}) \|\nabla f(w^r)\|^2 - \eta T \sqrt{c} \sigma_h \|\nabla f(w^r)\|. \quad (10)$$

1084

1085

1086

Substituting this inequality to equation 8, we have:

1087

1088

1089

1090

1091

1092

1093

Applying Jensen's inequality, the last term of the equation 11 will be:

1094

1095

1096

1097

And then we apply Cauchy-Schwarz inequality, the last term of the equation 11 will be:

1098

1099

1100

1101

1102

Substitute this inequaltiy to equation 11 We get:

1103

1104

1105

1106

1107

1108

1109

Taking Expectation and lemma C.7:

1110

1111

1112

1113

1114

1115

1116

1117

1118

According to the equation 7, we know that the client-average squared gradient has upper bound. We substitute the equation 7 to the above inequality last term we get:

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

$$\begin{aligned} f(w^{r+1}) &\leq f(w^r) - \eta T \sqrt{c} (1 + \sqrt{c_h}) \|\nabla f(w^r)\|^2 - \eta T \sqrt{c} \sigma_h \|\nabla f(w^r)\| \\ &\quad + \frac{L}{2} \eta^2 T \sum_{k=1}^K p_k \sum_{t=0}^{T-1} \|\hat{\nabla} f_k^{r,t}\|^2. \end{aligned}$$

According to the equation 7, we know that the client-average squared gradient has upper bound. We substitute the equation 7 to the above inequality last term we get:

$$\begin{aligned} \mathbb{E}_{\bar{z}} f(w^{r+1}) &\leq \mathbb{E}_{\bar{z}} f(w^r) - \eta T \sqrt{c} (1 + \sqrt{c_h}) \|\nabla f(w^r)\|^2 - \eta T \sqrt{c} \sigma_h \|\nabla f(w^r)\| \\ &\quad + \frac{L \eta^2 T (2+u) c}{2K} \sum_{k=1}^K \sum_{t=0}^{T-1} \|\nabla f_k(w^{r,t})\|^2. \end{aligned} \quad (12)$$

**Accumulating Over  $R$  Rounds.** Summing equation 12 over  $r = 0$  to  $R - 1$ , we get:

1134  
 1135  
 1136 
$$\mathbb{E}_{\bar{z}}[f(w^R)] - \mathbb{E}_{\bar{z}}[f(w^0)] \leq -\eta T \sqrt{c} (1 + \sqrt{c_h}) \sum_{r=0}^{R-1} \|\nabla f(w^r)\|^2$$
  
 1137  
 1138  
 1139 
$$- \eta T \sqrt{c} \sigma_h \sum_{r=0}^{R-1} \|\nabla f(w^r)\|$$
  
 1140  
 1141  
 1142 
$$+ \frac{L^2 \eta^2 T (2+u) (u+2)}{K} \sum_{r=0}^{R-1} \sum_{k=1}^K (f_k(w^r) - f_k^*)$$
  
 1143  
 1144  
 1145 
$$+ \frac{L \eta^2 T^2 (2+u) c}{2} \sigma^2 R.$$
  
 1146

1147 From the accumulated global descent inequality over  $R$  rounds:  
 1148

1149 First we set

1150 
$$S = \sum_{r=0}^{R-1} \|\nabla f(w^r)\|^2.$$
  
 1151  
 1152

1153 This represents the sum of squared gradient norms over  $R$  rounds. The second term in the inequality  
 1154 involves  $\sum_{r=0}^{R-1} \|\nabla f(w^r)\|$ , and we apply the Cauchy-Schwarz inequality to it. For the sequence  
 1155  $a_r = \|\nabla f(w^r)\|$  (with  $r = 0, 1, \dots, R-1$ ), we consider it as a vector in  $\mathbb{R}^R$  along with a vector of  
 1156 ones:

1157 
$$\sum_{r=0}^{R-1} \|\nabla f(w^r)\| = \sum_{r=0}^{R-1} \|\nabla f(w^r)\| \cdot 1 \leq \sqrt{\sum_{r=0}^{R-1} \|\nabla f(w^r)\|^2} \cdot \sqrt{\sum_{r=0}^{R-1} 1^2}.$$
  
 1158  
 1159

1160 Since  $\sum_{r=0}^{R-1} 1^2 = R$ , we obtain:  
 1161

1162 
$$\sum_{r=0}^{R-1} \|\nabla f(w^r)\| \leq \sqrt{\sum_{r=0}^{R-1} \|\nabla f(w^r)\|^2} \cdot \sqrt{R} = \sqrt{R} \sqrt{S} = \sqrt{RS}.$$
  
 1163  
 1164

1165 Substituting this into the second term, we have:

1166  
 1167 
$$\eta T \sqrt{c} \sigma_h \sum_{r=0}^{R-1} \|\nabla f(w^r)\| \leq \eta T \sqrt{c} \sigma_h \sqrt{RS}.$$
  
 1168  
 1169

1170 Thus, the inequality becomes:

1171 
$$\mathbb{E}_{\bar{z}}[f(w^R)] - \mathbb{E}_{\bar{z}}[f(w^0)] \leq -\eta T \sqrt{c} (1 + \sqrt{c_h}) S + \eta T \sqrt{c} \sigma_h \sqrt{RS}$$
  
 1172  
 1173 
$$+ \frac{L^2 \eta^2 T (2+u) (u+2)}{K} \sum_{r=0}^{R-1} \sum_{k=1}^K (f_k(w^r) - f_k^*)$$
  
 1174  
 1175 
$$+ \frac{L \eta^2 T^2 (2+u) c}{2} \sigma^2 R.$$
  
 1176  
 1177

1178 Second, we focus on the term  $\eta T \sqrt{c} \sigma_h \sqrt{RS}$  and apply Young's Inequality with  $\delta > 0$  and non-  
 1179 negative real numbers  $x$  and  $y$ ,

1180 
$$xy \leq \frac{x^2}{2\delta} + \frac{y^2\delta}{2}.$$
  
 1181  
 1182

We identify  $x = \sqrt{S}$  and  $y = \eta T \sqrt{c} \sigma_h \sqrt{R}$ , since:

1184 
$$\eta T \sqrt{c} \sigma_h \sqrt{RS} = (\eta T \sqrt{c} \sigma_h \sqrt{R}) \cdot \sqrt{S}.$$
  
 1185  
 1186

Applying Young's Inequality:

1187 
$$\sqrt{S} \cdot (\eta T \sqrt{c} \sigma_h \sqrt{R}) \leq \frac{(\sqrt{S})^2}{2\delta} + \frac{(\eta T \sqrt{c} \sigma_h \sqrt{R})^2 \delta}{2}.$$

1188 Therefore:

$$1189 \eta T \sqrt{c} \sigma_h \sqrt{RS} \leq \frac{S}{2\delta} + \frac{\eta^2 T^2 c \sigma_h^2 R \delta}{2}.$$

$$1192 -\eta T \sqrt{c} \sigma_h \sqrt{RS} \leq \frac{S}{2\delta} + \frac{\eta^2 T^2 c \sigma_h^2 R \delta}{2}.$$

1194 Finally we replace the second term in the inequality with the above result:

$$1197 \mathbb{E}_{\bar{z}}[f(w^R)] - \mathbb{E}_{\bar{z}}[f(w^0)] \leq -\eta T \sqrt{c} (1 + \sqrt{c_h}) S + \left( \frac{S}{2\delta} + \frac{\eta^2 T^2 c \sigma_h^2 R \delta}{2} \right) \\ 1198 + \frac{L^2 \eta^2 T (2+u)^2}{K} \sum_{r=0}^{R-1} \sum_{k=1}^K (f_k(w^r) - f_k^*) \\ 1199 + \frac{L \eta^2 T^2 (2+u) c}{2} \sigma^2 R.$$

1204 This inequality now depends on  $\delta$ .

$$1207 \left( \eta T \sqrt{c} (1 + \sqrt{c_h}) - \frac{1}{2\delta} \right) \sum_{r=0}^{R-1} \|\nabla f(w^r)\|^2 \leq \mathbb{E}_{\bar{z}}[f(w^0) - f(w^R)] + \eta^2 T^2 c \sigma_h^2 R \delta 2 \\ 1208 + \frac{L^2 \eta^2 T (2+u)^2}{K} \sum_{r=0}^{R-1} \sum_{k=1}^K (f_k(w^r) - f_k^*) \quad (14) \\ 1209 + \frac{L \eta^2 T^2 (2+u) c \sigma^2 R}{2}.$$

1216 According to Assumption C.1, we have:

$$1218 f_k(\mathbf{w}^*) \leq f_k(\mathbf{w}_k^*) + \langle \nabla f_k(\mathbf{w}_k^*), \mathbf{w}^* - \mathbf{w}_k^* \rangle + \frac{L}{2} \|\mathbf{w}^* - \mathbf{w}_k^*\|_2^2.$$

1220 Since  $\mathbf{w}_k^*$  is the minimizer of  $f_k(\mathbf{w})$ , the gradient at the local optimum must be zero:

$$1222 \nabla f_k(\mathbf{w}_k^*) = 0.$$

1224 Substituting this into the inner product term:

$$1225 \langle \nabla f_k(\mathbf{w}_k^*), \mathbf{w}^* - \mathbf{w}_k^* \rangle = \langle 0, \mathbf{w}^* - \mathbf{w}_k^* \rangle = 0.$$

1226 Thus, the inner product term disappears because the gradient at  $\mathbf{w}_k^*$  is zero, making the inner product with any vector (including  $\mathbf{w}^* - \mathbf{w}_k^*$ ) equal to zero.

1228 With the inner product term vanishing, the inequality simplifies to:

$$1230 f_k(\mathbf{w}^*) \leq f_k(\mathbf{w}_k^*) + \frac{L}{2} \Delta_k.$$

1233 This provides an upper bound on  $f_k(\mathbf{w}^*)$  in terms of the local optimal loss  $f_k^*$  and the optimality gap  $\Delta_k$ .

1235 The global optimal loss is defined as:

$$1237 f^* = f(\mathbf{w}^*) = \sum_{k=1}^K p_k f_k(\mathbf{w}^*).$$

1239 Using the bound derived for each local loss:

$$1241 f_k(\mathbf{w}^*) \leq f_k^* + \frac{L}{2} \Delta_k,$$

1242 we substitute this into the expression for  $f^*$ :  
 1243

$$1244 \quad f^* = \sum_{k=1}^K p_k f_k(\mathbf{w}^*) \leq \sum_{k=1}^K p_k \left( f_k^* + \frac{L}{2} \Delta_k \right).$$

1247 Expanding the right-hand side:  
 1248

$$1249 \quad f^* \leq \sum_{k=1}^K p_k f_k^* + \frac{L}{2} \sum_{k=1}^K p_k \Delta_k.$$

1252 From the above equation, we have:  
 1253

$$1254 \quad f^* - \frac{L}{2} \sum_{k=1}^K p_k \Delta_k \leq \sum_{k=1}^K p_k f_k^*.$$

$$1258 \quad -\frac{1}{K} \sum_{k=1}^K f_k^* \leq -f^* + \frac{L}{2K} \sum_{k=1}^K \Delta_k.$$

1261 From the equation 14, we have the term:  
 1262

$$1263 \quad \frac{L^2 \eta^2 T (2+u)^2}{K} \sum_{r=0}^{R-1} \sum_{k=1}^K (f_k(w^r) - f_k^*).$$

1266 First, we express the double sum as:  
 1267

$$1268 \quad \sum_{r=0}^{R-1} \sum_{k=1}^K (f_k(w^r) - f_k^*) = \sum_{r=0}^{R-1} \left( \sum_{k=1}^K f_k(w^r) - \sum_{k=1}^K f_k^* \right).$$

1271 Since  $p_k = \frac{1}{K}$ , we have:  
 1272

$$1273 \quad \sum_{k=1}^K f_k(w^r) = K f(w^r),$$

1275 where  $f(w^r) = \sum_{k=1}^K p_k f_k(w^r) = \frac{1}{K} \sum_{k=1}^K f_k(w^r)$ . Therefore:  
 1276

$$1277 \quad \sum_{r=0}^{R-1} \sum_{k=1}^K (f_k(w^r) - f_k^*) = \sum_{r=0}^{R-1} \left( K f(w^r) - \sum_{k=1}^K f_k^* \right).$$

1280 From the earlier derivation, we have the inequality:  
 1281

$$1282 \quad -\frac{1}{K} \sum_{k=1}^K f_k^* \leq -f^* + \frac{L}{2K} \sum_{k=1}^K \Delta_k.$$

1284 Substituting this into the expression above:  
 1285

$$1286 \quad \sum_{r=0}^{R-1} \sum_{k=1}^K (f_k(w^r) - f_k^*) \leq \sum_{r=0}^{R-1} \left( K f(w^r) - \left( K f^* - \frac{L}{2} \sum_{k=1}^K \Delta_k \right) \right).$$

1289 Thus:  
 1290

$$1291 \quad \sum_{r=0}^{R-1} \sum_{k=1}^K (f_k(w^r) - f_k^*) \leq \sum_{r=0}^{R-1} \left( K f(w^r) - K f^* + \frac{L}{2} \sum_{k=1}^K \Delta_k \right).$$

1293 Since  $\Delta_k$  is constant across iterations, we can factor it out:  
 1294

$$1295 \quad K \sum_{r=0}^{R-1} (f(w^r) - f^*) + \frac{L}{2} \sum_{r=0}^{R-1} \sum_{k=1}^K \Delta_k = K \sum_{r=0}^{R-1} (f(w^r) - f^*) + \frac{LR}{2} \sum_{k=1}^K \Delta_k.$$

1296 Now, multiply by the coefficient:  
 1297

$$1299 \frac{L^2\eta^2T(2+u)^2}{K} \sum_{r=0}^{R-1} \sum_{k=1}^K (f_k(w^r) - f_k^*) \leq \frac{L^2\eta^2T(2+u)^2}{K} \left[ K \sum_{r=0}^{R-1} (f(w^r) - f^*) + \frac{LR}{2} \sum_{k=1}^K \Delta_k \right].$$

1302  
 1303 Simplifying:  
 1304

$$1305 \quad 1306 L^2\eta^2T(2+u)^2 \sum_{r=0}^{R-1} (f(w^r) - f^*) + \frac{L^3\eta^2T(2+u)^2R}{2K} \sum_{k=1}^K \Delta_k.$$

1309 Substituting this result into the original target inequality, we get:  
 1310

$$1312 \quad 1313 \left( \eta T \sqrt{c} (1 + \sqrt{c_h}) - \frac{1}{2\delta} \right) \sum_{r=0}^{R-1} \|\nabla f(w^r)\|^2 \leq \mathbb{E}_{\bar{z}} [f(w^0) - f(w^R)] + \frac{\eta^2 T^2 c \sigma_h^2 R \delta}{2} \\ 1314 \quad 1315 + L^2 \eta^2 T (2+u)^2 \sum_{r=0}^{R-1} (f(w^r) - f^*) \\ 1316 \quad 1317 + \frac{L^3 \eta^2 T (2+u)^2 R}{2K} \sum_{k=1}^K \Delta_k \\ 1318 \quad 1319 + \frac{L \eta^2 T^2 (2+u) c \sigma^2 R}{2}.$$

1324 According to the Assumption C.2 we have:  
 1325

$$1327 \quad 1328 2\mu(f(\mathbf{w}^r) - f^*) \leq \|\nabla f(\mathbf{w}^r)\|^2, \quad \forall \mathbf{w}^r \in \mathbb{R}^d,$$

$$1332 \quad 1333 2\mu \sum_{r=0}^{R-1} (f(\mathbf{w}^r) - f^*) \leq \sum_{r=0}^{R-1} \|\nabla f(\mathbf{w}^r)\|^2, \quad \forall \mathbf{w}^r \in \mathbb{R}^d,$$

1337 We let  $\eta T \sqrt{c} (1 + \sqrt{c_h}) - \frac{1}{2\delta} > 0$  and substitute the above inequality, we have:  
 1338

$$1340 \quad 1341 2\mu(\eta T \sqrt{c} (1 + \sqrt{c_h}) - \frac{1}{2\delta}) \sum_{r=0}^{R-1} (f(\mathbf{w}^r) - f^*) \leq \mathbb{E}_{\bar{z}} [f(w^0) - f(w^R)] + \frac{\eta^2 T^2 c \sigma_h^2 R \delta}{2} \\ 1342 \quad 1343 + L^2 \eta^2 T (2+u)^2 \sum_{r=0}^{R-1} (f(w^r) - f^*) \\ 1344 \quad 1345 + \frac{L^3 \eta^2 T (2+u)^2 R}{2K} \sum_{k=1}^K \Delta_k \\ 1346 \quad 1347 + \frac{L \eta^2 T^2 (2+u) c \sigma^2 R}{2}.$$

1350  
 1351  
 1352 
$$\sum_{r=0}^{R-1} (f(w^r) - f^*) \leq \frac{\mathbb{E}_{\bar{z}}[f(w^0) - f(w^R)]}{2\mu(\eta T \sqrt{c}(1 + \sqrt{c_h}) - \frac{1}{2\delta}) - L^2 \eta^2 T (2 + u)^2}$$
  
 1353  
 1354  
 1355  
 1356 
$$+ \frac{\eta^2 T^2 c \sigma_h^2 R \delta}{2 \left[ 2\mu(\eta T \sqrt{c}(1 + \sqrt{c_h}) - \frac{1}{2\delta}) - L^2 \eta^2 T (2 + u)^2 \right]}$$
  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372 
$$\frac{1}{R} \sum_{r=0}^{R-1} (f(w^r) - f^*) \leq \frac{1}{R} \frac{\mathbb{E}_{\bar{z}}[f(w^0) - f(w^R)]}{2\mu(\eta T \sqrt{c}(1 + \sqrt{c_h}) - \frac{1}{2\delta}) - L^2 \eta^2 T (2 + u)^2}$$
  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390 We select  $\delta = \frac{1}{\eta T \sqrt{c}(1 + \sqrt{c_h})}$ , which leads to:  
 1391  
 1392  
 1393 
$$\frac{1}{2\delta} = \frac{\eta T \sqrt{c}(1 + \sqrt{c_h})}{2}$$
  
 1394  
 1395  
 1396 Substituting into the denominator:  
 1397  
 1398  
 1399  
 1400 
$$2\mu \left( \eta T \sqrt{c}(1 + \sqrt{c_h}) - \frac{\eta T \sqrt{c}(1 + \sqrt{c_h})}{2} \right) = \mu \eta T \sqrt{c}(1 + \sqrt{c_h})$$
  
 1401  
 1402  
 1403

With the chosen  $\delta$ , we have:

$$\begin{aligned}
& \frac{1}{R} \sum_{r=0}^{R-1} (f(w^r) - f^*) \leq \frac{1}{R} \cdot \frac{\mathbb{E}_{\bar{z}} [f(w^0) - f(w^R)]}{\mu \eta T \sqrt{c(1 + \sqrt{c_h})} - L^2 \eta^2 T (2 + u)^2} \\
& \quad + \frac{\sqrt{c} \sigma_h^2}{2(1 + \sqrt{c_h}) [\mu \sqrt{c}(1 + \sqrt{c_h}) - L^2 \eta (2 + u)^2]} \\
& \quad + \frac{L^3 \eta (2 + u)^2 \sum_{k=1}^K \Delta_k}{2K [\mu \sqrt{c}(1 + \sqrt{c_h}) - L^2 \eta (2 + u)^2]} \\
& \quad + \frac{L \eta T (2 + u) c \sigma^2}{2 [\mu \sqrt{c}(1 + \sqrt{c_h}) - L^2 \eta (2 + u)^2]}, \tag{17}
\end{aligned}$$

where the step-size  $\eta$  must satisfy:  $\eta < \frac{\mu \sqrt{c}(1 + \sqrt{c_h})}{L^2 (2 + u)^2}$  to ensure denominator positivity.

Plugging in a constant learning rate  $\eta = \min \left\{ \frac{1}{L(u+2)}, \frac{\mu \sqrt{c}(1 + \sqrt{c_h})}{2L^2(2+u)^2} \right\}$ . We substitute this  $\eta$  to equation 17 and get:

$$\begin{aligned}
& \frac{1}{R} \sum_{r=0}^{R-1} (f(w^r) - f^*) \leq \frac{4L^2(2+u)^2}{\mu^2 c (1 + \sqrt{c_h})^2 T R} \mathbb{E}_{\bar{z}} [f(w^0) - f^*] \\
& \quad + \frac{\sigma_h^2}{\mu (1 + \sqrt{c_h})^2} + \frac{L}{K} \sum_{k=1}^K \Delta_k + \frac{T c \sigma^2}{2L (2 + u)}. \\
& \frac{1}{R} \sum_{r=0}^{R-1} (f(w^r) - f^*) \leq \mathcal{O} \left( \frac{(2+u)^2}{T R} \cdot \mathbb{E} [f(w^0) - f(w^R)] \right) + \mathcal{O} \left( \frac{T}{2+u} \right) + \mathcal{O}(1). \tag{18}
\end{aligned}$$

□

### C.5 MEERKAT-VP CONVERGENCE ANALYSIS

We propose a Virtual Path Client Selection (MEERKAT-VP) mechanism that identifies clients with highly heterogeneous data distributions based on their optimization trajectories. Instead of excluding them, MEERKAT-VP applies early stopping to these clients to limit their adverse influence on global model updates while still preserving their participation.

*Proof. Motivation for Early Stopping:* In federated learning, clients perform local updates starting from the global model  $w^r$ . For  $T > 1$ , clients may drift towards their local optima, introducing bias into the global update due to data heterogeneity. By identifying "bad" clients and limiting them to one update step, we reduce their drift and align their contributions more closely with the global gradient.

We divide the  $K$  clients into two groups:

- **Balanced-distribution clients ( $K_g$ ):** Perform  $T$  local step updates.
- **Skewed-distribution clients ( $K_b$ ):** Perform only 1 local step update.

The global model update becomes:

$$w^{r+1} = w^r + \frac{1}{K} \sum_{k \in K_g} (w_k^{r,T} - w^r) + \frac{1}{K} \sum_{k \in K_b} (w_k^{r,1} - w^r)$$

where:

$$w_k^{r,T} - w^r = -\eta \sum_{t=0}^{T-1} \hat{\nabla} f_k(w^{r,t}), \quad w_k^{r,1} - w^r = -\eta \hat{\nabla} f_k(w^r)$$

1458 **Loss Descent Analysis** Using the  $L$ -smoothness property:

$$1460 \quad f(w^{r+1}) \leq f(w^r) + \langle \nabla f(w^r), w^{r+1} - w^r \rangle + \frac{L}{2} \|w^{r+1} - w^r\|^2$$

1462 We analyze the inner product term:

$$1465 \quad \langle \nabla f(w^r), w^{r+1} - w^r \rangle = \sum_{k=1}^K p_k \langle \nabla f(\mathbf{w}^r), \mathbf{w}_k^{r,T} - \mathbf{w}^r \rangle$$

$$1470 \quad \sum_{k=1}^K p_k \langle \nabla f(\mathbf{w}^r), \mathbf{w}_k^{r,T} - \mathbf{w}^r \rangle = -\eta \sum_{k=1}^K p_k \langle \nabla f(\mathbf{w}^r), \sum_{t=0}^{T-1} \hat{\nabla} f_k(\mathbf{w}^{r,t}, \bar{\mathbf{z}}_t) \rangle$$

$$1473 \quad = -\eta \sum_{k=1}^K p_k \sum_{t=0}^{T-1} \langle \nabla f(\mathbf{w}^r), \hat{\nabla} f_k(\mathbf{w}^{r,t}, \bar{\mathbf{z}}_t) \rangle$$

1477 Since we have balanced-distribution clients and skewed-distribution clients:

$$1480 \quad \langle \nabla f(w^r), w^{r+1} - w^r \rangle = \frac{1}{K} \sum_{k \in K_g} \langle \nabla f(w^r), w_k^{r,T} - w^r \rangle + \frac{1}{K} \sum_{k \in K_b} \langle \nabla f(w^r), w_k^{r,1} - w^r \rangle$$

$$1484 \quad \langle \nabla f(w^r), w^{r+1} - w^r \rangle = -\frac{\eta}{K} \sum_{k \in K_g} \sum_{t=0}^{T-1} \langle \nabla f(w^r), \hat{\nabla} f_k(w^{r,t}) \rangle$$

$$1488 \quad - \frac{\eta}{K} \sum_{k \in K_b} \langle \nabla f(w^r), \hat{\nabla} f_k(w^r) \rangle \tag{19}$$

1490 Since  $\hat{\nabla} f_k^t$  is unbiased, we have:

$$1492 \quad \sum_{k=1}^K \sum_{t=0}^{T-1} \langle \nabla f(w^r), \hat{\nabla} f_k(\mathbf{w}^{r,t}, \bar{\mathbf{z}}_t) \rangle = \sum_{k=1}^K \sum_{t=0}^{T-1} \langle \nabla f(w^r), \mathbb{E}_{\bar{\mathbf{z}}}[\hat{\nabla} f_k(\mathbf{w}^{r,t}, \bar{\mathbf{z}}_t)] \rangle$$

1495 We substitute the equation 4 and get:

$$1497 \quad \sum_{k=1}^K \sum_{t=0}^{T-1} \langle \nabla f(w^r), \mathbb{E}_{\bar{\mathbf{z}}}[\hat{\nabla} f_k(\mathbf{w}^{r,t}, \bar{\mathbf{z}}_t)] \rangle = \sum_{k=1}^K \sum_{t=0}^{T-1} \langle \nabla f(w^r), \mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t}) \rangle.$$

1500 Thus taking expectation of equation 19, we can get:

$$1503 \quad \mathbb{E}_{\bar{\mathbf{z}}} \langle \nabla f(\mathbf{w}^r), \mathbf{w}^{r+1} - \mathbf{w}^r \rangle = -\frac{\eta}{K} \left( \sum_{k \in K_g} \sum_{t=0}^{T-1} \langle \nabla f(\mathbf{w}^r), \mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t}) \rangle \right. \\ 1505 \quad \left. + \sum_{k \in K_b} \langle \nabla f(\mathbf{w}^r), \mathbf{m} \odot \nabla f_k(\mathbf{w}^r) \rangle \right) \tag{20}$$

1509 Under the Cauchy–Schwarz inequality, we have:

$$1511 \quad \langle \nabla f(\mathbf{w}^r), \mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t}) \rangle \leq \|\nabla f(\mathbf{w}^r)\| \|\mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t})\|$$

1512 We substitute Assumption C.6 get:  
 1513

$$1514 \|\nabla f(\mathbf{w}^r)\| \|\mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t})\| = \sqrt{c} \|\nabla f(\mathbf{w}^r)\| \|\nabla f_k(\mathbf{w}^{r,t})\|.$$

1515 Thus we get:  
 1516

$$1517 \langle \nabla f(\mathbf{w}^r), \mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t}) \rangle \leq \sqrt{c} \|\nabla f(\mathbf{w}^r)\| \|\nabla f_k(\mathbf{w}^{r,t})\|.$$

1518 By the triangle inequality, we have  
 1519

$$1520 \|\nabla f_k(\mathbf{w}^{r,t})\| \leq \|\nabla f(\mathbf{w}^r)\| + \|\nabla f_k(\mathbf{w}^{r,t}) - \nabla f(\mathbf{w}^r)\|$$

1521 We substitute Assumption C.3 and use the properties of square roots we get:  
 1522

$$1523 \|\nabla f(\mathbf{w}^r)\| + \|\nabla f_k(\mathbf{w}^{r,t}) - \nabla f(\mathbf{w}^r)\| \leq \|\nabla f(\mathbf{w}^r)\| + \sqrt{c_h \|\nabla f(\mathbf{w}^r)\|^2 + \sigma_h^2}$$

$$1524 \leq (1 + \sqrt{c_h}) \|\nabla f(\mathbf{w}^r)\| + \sigma_h.$$

1525 Using the bound  $\langle \nabla f(\mathbf{w}^r), m \odot \nabla f_k(\mathbf{w}^{r,t}) \rangle \leq \sqrt{c} \|\nabla f(\mathbf{w}^r)\| \|\nabla f_k\|$  from Cauchy–Schwarz and  
 1526 Assumption C.6, and then plugging in the above, we obtain

$$1527 \langle \nabla f(\mathbf{w}^r), m \odot \nabla f_k(\mathbf{w}^{r,t}) \rangle \leq \sqrt{c} \|\nabla f(\mathbf{w}^r)\| [(1 + \sqrt{c_h}) \|\nabla f(\mathbf{w}^r)\| + \sigma_h]$$

$$1528 \leq \sqrt{c} (1 + \sqrt{c_h}) \|\nabla f(\mathbf{w}^r)\|^2 + \sqrt{c} \sigma_h \|\nabla f(\mathbf{w}^r)\|.$$

1529 Since this bound holds uniformly for all  $k$  and  $t$ , and based on the equation 20 we get:  
 1530

$$1531 \sum_{k \in K_g} \sum_{t=0}^{T-1} \langle \nabla f(\mathbf{w}^r), \mathbf{m} \odot \nabla f_k(\mathbf{w}^{r,t}) \rangle + \sum_{k \in K_b} \langle \nabla f(\mathbf{w}^r), \mathbf{m} \odot \nabla f_k(\mathbf{w}^r) \rangle$$

$$1534 \leq (|K_g|T + |K_b|) [\sqrt{c}(1 + \sqrt{c_h}) \|\nabla f(\mathbf{w}^r)\|^2 + \sqrt{c} \sigma_h \|\nabla f(\mathbf{w}^r)\|].$$

1535 We get:  
 1536

$$1537 \mathbb{E}_{\bar{z}}[f(w^{r+1})] \leq \mathbb{E}_{\bar{z}}[f(w^r)] - \frac{\eta \sqrt{c} \alpha}{K} (1 + \sqrt{c_h}) \|\nabla f(w^r)\|^2$$

$$1539 - \frac{\eta \sqrt{c} \alpha}{K} \sigma_h \|\nabla f(w^r)\| + \frac{L}{2} \mathbb{E}_{\bar{z}} \|w^{r+1} - w^r\|^2 \quad (21)$$

1541 where  $\alpha = |K_g|T + |K_b|$ .  
 1542

1543 Since the global model update is given by:  
 1544

$$1544 w^{r+1} = w^r + \frac{1}{K} \sum_{k \in K_g} (w_k^{r,T} - w^r) + \frac{1}{K} \sum_{k \in K_b} (w_k^{r,1} - w^r)$$

1546 We substitute the local updates and the squared norm is:  
 1547

$$1548 \|w^{r+1} - w^r\|^2 = \frac{\eta^2}{K^2} \left\| \sum_{k \in K_g} \sum_{t=0}^{T-1} \hat{\nabla} f_k(w^{r,t}) + \sum_{k \in K_b} \hat{\nabla} f_k(w^r) \right\|^2$$

1552 Define the update contribution per client:  
 1553

$$1554 \hat{\Delta}_k = \begin{cases} -\eta \sum_{t=0}^{T-1} \hat{\nabla} f_k(w^{r,t}) & \text{if } k \in K_g, \\ -\eta \hat{\nabla} f_k(w^r) & \text{if } k \in K_b. \end{cases}$$

1556 Then:  
 1557

$$1558 w^{r+1} - w^r = \frac{1}{K} \sum_{k=1}^K \hat{\Delta}_k$$

$$1560 \|w^{r+1} - w^r\|^2 = \frac{1}{K^2} \left\| \sum_{k=1}^K \hat{\Delta}_k \right\|^2$$

1562 Using the Cauchy-Schwarz inequality:  
 1563

$$1564 \left\| \sum_{k=1}^K \hat{\Delta}_k \right\|^2 \leq K \sum_{k=1}^K \|\hat{\Delta}_k\|^2, \quad \text{where } \hat{\Delta}_k \text{ denotes the actual model update on client } k.$$

1566

So:

1567

1568

1569

$$\|w^{r+1} - w^r\|^2 \leq \frac{1}{K} \sum_{k=1}^K \|\hat{\Delta}_k\|^2$$

1570

Now compute  $\|\hat{\Delta}_k\|^2$ :

1571

$$\|\hat{\Delta}_k\|^2 = \eta^2 \left\| \sum_{t=0}^{T-1} \hat{\nabla} f_k(w^{r,t}) \right\|^2 \quad \text{for } k \in K_g,$$

1572

1573

$$\|\hat{\Delta}_k\|^2 = \eta^2 \left\| \hat{\nabla} f_k(w^r) \right\|^2 \quad \text{for } k \in K_b.$$

1574

Thus:

1575

1576

1577

$$\|w^{r+1} - w^r\|^2 \leq \frac{\eta^2}{K} \left( \sum_{k \in K_g} \left\| \sum_{t=0}^{T-1} \hat{\nabla} f_k(w^{r,t}) \right\|^2 + \sum_{k \in K_b} \left\| \hat{\nabla} f_k(w^r) \right\|^2 \right)$$

1578

We take the expectation:

1579

1580

1581

1582

1583

1584

1585

1586

For  $k \in K_b$ :

1587

$$\mathbb{E}_{\bar{z}} \left\| \hat{\nabla} f_k(w^r) \right\|^2 = (2+u)c \left\| \nabla f_k(w^r) \right\|^2$$

1588

1589

For  $k \in K_g$ :

1590

1591

1592

1593

$$\mathbb{E}_{\bar{z}} \left\| \sum_{t=0}^{T-1} \hat{\nabla} f_k(w^{r,t}) \right\|^2$$

1594

Using the Cauchy-Schwarz inequality:

1595

1596

1597

1598

$$\mathbb{E}_{\bar{z}} \left\| \sum_{t=0}^{T-1} \hat{\nabla} f_k(w^{r,t}) \right\|^2 \leq T \sum_{t=0}^{T-1} \mathbb{E}_{\bar{z}} \left\| \hat{\nabla} f_k(w^{r,t}) \right\|^2$$

1599

According to the lemma C.7:

1600

1601

1602

So:

1603

1604

1605

$$\mathbb{E}_{\bar{z}} \left\| \sum_{t=0}^{T-1} \hat{\nabla} f_k(w^{r,t}) \right\|^2 \leq T(2+u)c \sum_{t=0}^{T-1} \left\| \nabla f_k(w^{r,t}) \right\|^2$$

1606

Combine the terms we get:

1607

1608

1609

1610

1611

$$\mathbb{E}_{\bar{z}} \|w^{r+1} - w^r\|^2 \leq \frac{\eta^2(2+u)c}{K} \left( T \sum_{k \in K_g} \sum_{t=0}^{T-1} \left\| \nabla f_k(w^{r,t}) \right\|^2 + \sum_{k \in K_b} \left\| \nabla f_k(w^r) \right\|^2 \right)$$

1612

We substitute this inequality to the equation 21.

1613

1614

1615

1616

1617

1618

1619

$$\begin{aligned} \mathbb{E}_{\bar{z}}[f(w^{r+1})] &\leq \mathbb{E}_{\bar{z}}[f(w^r)] - \frac{\eta\sqrt{c}\alpha}{K} (1 + \sqrt{c_h}) \left\| \nabla f(w^r) \right\|^2 \\ &\quad - \frac{\eta\sqrt{c}\alpha}{K} \sigma_h \left\| \nabla f(w^r) \right\|^2 \\ &\quad + \frac{\eta^2(2+u)cL}{2K} \left( T \sum_{k \in K_g} \sum_{t=0}^{T-1} \left\| \nabla f_k(w^{r,t}) \right\|^2 + \sum_{k \in K_b} \left\| \nabla f_k(w^r) \right\|^2 \right) \end{aligned} \tag{22}$$

$$\begin{aligned}
& \mathbb{E}_{\bar{z}}[f(w^{r+1})] \leq \mathbb{E}_{\bar{z}}[f(w^r)] - \frac{\eta\sqrt{c}\alpha}{K} (1 + \sqrt{c_h}) \|\nabla f(w^r)\|^2 - \frac{\eta\sqrt{c}\alpha}{K} \sigma_h \|\nabla f(w^r)\| \\
& + \frac{\eta^2(2+u)cLT}{2K} \sum_{k \in K_g} \sum_{t=0}^{T-1} \|\nabla f_k(w^{r,t})\|^2 + \frac{\eta^2(2+u)cL}{2K} \sum_{k \in K_b} \|\nabla f_k(w^r)\|^2
\end{aligned}$$

According to the equation 7, we know that the client-average squared gradient has upper bound.

$$\begin{aligned}
& \mathbb{E}_{\bar{z}}[f(w^{r+1})] \leq \mathbb{E}_{\bar{z}}[f(w^r)] - \frac{\eta\sqrt{c}\alpha}{K} (1 + \sqrt{c_h}) \|\nabla f(w^r)\|^2 - \frac{\eta\sqrt{c}\alpha}{K} \sigma_h \|\nabla f(w^r)\| \\
& + \frac{\eta^2(2+u)cLT}{2K} \sum_{k \in K_g} \left[ \frac{2L(2+u)}{c} (f_k(w_k^{0,r}) - f_k^*) + T\sigma^2 \right] \\
& + \frac{\eta^2(2+u)cL}{2K} \sum_{k \in K_b} \|\nabla f_k(w^r)\|^2.
\end{aligned} \tag{23}$$

Using Assumption C.3, which states that for any  $\theta \in \mathbb{R}^d$ ,

$$\|\nabla f(\theta) - \nabla f_i(\theta)\|^2 \leq c_h \|\nabla f(\theta)\|^2 + \sigma_h^2,$$

we can bound the squared norm of the local gradient  $\|\nabla f_k(w^r)\|^2$ . Specifically, by the inequality  $(x+y)^2 \leq 2x^2 + 2y^2$ , we have:

$$\|\nabla f_k(w^r)\|^2 = \|\nabla f(w^r) + (\nabla f_k(w^r) - \nabla f(w^r))\|^2 \leq 2 \|\nabla f(w^r)\|^2 + 2 \|\nabla f_k(w^r) - \nabla f(w^r)\|^2.$$

Then, applying Assumption C.3 with  $\theta = w^r$  and  $i = k$ :

$$\|\nabla f_k(w^r) - \nabla f(w^r)\|^2 \leq c_h \|\nabla f(w^r)\|^2 + \sigma_h^2.$$

Therefore,

$$\|\nabla f_k(w^r)\|^2 \leq 2 \|\nabla f(w^r)\|^2 + 2 \left( c_h \|\nabla f(w^r)\|^2 + \sigma_h^2 \right) = (2 + 2c_h) \|\nabla f(w^r)\|^2 + 2\sigma_h^2.$$

Thus, we obtain the bound:

$$\|\nabla f_k(w^r)\|^2 \leq (2 + 2c_h) \|\nabla f(w^r)\|^2 + 2\sigma_h^2.$$

We substitute the bound to the inequality 23, according to the Assumption C.3, we substitute the last term:

$$\begin{aligned}
& \mathbb{E}_{\bar{z}}[f(w^{r+1})] \leq \mathbb{E}_{\bar{z}}[f(w^r)] - \frac{\eta\sqrt{c}\alpha}{K} (1 + \sqrt{c_h}) \|\nabla f(w^r)\|^2 - \frac{\eta\sqrt{c}\alpha}{K} \sigma_h \|\nabla f(w^r)\| \\
& + \frac{\eta^2(2+u)cLT}{2K} \sum_{k \in K_g} \left[ \frac{2L(2+u)}{c} (f_k(w_k^{0,r}) - f_k^*) + T\sigma^2 \right] \\
& + \frac{\eta^2(2+u)cL}{2K} \sum_{k \in K_b} \left[ (2 + 2c_h) \|\nabla f(w^r)\|^2 + \sigma_h^2 \right]. \\
& \mathbb{E}_{\bar{z}}[f(w^{r+1})] \leq \mathbb{E}_{\bar{z}}[f(w^r)] - \frac{\eta\sqrt{c}\alpha}{K} (1 + \sqrt{c_h}) \|\nabla f(w^r)\|^2 - \frac{\eta\sqrt{c}\alpha}{K} \sigma_h \|\nabla f(w^r)\| \\
& + \frac{\eta^2(2+u)cLT}{2K} \sum_{k \in K_g} \left[ \frac{2L(2+u)}{c} (f_k(w_k^{0,r}) - f_k^*) + T\sigma^2 \right] \\
& + \frac{\eta^2(2+u)cL|K_b|(1+c_h)}{K} \|\nabla f(w^r)\|^2 + \frac{\eta^2(2+u)cL|K_b|\sigma_h^2}{K}.
\end{aligned}$$

$$\begin{aligned}
& \mathbb{E}_{\bar{z}}[f(w^{r+1})] \leq \mathbb{E}_{\bar{z}}[f(w^r)] \\
& + \frac{\eta^2(2+u)cLK_b(1+c_h) - \eta\sqrt{c}\alpha\|\nabla f(w^r)\|^2}{K} \\
& - \frac{\eta\sqrt{c}\alpha\sigma_h}{K}\|\nabla f(w^r)\| \\
& + \frac{\eta^2(2+u)^2L^2T}{K} \sum_{k \in K_g} (f_k(w_k^{0,r}) - f_k^*) \\
& + \frac{\eta^2(2+u)cL}{2K} \left( T^2\sigma^2K_g + 2K_b\sigma_h^2 \right).
\end{aligned} \tag{24}$$

**Accumulating Over  $R$  Rounds.** Summing equation 24 over  $r = 0$  to  $R - 1$ ,

$$\begin{aligned}
& \mathbb{E}_{\bar{z}}[f(w^R)] - \mathbb{E}_{\bar{z}}[f(w^0)] \leq \frac{\eta^2(2+u)cLK_b(1+c_h) - \eta\sqrt{c}\alpha\sum_{r=0}^{R-1}\|\nabla f(w^r)\|^2}{K} \\
& - \frac{\eta\sqrt{c}\alpha\sigma_h}{K}\sum_{r=0}^{R-1}\|\nabla f(w^r)\| \\
& + \frac{\eta^2(2+u)^2L^2T}{K} \sum_{r=0}^{R-1} \sum_{k \in K_g} (f_k(w_k^{0,r}) - f_k^*) \\
& + \frac{\eta^2(2+u)cLR}{2K} \left( T^2\sigma^2K_g + 2K_b\sigma^2 \right).
\end{aligned} \tag{25}$$

According to our previous derivation, we know that:

$$\sum_{r=0}^{R-1} \|\nabla f(w^r)\| \leq \sqrt{R} \sqrt{\sum_{r=0}^{R-1} \|\nabla f(w^r)\|^2}. \tag{26}$$

Apply Young's inequality with  $\delta > 0$  and nonnegative real numbers  $x$  and  $y$ ,

$$xy \leq \frac{x^2}{2\delta} + \frac{y^2\delta}{2}.$$

$$\begin{aligned}
& \frac{\eta\sqrt{c}\sigma_h\alpha}{K} \sum_{r=0}^R \|\nabla f(w^r)\| \leq \frac{\eta\sqrt{c}\alpha\sigma_h\sqrt{R}}{K} \sqrt{\sum_{r=0}^R \|\nabla f(w^r)\|^2} \\
& \leq \frac{1}{2\delta} \sum_{r=0}^R \|\nabla f(w^r)\|^2 + \frac{\eta^2 c \alpha^2 \sigma_h^2 R \delta}{2K^2} \\
& - \frac{\eta\sqrt{c}\sigma_h\alpha}{K} \sum_{r=0}^R \|\nabla f(w^r)\| \leq \frac{1}{2\delta} \sum_{r=0}^R \|\nabla f(w^r)\|^2 + \frac{\eta^2 c \alpha^2 \sigma_h^2 R \delta}{2K^2}.
\end{aligned}$$

We substitute this to the equation 25.

$$\begin{aligned}
& \mathbb{E}_{\bar{z}}[f(w^R)] - \mathbb{E}_{\bar{z}}[f(w^0)] \leq \left( \frac{\eta^2(2+u)cLK_b(1+c_h)}{K} - \frac{\eta\sqrt{c}\alpha}{K} + \frac{1}{2\delta} \right) \sum_{r=0}^{R-1} \|\nabla f(w^r)\|^2 \\
& \quad + \frac{\eta^2 c \alpha^2 \sigma_h^2 R \delta}{2 K^2} + \frac{\eta^2(2+u)^2 L^2 T}{K} \sum_{r=0}^R \sum_{k \in K_g} (f_k(w_k^{0,r}) - f_k^*) \\
& \quad + \frac{\eta^2(2+u)cL R}{2K} (T^2 \sigma^2 K_g + 2 K_b \sigma_h^2).
\end{aligned} \tag{27}$$

Given that  $w_k^{0,r} = w^r$ , this term is equivalent to  $\sum_{r=0}^R \sum_{k \in K_g} (f_k(w^r) - f_k^*)$ .

From our previous discussion, we have the inequality for a single round  $r$ :

$$\sum_{k \in K_g} (f_k(w^r) - f_k^*) \leq \sum_{k=1}^K (f_k(w^r) - f_k^*)$$

and the inequality used in Part 2 of the proof:

$$\sum_{r=0}^{R-1} \sum_{k=1}^K (f_k(w^r) - f_k^*) \leq \sum_{r=0}^{R-1} \left( K f(w^r) - K f^* + \frac{L}{2} \sum_{k=1}^K \Delta_k \right).$$

Combining these two inequalities, we obtain a bound for the sum over the set  $K_g$ :

We set  $\gamma \leq 1$  which means that the subset clients the effect to the global:

$$\sum_{r=0}^{R-1} \sum_{k \in K_g} (f_k(w^r) - f_k^*) \leq \gamma \sum_{r=0}^{R-1} \left( K (f(w^r) - f^*) + \frac{L}{2} \sum_{k=1}^K \Delta_k \right)$$

We substitute this to the above inequality get:

$$\begin{aligned}
& \mathbb{E}_{\bar{z}}[f(w^R)] - \mathbb{E}_{\bar{z}}[f(w^0)] \leq \left( \frac{\eta^2(2+u)cLK_b(1+c_h)}{K} - \frac{\eta\sqrt{c}\alpha}{K} + \frac{1}{2\delta} \right) \sum_{r=0}^{R-1} \|\nabla f(w^r)\|^2 \\
& \quad + \frac{\eta^2 c \alpha^2 \sigma_h^2 R \delta}{2 K^2} + \eta^2(2+u)^2 L^2 T \gamma \sum_{r=0}^{R-1} (f(w^r) - f^*) \\
& \quad + \frac{\eta^2(2+u)^2 L^3 T R \gamma}{2K} \sum_{k=1}^K \Delta_k \\
& \quad + \frac{\eta^2(2+u)cL R}{2K} (T^2 \sigma^2 K_g + 2 K_b \sigma_h^2).
\end{aligned}$$

We substitute  $\alpha$ :

$$\begin{aligned}
& \mathbb{E}_{\bar{z}}[f(w^R)] - \mathbb{E}_{\bar{z}}[f(w^0)] \leq \left( \frac{\eta^2(2+u)cLK_b(1+c_h)}{K} - \frac{\eta\sqrt{c}(K_g T + K_b)}{K} + \frac{1}{2\delta} \right) \sum_{r=0}^{R-1} \|\nabla f(w^r)\|^2 \\
& \quad + \frac{\eta^2 c (K_g^2 T^2 + 2K_g T K_b + K_b^2) \sigma_h^2 R \delta}{2 K^2} + \eta^2(2+u)^2 L^2 T \gamma \sum_{r=0}^{R-1} (f(w^r) - f^*) \\
& \quad + \frac{\eta^2(2+u)^2 L^3 T R \gamma}{2K} \sum_{k=1}^K \Delta_k + \frac{\eta^2(2+u)cL R}{2K} (T^2 \sigma^2 K_g + 2 K_b \sigma_h^2).
\end{aligned}$$

1782 To simplify the inequality, we solve for  $\delta$ :  
 1783

$$\begin{aligned} \frac{1}{2\delta} &= -\frac{\eta\sqrt{c}(K_gT + K_b)}{2K} - \frac{\eta^2(2+u)cLK_b(1+c_h)}{K} + \frac{\eta\sqrt{c}(K_gT + K_b)}{K}, \\ \frac{1}{2\delta} &= \frac{\eta\sqrt{c}(K_gT + K_b)}{2K} - \frac{\eta^2(2+u)cLK_b(1+c_h)}{K}, \\ \delta &= \frac{K}{\eta\sqrt{c}(K_gT + K_b) - 2\eta^2(2+u)cLK_b(1+c_h)}. \end{aligned}$$

1791 For  $\delta > 0$ , the denominator must be positive:  
 1792

$$\eta\sqrt{c}(K_gT + K_b) - 2\eta^2(2+u)cLK_b(1+c_h) > 0,$$

1793 yielding the condition:  
 1794

$$\eta < \frac{\sqrt{c}(K_gT + K_b)}{2(2+u)cLK_b(1+c_h)}.$$

1795 Substitute  $\delta$ :  
 1796

$$\begin{aligned} \mathbb{E}_{\bar{z}}[f(w^R)] - \mathbb{E}_{\bar{z}}[f(w^0)] &\leq -\frac{\eta\sqrt{c}(K_gT + K_b)}{2K} \sum_{r=0}^{R-1} \|\nabla f(w^r)\|^2 \\ &\quad + \frac{\eta^2 c(K_gT + K_b)^2 \sigma_h^2 R}{2K (\eta\sqrt{c}(K_gT + K_b) - 2\eta^2(2+u)cLK_b(1+c_h))} \\ &\quad + \eta^2(2+u)^2 L^2 T \gamma \sum_{r=0}^{R-1} (f(w^r) - f^*) \\ &\quad + \frac{\eta^2(2+u)^2 L^3 T R \gamma}{2K} \sum_{k=1}^K \Delta_k \\ &\quad + \frac{\eta^2(2+u)cLR}{2K} (T^2 \sigma^2 K_g + 2K_b \sigma_h^2). \end{aligned}$$

1813 According to the Assumption C.2 we have:  
 1814

$$2\mu(f(\mathbf{w}^r) - f^*) \leq \|\nabla f(\mathbf{w}^r)\|^2, \quad \forall \mathbf{w}^r \in \mathbb{R}^d,$$

$$2\mu \sum_{r=0}^{R-1} (f(\mathbf{w}^r) - f^*) \leq \sum_{r=0}^{R-1} \|\nabla f(\mathbf{w}^r)\|^2, \quad \forall \mathbf{w}^r \in \mathbb{R}^d,$$

1821 Combine the PL inequality to the above function we get:  
 1822

$$\begin{aligned} \frac{\eta\sqrt{c}(K_gT + K_b)}{2K} \sum_{r=0}^{R-1} \|\nabla f(w^r)\|^2 &\leq \mathbb{E}_{\bar{z}}[f(w^0)] - \mathbb{E}_{\bar{z}}[f(w^R)] \\ &\quad + \frac{\eta^2 c(K_gT + K_b)^2 \sigma_h^2 R}{2K (\eta\sqrt{c}(K_gT + K_b) - 2\eta^2(2+u)cLK_b(1+c_h))} \\ &\quad + \eta^2(2+u)^2 L^2 T \gamma \sum_{r=0}^{R-1} (f(w^r) - f^*) \\ &\quad + \frac{\eta^2(2+u)^2 L^3 T R \gamma}{2K} \sum_{k=1}^K \Delta_k \\ &\quad + \frac{\eta^2(2+u)cLR}{2K} (T^2 \sigma^2 K_g + 2K_b \sigma_h^2). \end{aligned}$$

1836 Let  $S_E = \sum_{r=0}^{R-1} \mathbb{E}[f(w^r) - f^*]$ ,  $D_\delta = \eta\sqrt{c}(K_g T + K_b) - 2\eta^2(2+u)cLK_b(1+c_h)$ . We require  
 1837  $D_\delta > 0$ .

1838 Substituting this back into the original inequality:

$$\begin{aligned} 1840 \mathbb{E}[f(w^R)] - \mathbb{E}[f(w^0)] &\leq -\frac{\eta\mu\sqrt{c}(K_g T + K_b)}{K} S_E + \eta^2(2+u)^2 L^2 T \gamma S_E \\ 1841 &\quad + \frac{\eta^2 c(K_g T + K_b)^2 \sigma_h^2 R}{2K D_\delta} \\ 1842 &\quad + \frac{\eta^2(2+u)^2 L^3 T R \gamma}{2K} \sum_{k=1}^K \Delta_k \\ 1843 &\quad + \frac{\eta^2(2+u)cLR}{2K} (T^2 \sigma^2 K_g + 2K_b \sigma_h^2). \\ 1844 \\ 1845 \\ 1846 \\ 1847 \\ 1848 \\ 1849 \end{aligned}$$

1850 Collecting terms involving  $S_E$ :

$$\begin{aligned} 1851 \mathbb{E}[f(w^R)] - \mathbb{E}[f(w^0)] &\leq \left( \eta^2(2+u)^2 L^2 T \gamma - \frac{\eta\mu\sqrt{c}(K_g T + K_b)}{K} \right) S_E + \text{other terms}. \\ 1852 \\ 1853 \end{aligned}$$

1854 Moving  $S_E$  to the left side:

$$\begin{aligned} 1855 \left( \frac{\eta\mu\sqrt{c}(K_g T + K_b)}{K} - \eta^2(2+u)^2 L^2 T \gamma \right) S_E &\leq \mathbb{E}[f(w^0)] - \mathbb{E}[f(w^R)] \\ 1856 &\quad + \frac{\eta^2 c(K_g T + K_b)^2 \sigma_h^2 R}{2K D_\delta} + \frac{\eta^2(2+u)^2 L^3 T R \gamma}{2K} \sum_{k=1}^K \Delta_k \\ 1857 \\ 1858 \\ 1859 \\ 1860 \\ 1861 \\ 1862 \\ 1863 \\ 1864 \\ 1865 \\ 1866 \\ 1867 \end{aligned} \sum_{k=1}^K \Delta_k + \frac{\eta^2(2+u)cLR}{2K} (T^2 \sigma^2 K_g + 2K_b \sigma_h^2). \quad (28)$$

1863 Since  $\mathbb{E}[f(w^R)] \geq f^*$  (typically  $f^*$  is the minimum), we have  $\mathbb{E}[f(w^0)] - \mathbb{E}[f(w^R)] \leq \mathbb{E}[f(w^0) - f^*]$ . Let  $f_0^* = \mathbb{E}[f(w^0) - f^*]$  (the initial expected suboptimality). Let the coefficient of  $S_E$  be  
 1864  $C'_S = \frac{\eta\mu\sqrt{c}(K_g T + K_b)}{K} - \eta^2(2+u)^2 L^2 T \gamma$ . To ensure  $C'_S > 0$ , we need  $\eta$  sufficiently small such that  
 1865  $\eta < \frac{\mu\sqrt{c}(K_g T + K_b)}{K(2+u)^2 L^2 T \gamma}$ . Then:

$$\begin{aligned} 1868 C'_S S_E &\leq f_0^* + \frac{\eta^2 c(K_g T + K_b)^2 \sigma_h^2 R}{2K D_\delta} + \frac{\eta^2(2+u)^2 L^3 T R \gamma}{2K} \sum_{k=1}^K \Delta_k \\ 1869 \\ 1870 \\ 1871 \\ 1872 \\ 1873 \end{aligned} + \frac{\eta^2(2+u)cLR}{2K} (T^2 \sigma^2 K_g + 2K_b \sigma_h^2).$$

1874 Our goal is  $\frac{1}{R} S_E = \frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E}[f(w^r) - f^*]$ . Dividing both sides by  $R$ :

$$\begin{aligned} 1876 C'_S \frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E}[f(w^r) - f^*] &\leq \frac{f_0^*}{R} + \frac{\eta^2 c(K_g T + K_b)^2 \sigma_h^2}{2K D_\delta} \\ 1877 &\quad + \frac{\eta^2(2+u)^2 L^3 T \gamma}{2K} \sum_{k=1}^K \Delta_k \\ 1878 \\ 1879 \\ 1880 \\ 1881 \\ 1882 \\ 1883 \\ 1884 \end{aligned} + \frac{\eta^2(2+u)cL}{2K} (T^2 \sigma^2 K_g + 2K_b \sigma_h^2).$$

Finally, dividing both sides by  $C'_S$  (assuming  $C'_S > 0$ ):

$$\begin{aligned} 1885 \frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E}_{\bar{z}}[f(w^r) - f^*] &\leq \frac{1}{C'_S} \left[ \frac{f_0^*}{R} + \frac{\eta^2 c(K_g T + K_b)^2 \sigma_h^2}{2K(\eta\sqrt{c}(K_g T + K_b) - 2\eta^2(2+u)cLK_b(1+c_h))} \right. \\ 1886 &\quad \left. + \frac{\eta^2(2+u)^2 L^3 T \gamma}{2K} \sum_{k=1}^K \Delta_k + \frac{\eta^2(2+u)cL}{2K} (T^2 \sigma^2 K_g + 2K_b \sigma_h^2) \right], \\ 1887 \\ 1888 \\ 1889 \end{aligned}$$

1890

where

1891

$$1892 \quad C'_S = \frac{\eta \mu \sqrt{c} (K_g T + K_b)}{K} - \eta^2 (2+u)^2 L^2 T \gamma, \quad D_\delta = \eta \sqrt{c} (K_g T + K_b) - 2\eta^2 (2+u) c L K_b (1+c_h).$$

1893

To ensure both  $C'_S > 0$  and  $D_\delta > 0$ , we require

1894

1895

1896

1897

1898

1899

Let

1900

$$\eta_{\max} = \min\{\bar{\eta}_\delta, \bar{\eta}_S\}, \quad \theta \in (0, \frac{1}{2}].$$

1901

Choosing  $\theta = \frac{1}{2}$  gives

1902

1903

1904

$$\eta = \frac{1}{2} \eta_{\max} = \frac{\mu \sqrt{c} (K_g T + K_b)}{2 K (2+u)^2 L^2 T \gamma}.$$

We select

1905

1906

1907

And from previous client convergence conclusion, we pick a constant local learning rate

1908

1909

1910

$$\eta_{\text{client}} = \frac{c}{\alpha} = \frac{1}{L(u+2)} < \frac{2c}{\alpha}$$

1911

1912

1913

1914

Substituting the learning rate  $\eta = \min\left\{\frac{1}{L(u+2)}, \frac{\mu \sqrt{c} (K_g T + K_b)}{2 K (2+u)^2 L^2 T \gamma}\right\}$ , since  $\eta$  is a small value, we neglect  $\eta^2$ .

1915

1916

1917

$$\begin{aligned} R \sum_{r=0}^{R-1} \mathbb{E}_{\bar{z}}[f(w^r) - f^*] &\leq \frac{4K^2(2+u)^2 L^2 T \gamma \mathbb{E}[f(w^0) - f^*]}{\mu^2 c (K_g T + K_b)^2 R} + \frac{\sigma_h^2}{2} + \frac{(2+u) L}{4K} \sum_{k=1}^K \Delta_k \\ &\quad + \frac{c}{4K(2+u)L T \gamma} (T^2 \sigma^2 K_g + 2K_b \sigma_h^2). \end{aligned}$$

1918

1919

1920

1921

1922

1923

1924

1925

1926

1927

1928

1929

1930

1931

1932

1933

1934

$$\begin{aligned} \frac{1}{R} \sum_{r=0}^{R-1} \mathbb{E}_{\bar{z}}[f(w^r) - f^*] &\leq O\left(\frac{K^2 (2+u)^2 \gamma T}{c (K_g T + K_b)^2 R}\right) \\ &\quad + O\left(\frac{1+u}{K} \left(\sum_{k=1}^{K_g} \Delta_{kg} + \sum_{k=1}^{K_b} \Delta_{kb}\right)\right) \\ &\quad + O\left(\frac{c T K_g}{K(1+u) \gamma}\right) \\ &\quad + O\left(\frac{c K_b \sigma_h^2}{K(1+u) T \gamma}\right) + O(1). \end{aligned} \tag{29}$$

□

1935

Define the error upper-bounds for MEERKAT-VP and the baseline MEERKAT as follows:

1936

1937

1938

1939

1940

1941

1942

1943

$$\begin{aligned} E_{\text{MEERKAT-VP}} &= \underbrace{\frac{4K^2(2+u)^2 L^2 T \gamma}{\mu^2 c (K_g T + K_b)^2} \frac{\mathbb{E}[f(w^0) - f^*]}{R} + \left[ \frac{\sigma_h^2}{2} + \frac{(2+u)L}{4K} \sum_{k=1}^K \Delta_k + \frac{c(T^2 \sigma^2 K_g + 2K_b \sigma_h^2)}{4K(2+u)L T \gamma} \right]}_{\text{(I) Transient term}} \\ &\quad + \underbrace{\left[ \frac{\sigma_h^2}{2} + \frac{(2+u)L}{4K} \sum_{k=1}^K \Delta_k + \frac{c(T^2 \sigma^2 K_g + 2K_b \sigma_h^2)}{4K(2+u)L T \gamma} \right]}_{\text{(II) Steady-state term}}, \\ E_{\text{MEERKAT}} &= \underbrace{\frac{4L^2(2+u)^2}{\mu^2 c (1 + \sqrt{c_h})^2 T} \frac{\mathbb{E}[f(w^0) - f^*]}{R} + \left[ \frac{\sigma_h^2}{2} + \frac{(2+u)L}{4K} \sum_{k=1}^K \Delta_k + \frac{c(T^2 \sigma^2 K_g + 2K_b \sigma_h^2)}{4K(2+u)L T \gamma} \right]}_{\text{(I') Transient term}} \\ &\quad + \underbrace{\left[ \frac{\sigma_h^2}{2} + \frac{(2+u)L}{4K} \sum_{k=1}^K \Delta_k + \frac{c(T^2 \sigma^2 K_g + 2K_b \sigma_h^2)}{4K(2+u)L T \gamma} \right]}_{\text{(II') Steady-state term}}. \end{aligned}$$

1944

- **Transient term ratio:**

1945

1946

1947

$$\frac{(I)}{(I')} \approx \gamma(1 + \sqrt{c_h})^2 < 1, \quad \text{and as } c_h \rightarrow 1, \gamma(1 + \sqrt{c_h})^2 \rightarrow 0.$$

1948

- **Noise term ratio:**

1949

1950

1951

1952

1953

1954

$$\frac{\sigma_h^2/2}{\sigma_h^2/(\mu(1 + \sqrt{c_h})^2)} = \frac{\mu(1 + \sqrt{c_h})^2}{2}, \quad \text{which is } < 1 \text{ when } \mu(1 + \sqrt{c_h})^2 < 2.$$

Empirically  $\mu < 1$ , thus  $\mu(1 + \sqrt{c_h})^2 < 2$  is True. Additionally, VPCS includes an extra term  $\frac{cK_b\sigma_h^2}{2K(2+u)LT\gamma}$ , which decays as  $\frac{1}{T}$  and becomes negligible for large  $T$ .

1955

- **Heterogeneity and variance terms:**

1956

1957

1958

$$\frac{(2+u)L}{4K} \sum_{k=1}^K \Delta_k < \frac{L}{K} \sum_{k=1}^K \Delta_k, \quad \text{and the extra variance term decays as } 1/K.$$

1959

1960

Therefore, under the same  $T$  and  $R$ ,  $E_{\text{MEERKAT-VP}} < E_{\text{MEERKAT}}$  and this gap widens as data heterogeneity  $c_h$  increases.

1961

1962

## REMARKS

1963

1964

1965

1966

The analysis of the upper bound in Equation 17 reveals how the local training step  $T$ , density level  $u$ , and communication rounds  $R$  collectively influence the optimization dynamics through a balance of convergence rate, bias-variance trade-offs, and steady-state error control:

1967

1968

1969

1970

1971

1972

1973

1974

1975

1976

1977

1978

1979

1980

1981

- **Impact of Local Update Steps  $T$ :** A smaller  $T$  amplifies the term  $\mathcal{O}\left(\frac{(2+u)^2}{TR} \cdot \mathbb{E}[f(w^0) - f(w^R)]\right)$ , increasing the average optimality gap after  $R$  communication rounds when  $R$  is fixed. However, this effect can be mitigated by increasing  $R$ , as the scaling factor  $\frac{1}{R}$  reduces the term's impact. Conversely, reducing  $T$  diminishes the variance term  $\mathcal{O}\left(\frac{T}{2+u}\right)$ , leading to a smaller steady-state error. Thus, a smaller  $T$  may prolong the transient phase but ultimately achieves a tighter optimality gap relative to  $f^*$  after sufficient rounds.

- **Density Level  $u$ .** Reducing  $u$  (i.e., increasing sparsity) quadratically benefits the transient term, yet it also inflates the steady-state term through the denominator  $2+u$ . Choosing  $u$  therefore amounts to balancing communication savings against the plateau error; aggressive sparsification should be coupled with smaller  $T$  to avoid performance degradation.

- **MEERKAT-VP Client Selection Strategy:** By early-stopping extreme data-imbalance clients with a single local training step, MEERKAT-VP effectively reduces Non-IID drift in zeroth-order federated learning fine-tuning. This strategy lowers the coefficient of the transient term and further reduces heterogeneity- and variance-induced steady-state error. Under fixed  $T$  and  $R$ , these effects yield strictly faster convergence and a tighter optimality gap in Non-IID settings.

1982

1983

1984

These conclusions illustrate how tuning  $T$ ,  $R$ ,  $u$ , and the MEERKAT-VP client selection strategy can optimize performance in federated, sparse, and Non-IID learning scenarios.

1985

1986

## C.6 EMPIRICAL ANALYSIS OF THE GRADIP PHENOMENON

1987

1988

1989

1990

By Lemma C.8, the masked sparse zeroth-order (ZO) surrogate gradient is an *unbiased* estimator of the masked first-order gradient. Building on this fact, we define the vector  $g_c(w; x, y)$  is obtained by computing the gradient of the cross-entropy loss for a single sample with respect to a small subset of parameters selected by a mask.

1991

1992

1993

From logits to Softmax Probabilities we have:

1994

- The model's final layer outputs a *logit* for each class:

1995

1996

1997

$$h(x; w) = (h_1, \dots, h_C) \in \mathbb{R}^C.$$

- The softmax probabilities are given by:

$$p_j(x; w) = \frac{e^{h_j}}{\sum_{r=1}^C e^{h_r}}.$$

1998 The cross-entropy loss for a single sample is:  
 1999

2000  $\ell(w; x, y) = -\log p_y(x; w), \quad \text{where } y \in \{1, \dots, C\}.$   
 2001

2002 For each logit  $h_j$ , the partial derivative is:  
 2003

2004  $\frac{\partial \ell}{\partial h_j} = p_j - \mathbf{1}_{\{y=j\}} = p_j - (e_y)_j,$   
 2005

2006 where  $e_y$  is the one-hot vector with 1 in the  $y$ -th component.  
 2007

2008 Since we are only interested in the sensitive parameters selected by the mask  $m$ , the gradient with  
 2009 respect to the parameters can be written as:  
 2010

2011 
$$\begin{aligned} g_c(w; x, y) &= \nabla_{w_m} \ell(w; x, y) \\ 2012 &= \sum_{j=1}^C \frac{\partial \ell}{\partial h_j} \nabla_{w_m} h_j(x; w) \\ 2013 &= (p(x; w) - e_y)^\top \nabla_{w_m} h(x; w). \end{aligned}$$
  
 2014

2015 Here:  
 2016

2017 •  $\nabla_{w_m} h_j(x; w)$  is the gradient/Jacobian of the logit  $h_j$  with respect to the masked parameter  $w_m$ .  
 2018 • By collecting the coefficients  $p_j - \mathbf{1}_{y=j}$  into a vector, we obtain the compact form:  
 2019

2020 
$$g_c(w; x, y) = (p - e_y)^\top \nabla_{w_m} h(x; w).$$
  
 2021

2022 In our existing local client convergence inequality and from the assumption C.4, we can empirically  
 2023 write the key constant estimator variance:  
 2024

2025 
$$\sigma_k^2 = \frac{1}{d} \text{Var}_{(x,y) \sim D_k} [g_c(w; x, y)].$$
  
 2026

2027 We write  $g_c$  in matrix form: Define:  
 2028

2029 
$$\mathbf{J}(x; w) = \nabla_{w_m} h(x; w) \in \mathbb{R}^{d_m \times C}, \quad \mathbf{a}(x, y; w) = p(x; w) - e_y \in \mathbb{R}^C.$$
  
 2030

2031 Thus:  
 2032

2033 
$$g_c(w; x, y) = \mathbf{J}^\top(x; w) \mathbf{a}(x, y; w) \in \mathbb{R}^{d_m}.$$
  
 2034

2035 We substitute this equation to the above estimator variance:  
 2036

2037 
$$\sigma_k^2 = \frac{1}{d_m} \underbrace{\mathbb{E}_{(x,y)} \|g(w; x, y) - \nabla f_k(w)\|^2}_{\text{total variance}} = \frac{1}{d_m} \text{tr} \left( \mathbf{J}^\top \underbrace{\text{Cov}_{(x,y)} [\mathbf{a}(x, y; w)]}_{\Sigma_a} \mathbf{J} \right). \quad (1)$$
  
 2038

2039 Note:  
 2040

2041 •  $\Sigma_a \in \mathbb{R}^{C \times C}$  is determined solely by the **label distribution and prediction probabilities**.  
 2042 •  $\mathbf{J}$  reflects the network structure and influences only a similarity coefficient.  
 2043

2044 **Analysis of Extreme Non-IID (Single Label  $y^\dagger$ ):**  
 2045

2046 • The label is fixed, so  $\mathbf{1}_{y=j}$  is constant.  
 2047 • If the model is mostly correct:  $p \approx e_{y^\dagger}$ , then  $\mathbf{a}(x, y; w) \approx \mathbf{0}$ , yielding:  
 2048

2049 
$$\Sigma_a \approx \mathbf{0} \implies \sigma_{\text{non}}^2 \approx \frac{1}{d_m} \text{tr}(\mathbf{0}) = 0.$$
  
 2050

2051 **Analysis of Approximate IID (Balanced Multi-Label)**  
 2052

2053 • The label  $y$  varies across  $\{1, \dots, C\}$ .  
 2054

2052 • Even as the loss decreases,  $p_j$  differs across classes. The covariance is:  
 2053  
 2054  $(\Sigma_a)_{rs} = \mathbb{E}[(p_r - \mathbf{1}_{y=r})(p_s - \mathbf{1}_{y=s})] - \underbrace{(\mathbb{E}[p_r - \mathbf{1}_{y=r}])}_{=0} \underbrace{(\mathbb{E}[p_s - \mathbf{1}_{y=s}])}_{=0}.$   
 2055

2056 This matrix has diagonal elements  $\mathbb{E}[(p_r - \mathbf{1}_{y=r})^2] > 0$ , making  $\Sigma_a$  positive definite or semi-  
 2057 definite but non-zero. Thus:

$$2058 \sigma_{\text{iid}}^2 = \frac{1}{d_m} \text{tr}(\mathbf{J}^\top \Sigma_a \mathbf{J}) > 0. \\ 2059$$

2060 Our local convergence bound is:  
 2061

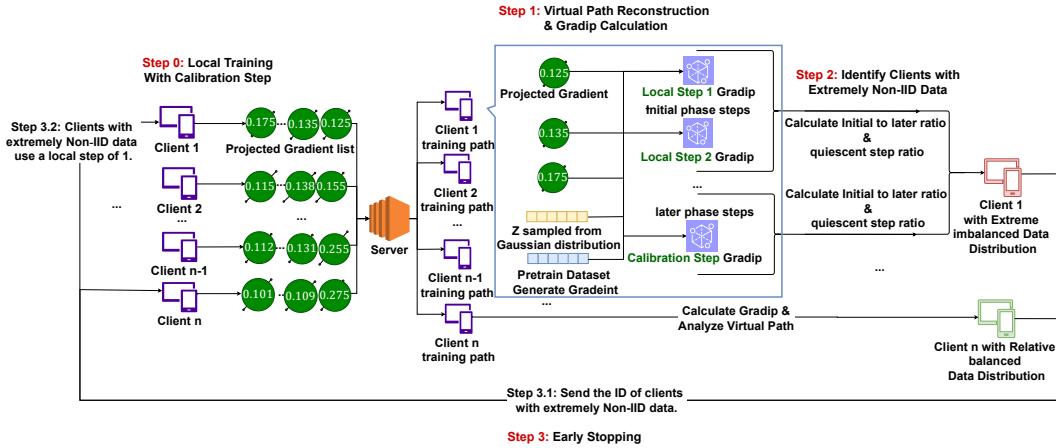
$$2062 \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f_k(w_k^t)\|^2 \leq O\left(\frac{1}{T}\right) + \sigma_k^2, \\ 2063 2064$$

2065 which indicates that in the steady state, the upper bound of the gradient norm is determined by  $\sigma_k^2$ .  
 2066 Therefore,  
 2067

2068  $\sigma_{\text{iid}}^2 \gg \sigma_{\text{non-iid}}^2 \approx 0 \implies \begin{cases} \text{IID clients: Gradient Norm oscillates significantly;} \\ \text{Non-IID clients: Gradient Norm decreases monotonically and approaches 0.} \end{cases}$   
 2069  
 2070  
 2071

## 2072 REMARKS

2073 In summary, by substituting the explicit form of the cross-entropy gradient into our sparse ZO  
 2074 convergence formula, we can empirically explain that due to the variance differences caused by  
 2075 label distributions, the Gradient Norms of IID clients maintains significant fluctuations, while that of  
 2076 extremely Non-IID clients rapidly decays and converges to zero.  
 2077



2093 Figure 5: MEERKAT-VP: Each client locally trains with a prescribed statistic step, yielding a  
 2094 sequence of projected gradients. The server leverages a randomly sampled vector  $z_k^t$  from the  
 2095 Gaussian distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  to reconstruct  $\nabla f_k^t$ , and then computes GradIP (see Definition 2.3)  
 2096 at every local training step. By analyzing the resulting GradIP values across all clients, the system  
 2097 distinguishes those clients with extremely Non-IID data from those that are relatively balanced. For  
 2098 the parameters **later phase steps**, **initial phase steps**, **quiescent step ratio**, and **initial to later ratio**,  
 2099 please refer to Table 3 in Appendix D.1

## 2100 D MORE EXPERIMENTAL DETAILS

### 2101 D.1 ADDITIONAL EXPERIMENTAL SETTINGS

2102 **Testbed.** All experiments are run on servers with the following configurations: RTX A6000 Setup:  
 2103 Ubuntu 18.04.6 LTS with 2 NVIDIA RTX A6000 GPUs (each with 48GB GPU memory). GH200

---

2106  
 2107  
 2108  
 2109  
 2110  
 2111

**Algorithm 2** MEERKAT: Sparse Zeroth-Order Optimization for Federated LLM Fine-Tuning

---

2112   **Input:** pre-trained weight  $\mathbf{w}_0$ , sparse mask  $\mathbf{m}$ , learning rate  $\eta$ , perturbation scale  $\epsilon$ , number of  
 2113    rounds  $R$ , total number of *clients*  $K$  number of local steps  $T$   
 2114    *Server* initiate seed list  $\{s_1^1, \dots, s_1^T\}$   
 2115   **for** Round  $r = 1$  to  $R$  **do**  
 2116     **Step 1. Local ZO update.**  
 2117     **for** each *client*  $k$  in **parallel** **do**  
 2118       Download model from *server*:  $\mathbf{w}_k \leftarrow \mathbf{w}_{r-1}$   
 2119       Download seed list  $\{s_r^1, \dots, s_r^T\}$  from *server*  
 2120       **for** local step  $t = 1$  to  $T$  **do**  
 2121         Initialize  $\mathbf{z}_k^t$  with seed  $s_r^t$ .  
 2122         Sample a batch  $\mathcal{B}$  on *client* dataset.  
 2123          $\tilde{\mathbf{w}}_k^t \leftarrow \mathbf{w}_k^t + \epsilon \cdot (\mathbf{z}_k^t \odot \mathbf{m})$   
 2124         Compute loss  $f_+ \leftarrow f(\tilde{\mathbf{w}}_k^t; \mathcal{B})$   
 2125          $\tilde{\mathbf{w}}_k^t \leftarrow \mathbf{w}_k^t - 2\epsilon \cdot (\mathbf{z}_k^t \odot \mathbf{m})$   
 2126         Compute loss:  $f_- \leftarrow f(\tilde{\mathbf{w}}_k^t; \mathcal{B})$   
 2127         Compute projected gradient:  
 2128         
$$g_k^t \leftarrow (f_+ - f_-)/2\epsilon$$
  
 2129        Update *client* model:  
 2130        
$$\hat{\nabla} f_k^t \leftarrow g_k^t \cdot (\mathbf{z}_k^t \odot \mathbf{m})$$
  
 2131        
$$\mathbf{w}_k^{t+1} \leftarrow \mathbf{w}_k^t - \eta \hat{\nabla} f_k^t$$
  
 2132       **end for**  
 2133       Send projected gradients  $\{g_k^1, g_k^2, \dots, g_k^T\}$  to *server*.  
 2134     **end for**  
 2135     **Step 2. Server recover each client's update with virtual path.**  
 2136     **for**  $k = 1$  to  $K$  **do**  
 2137       **for** local step  $t = 1$  to  $T$  **do**  
 2138         Generate  $\mathbf{z}_k^t$  with seed  $s_r^t$ .  
 2139         Perform *virtual path*:  
 2140         
$$\hat{\nabla} f_k^t = g_k^t \cdot (\mathbf{z}_k^t \odot \mathbf{m})$$
  
 2141         
$$\mathbf{w}_k^{t+1} \leftarrow \mathbf{w}_k^t - \eta \hat{\nabla} f_k^t$$
  
 2142       **end for**  
 2143       Store recover *client* model parameters  $\mathbf{w}_k^T$   
 2144     **end for**  
 2145     **end for**  
 2146     **Step 3. Server Aggregate reconstructed sparse model update.**  
 2147     
$$\mathbf{w}_r \leftarrow \frac{1}{K} \sum_{k=1}^K \mathbf{w}_k^T$$
  
 2148  
 2149     Generate new seed list  $\{s_{r+1}^1, \dots, s_{r+1}^T\}$ .  
 2150     **end for**  
 2151     **Output:**  $\mathbf{w}_R$

---

2152  
 2153  
 2154  
 2155  
 2156  
 2157  
 2158  
 2159

|      |                                                                                                                                    |                    |
|------|------------------------------------------------------------------------------------------------------------------------------------|--------------------|
| 2160 | <b>Algorithm 3</b> MEERKAT with high frequency server-client synchronization                                                       |                    |
| 2161 |                                                                                                                                    |                    |
| 2162 | <b>Input:</b> Seed $s$ and projected gradients $g_k^t$ from all clients, global model $\mathbf{m}$ , learning rate $\eta$ , number |                    |
| 2163 | of clients $K$ , sparse mask $\mathbf{m}$                                                                                          |                    |
| 2164 | <b>Aggregate projected gradients from all clients with same seed:</b>                                                              |                    |
| 2165 | $g \leftarrow \frac{1}{K} \sum_{k=1}^K g_k$                                                                                        |                    |
| 2166 |                                                                                                                                    |                    |
| 2167 |                                                                                                                                    |                    |
| 2168 | <b>Calculate Zeroth-Order Gradients:</b>                                                                                           |                    |
| 2169 | $\hat{\nabla}f \leftarrow g \cdot (\mathbf{z} \odot \mathbf{m})$                                                                   |                    |
| 2170 |                                                                                                                                    |                    |
| 2171 | <b>Update global model parameters:</b>                                                                                             |                    |
| 2172 | $\mathbf{w}_{r+1} \leftarrow \mathbf{w}_r - \eta (\hat{\nabla}f \odot \mathbf{m})$                                                 |                    |
| 2173 |                                                                                                                                    |                    |
| 2174 |                                                                                                                                    |                    |
| 2175 | <b>Generate new seed <math>s_{new}</math></b>                                                                                      |                    |
| 2176 | <b>Output:</b> Send aggregated global projected gradients $g$ and seed $s_{new}$ to all clients.                                   |                    |
| 2177 |                                                                                                                                    |                    |
| 2178 |                                                                                                                                    |                    |
| 2179 |                                                                                                                                    |                    |
| 2180 | Setup: Ubuntu 20.04 with 1 NVIDIA GH200 GPU (480GB GPU memory). A100 Setup: Ubuntu                                                 |                    |
| 2181 | 22.04 with 1 NVIDIA A100 GPU (40GB GPU memory).                                                                                    |                    |
| 2182 | <b>Dataset.</b> We conducted experiments using datasets from the GLUE and SuperGLUE benchmarks,                                    |                    |
| 2183 | including SST2, AgNews, Yelp, BoolQ, RTE, WSC, and WiC. To create IID clients data, we shuffle                                     |                    |
| 2184 | the entire dataset and evenly divide it among the clients. To create Non-IID clients data, we split the                            |                    |
| 2185 | data using a Dirichlet distribution. For all tasks, the Dirichlet $\alpha$ parameter is set to 0.5 to control the                  |                    |
| 2186 | degree of data heterogeneity.                                                                                                      |                    |
| 2187 | <b>Evaluation metric.</b> In our experiments, test accuracy is used as the primary evaluation metric.                              |                    |
| 2188 | Accuracy is computed as the proportion of correctly predicted labels across all evaluation samples.                                |                    |
| 2189 | Additionally, we incorporate the GradIP score (see Definition 2.3) to analyze further the dynamics of                              |                    |
| 2190 | local model training under IID and Non-IID client data settings. GradIP provides a metric to measure                               |                    |
| 2191 | the quality of client training trajectories, particularly in heterogeneous data distributions.                                     |                    |
| 2192 | <b>Notations.</b> We present the parameters definition used in MEERKAT-VP in Table 3.                                              |                    |
| 2193 |                                                                                                                                    |                    |
| 2194 |                                                                                                                                    |                    |
| 2195 |                                                                                                                                    |                    |
| 2196 | Table 3: MEERKAT-VP Parameters Notation                                                                                            |                    |
| 2197 |                                                                                                                                    |                    |
| 2198 |                                                                                                                                    |                    |
| 2199 |                                                                                                                                    |                    |
| 2200 |                                                                                                                                    |                    |
| 2201 |                                                                                                                                    |                    |
| 2202 |                                                                                                                                    |                    |
| 2203 |                                                                                                                                    |                    |
| 2204 |                                                                                                                                    |                    |
| 2205 | <b>Term</b>                                                                                                                        | <b>Explanation</b> |
| 2206 |                                                                                                                                    |                    |
| 2207 |                                                                                                                                    |                    |
| 2208 |                                                                                                                                    |                    |
| 2209 |                                                                                                                                    |                    |
| 2210 |                                                                                                                                    |                    |
| 2211 |                                                                                                                                    |                    |
| 2212 |                                                                                                                                    |                    |
| 2213 |                                                                                                                                    |                    |
| 2214 |                                                                                                                                    |                    |
| 2215 |                                                                                                                                    |                    |
| 2216 |                                                                                                                                    |                    |
| 2217 |                                                                                                                                    |                    |
| 2218 |                                                                                                                                    |                    |
| 2219 |                                                                                                                                    |                    |
| 2220 |                                                                                                                                    |                    |
| 2221 |                                                                                                                                    |                    |
| 2222 |                                                                                                                                    |                    |
| 2223 |                                                                                                                                    |                    |
| 2224 |                                                                                                                                    |                    |
| 2225 |                                                                                                                                    |                    |
| 2226 |                                                                                                                                    |                    |
| 2227 |                                                                                                                                    |                    |
| 2228 |                                                                                                                                    |                    |
| 2229 |                                                                                                                                    |                    |
| 2230 |                                                                                                                                    |                    |
| 2231 |                                                                                                                                    |                    |
| 2232 |                                                                                                                                    |                    |
| 2233 |                                                                                                                                    |                    |
| 2234 |                                                                                                                                    |                    |
| 2235 |                                                                                                                                    |                    |
| 2236 |                                                                                                                                    |                    |
| 2237 |                                                                                                                                    |                    |
| 2238 |                                                                                                                                    |                    |
| 2239 |                                                                                                                                    |                    |
| 2240 |                                                                                                                                    |                    |
| 2241 |                                                                                                                                    |                    |
| 2242 |                                                                                                                                    |                    |
| 2243 |                                                                                                                                    |                    |
| 2244 |                                                                                                                                    |                    |
| 2245 |                                                                                                                                    |                    |
| 2246 |                                                                                                                                    |                    |
| 2247 |                                                                                                                                    |                    |
| 2248 |                                                                                                                                    |                    |
| 2249 |                                                                                                                                    |                    |
| 2250 |                                                                                                                                    |                    |
| 2251 |                                                                                                                                    |                    |
| 2252 |                                                                                                                                    |                    |
| 2253 |                                                                                                                                    |                    |
| 2254 |                                                                                                                                    |                    |
| 2255 |                                                                                                                                    |                    |
| 2256 |                                                                                                                                    |                    |
| 2257 |                                                                                                                                    |                    |
| 2258 |                                                                                                                                    |                    |
| 2259 |                                                                                                                                    |                    |
| 2260 |                                                                                                                                    |                    |
| 2261 |                                                                                                                                    |                    |
| 2262 |                                                                                                                                    |                    |
| 2263 |                                                                                                                                    |                    |
| 2264 |                                                                                                                                    |                    |
| 2265 |                                                                                                                                    |                    |
| 2266 |                                                                                                                                    |                    |
| 2267 |                                                                                                                                    |                    |
| 2268 |                                                                                                                                    |                    |
| 2269 |                                                                                                                                    |                    |
| 2270 |                                                                                                                                    |                    |
| 2271 |                                                                                                                                    |                    |
| 2272 |                                                                                                                                    |                    |
| 2273 |                                                                                                                                    |                    |
| 2274 |                                                                                                                                    |                    |
| 2275 |                                                                                                                                    |                    |
| 2276 |                                                                                                                                    |                    |
| 2277 |                                                                                                                                    |                    |
| 2278 |                                                                                                                                    |                    |
| 2279 |                                                                                                                                    |                    |
| 2280 |                                                                                                                                    |                    |
| 2281 |                                                                                                                                    |                    |
| 2282 |                                                                                                                                    |                    |
| 2283 |                                                                                                                                    |                    |
| 2284 |                                                                                                                                    |                    |
| 2285 |                                                                                                                                    |                    |
| 2286 |                                                                                                                                    |                    |
| 2287 |                                                                                                                                    |                    |
| 2288 |                                                                                                                                    |                    |
| 2289 |                                                                                                                                    |                    |
| 2290 |                                                                                                                                    |                    |
| 2291 |                                                                                                                                    |                    |
| 2292 |                                                                                                                                    |                    |
| 2293 |                                                                                                                                    |                    |
| 2294 |                                                                                                                                    |                    |
| 2295 |                                                                                                                                    |                    |
| 2296 |                                                                                                                                    |                    |
| 2297 |                                                                                                                                    |                    |
| 2298 |                                                                                                                                    |                    |
| 2299 |                                                                                                                                    |                    |
| 2300 |                                                                                                                                    |                    |
| 2301 |                                                                                                                                    |                    |
| 2302 |                                                                                                                                    |                    |
| 2303 |                                                                                                                                    |                    |
| 2304 |                                                                                                                                    |                    |
| 2305 |                                                                                                                                    |                    |
| 2306 |                                                                                                                                    |                    |
| 2307 |                                                                                                                                    |                    |
| 2308 |                                                                                                                                    |                    |
| 2309 |                                                                                                                                    |                    |
| 2310 |                                                                                                                                    |                    |
| 2311 |                                                                                                                                    |                    |
| 2312 |                                                                                                                                    |                    |
| 2313 |                                                                                                                                    |                    |
| 2314 |                                                                                                                                    |                    |
| 2315 |                                                                                                                                    |                    |
| 2316 |                                                                                                                                    |                    |
| 2317 |                                                                                                                                    |                    |
| 2318 |                                                                                                                                    |                    |
| 2319 |                                                                                                                                    |                    |
| 2320 |                                                                                                                                    |                    |
| 2321 |                                                                                                                                    |                    |
| 2322 |                                                                                                                                    |                    |
| 2323 |                                                                                                                                    |                    |
| 2324 |                                                                                                                                    |                    |
| 2325 |                                                                                                                                    |                    |
| 2326 |                                                                                                                                    |                    |
| 2327 |                                                                                                                                    |                    |
| 2328 |                                                                                                                                    |                    |
| 2329 |                                                                                                                                    |                    |
| 2330 |                                                                                                                                    |                    |
| 2331 |                                                                                                                                    |                    |
| 2332 |                                                                                                                                    |                    |
| 2333 |                                                                                                                                    |                    |
| 2334 |                                                                                                                                    |                    |
| 2335 |                                                                                                                                    |                    |
| 2336 |                                                                                                                                    |                    |
| 2337 |                                                                                                                                    |                    |
| 2338 |                                                                                                                                    |                    |
| 2339 |                                                                                                                                    |                    |
| 2340 |                                                                                                                                    |                    |
| 2341 |                                                                                                                                    |                    |
| 2342 |                                                                                                                                    |                    |
| 2343 |                                                                                                                                    |                    |
| 2344 |                                                                                                                                    |                    |
| 2345 |                                                                                                                                    |                    |
| 2346 |                                                                                                                                    |                    |
| 2347 |                                                                                                                                    |                    |
| 2348 |                                                                                                                                    |                    |
| 2349 |                                                                                                                                    |                    |
| 2350 |                                                                                                                                    |                    |
| 2351 |                                                                                                                                    |                    |
| 2352 |                                                                                                                                    |                    |
| 2353 |                                                                                                                                    |                    |
| 2354 |                                                                                                                                    |                    |
| 2355 |                                                                                                                                    |                    |
| 2356 |                                                                                                                                    |                    |
| 2357 |                                                                                                                                    |                    |
| 2358 |                                                                                                                                    |                    |
| 2359 |                                                                                                                                    |                    |
| 2360 |                                                                                                                                    |                    |
| 2361 |                                                                                                                                    |                    |
| 2362 |                                                                                                                                    |                    |
| 2363 |                                                                                                                                    |                    |
| 2364 |                                                                                                                                    |                    |
| 2365 |                                                                                                                                    |                    |
| 2366 |                                                                                                                                    |                    |
| 2367 |                                                                                                                                    |                    |
| 2368 |                                                                                                                                    |                    |
| 2369 |                                                                                                                                    |                    |
| 2370 |                                                                                                                                    |                    |
| 2371 |                                                                                                                                    |                    |
| 2372 |                                                                                                                                    |                    |
| 2373 |                                                                                                                                    |                    |
| 2374 |                                                                                                                                    |                    |
| 2375 |                                                                                                                                    |                    |
| 2376 |                                                                                                                                    |                    |
| 2377 |                                                                                                                                    |                    |
| 2378 |                                                                                                                                    |                    |
| 2379 |                                                                                                                                    |                    |
| 2380 |                                                                                                                                    |                    |
| 2381 |                                                                                                                                    |                    |
| 2382 |                                                                                                                                    |                    |
| 2383 |                                                                                                                                    |                    |
| 2384 |                                                                                                                                    |                    |
| 2385 |                                                                                                                                    |                    |
| 2386 |                                                                                                                                    |                    |
| 2387 |                                                                                                                                    |                    |
| 2388 |                                                                                                                                    |                    |
| 2389 |                                                                                                                                    |                    |
| 2390 |                                                                                                                                    |                    |
| 2391 |                                                                                                                                    |                    |
| 2392 |                                                                                                                                    |                    |
| 2393 |                                                                                                                                    |                    |
| 2394 |                                                                                                                                    |                    |
| 2395 |                                                                                                                                    |                    |
| 2396 |                                                                                                                                    |                    |
| 2397 |                                                                                                                                    |                    |
| 2398 |                                                                                                                                    |                    |
| 2399 |                                                                                                                                    |                    |
| 2400 |                                                                                                                                    |                    |
| 2401 |                                                                                                                                    |                    |
| 2402 |                                                                                                                                    |                    |
| 2403 |                                                                                                                                    |                    |
| 2404 |                                                                                                                                    |                    |
| 2405 |                                                                                                                                    |                    |
| 2406 |                                                                                                                                    |                    |
| 2407 |                                                                                                                                    |                    |
| 2408 |                                                                                                                                    |                    |
| 2409 |                                                                                                                                    |                    |
| 2410 |                                                                                                                                    |                    |
| 2411 |                                                                                                                                    |                    |
| 2412 |                                                                                                                                    |                    |
| 2413 |                                                                                                                                    |                    |
| 2414 |                                                                                                                                    |                    |
| 2415 |                                                                                                                                    |                    |
| 2416 |                                                                                                                                    |                    |
| 2417 |                                                                                                                                    |                    |
| 2418 |                                                                                                                                    |                    |
| 2419 |                                                                                                                                    |                    |
| 2420 |                                                                                                                                    |                    |
| 2421 |                                                                                                                                    |                    |
| 2422 |                                                                                                                                    |                    |
| 2423 |                                                                                                                                    |                    |
| 2424 |                                                                                                                                    |                    |
| 2425 |                                                                                                                                    |                    |
| 2426 |                                                                                                                                    |                    |
| 2427 |                                                                                                                                    |                    |
| 2428 |                                                                                                                                    |                    |
| 2429 |                                                                                                                                    |                    |
| 2430 |                                                                                                                                    |                    |
| 2431 |                                                                                                                                    |                    |
| 2432 |                                                                                                                                    |                    |
| 2433 |                                                                                                                                    |                    |
| 2434 |                                                                                                                                    |                    |
| 2435 |                                                                                                                                    |                    |
| 2436 |                                                                                                                                    |                    |
| 2437 |                                                                                                                                    |                    |
| 2438 |                                                                                                                                    |                    |
| 2439 |                                                                                                                                    |                    |
| 2440 |                                                                                                                                    |                    |
| 2441 |                                                                                                                                    |                    |
| 2442 |                                                                                                                                    |                    |
| 2443 |                                                                                                                                    |                    |
| 2444 |                                                                                                                                    |                    |
| 2445 |                                                                                                                                    |                    |
| 2446 |                                                                                                                                    |                    |
| 2447 |                                                                                                                                    |                    |
| 2448 |                                                                                                                                    |                    |
| 2449 |                                                                                                                                    |                    |
| 2450 |                                                                                                                                    |                    |
| 2451 |                                                                                                                                    |                    |
| 2452 |                                                                                                                                    |                    |
| 2453 |                                                                                                                                    |                    |
| 2454 |                                                                                                                                    |                    |
| 2455 |                                                                                                                                    |                    |
| 2456 |                                                                                                                                    |                    |
| 2457 |                                                                                                                                    |                    |
| 2458 |                                                                                                                                    |                    |
| 2459 |                                                                                                                                    |                    |
| 2460 |                                                                                                                                    |                    |
| 2461 |                                                                                                                                    |                    |
| 2462 |                                                                                                                                    |                    |
| 2463 |                                                                                                                                    |                    |
| 2464 |                                                                                                                                    |                    |
| 2465 |                                                                                                                                    |                    |
| 2466 |                                                                                                                                    |                    |
| 2467 |                                                                                                                                    |                    |
| 2468 |                                                                                                                                    |                    |
| 2469 |                                                                                                                                    |                    |
| 2470 |                                                                                                                                    |                    |
| 2471 |                                                                                                                                    |                    |
| 2472 |                                                                                                                                    |                    |
| 2473 |                                                                                                                                    |                    |
| 2474 |                                                                                                                                    |                    |
| 2475 |                                                                                                                                    |                    |
| 2476 |                                                                                                                                    |                    |
| 2477 |                                                                                                                                    |                    |
| 2478 |                                                                                                                                    |                    |
| 2479 |                                                                                                                                    |                    |
| 2480 |                                                                                                                                    |                    |
| 2481 |                                                                                                                                    |                    |
| 2482 |                                                                                                                                    |                    |
| 2483 |                                                                                                                                    |                    |
| 2484 |                                                                                                                                    |                    |
| 2485 |                                                                                                                                    |                    |
| 2486 |                                                                                                                                    |                    |
| 2487 |                                                                                                                                    |                    |
| 2488 |                                                                                                                                    |                    |
| 2489 |                                                                                                                                    |                    |
| 2490 |                                                                                                                                    |                    |
| 2491 |                                                                                                                                    |                    |
| 2492 |                                                                                                                                    |                    |
| 2493 |                                                                                                                                    |                    |
| 2494 |                                                                                                                                    |                    |
| 2495 |                                                                                                                                    |                    |
| 2496 |                                                                                                                                    |                    |
| 2497 |                                                                                                                                    |                    |
| 2498 |                                                                                                                                    |                    |
| 2499 |                                                                                                                                    |                    |
| 2500 |                                                                                                                                    |                    |
| 2501 |                                                                                                                                    |                    |
| 2502 |                                                                                                                                    |                    |
| 2503 |                                                                                                                                    |                    |
| 2504 |                                                                                                                                    |                    |
| 2505 |                                                                                                                                    |                    |
| 2506 |                                                                                                                                    |                    |
| 2507 |                                                                                                                                    |                    |
| 2508 |                                                                                                                                    |                    |
| 2509 |                                                                                                                                    |                    |
| 2510 |                                                                                                                                    |                    |
| 2511 |                                                                                                                                    |                    |
| 2512 |                                                                                                                                    |                    |
| 2513 |                                                                                                                                    |                    |
| 2514 |                                                                                                                                    |                    |
| 2515 |                                                                                                                                    |                    |
| 2516 |                                                                                                                                    |                    |
| 2517 |                                                                                                                                    |                    |
| 2518 |                                                                                                                                    |                    |
| 2519 |                                                                                                                                    |                    |
| 2520 |                                                                                                                                    |                    |
| 2521 |                                                                                                                                    |                    |
| 2522 |                                                                                                                                    |                    |
| 2523 |                                                                                                                                    |                    |
| 2524 |                                                                                                                                    |                    |
| 2525 |                                                                                                                                    |                    |
| 2526 |                                                                                                                                    |                    |
| 2527 |                                                                                                                                    |                    |
| 2528 |                                                                                                                                    |                    |
| 2529 |                                                                                                                                    |                    |
| 2530 |                                                                                                                                    |                    |
| 2531 |                                                                                                                                    |                    |
| 2532 |                                                                                                                                    |                    |
| 2533 |                                                                                                                                    |                    |
| 2534 |                                                                                                                                    |                    |
| 2535 |                                                                                                                                    |                    |
| 2536 |                                                                                                                                    |                    |
| 2537 |                                                                                                                                    |                    |
| 2538 |                                                                                                                                    |                    |
| 2539 |                                                                                                                                    |                    |
| 2540 |                                                                                                                                    |                    |
| 2541 |                                                                                                                                    |                    |
| 2542 |                                                                                                                                    |                    |
| 2543 |                                                                                                                                    |                    |
| 2544 |                                                                                                                                    |                    |
| 2545 |                                                                                                                                    |                    |
| 2546 |                                                                                                                                    |                    |
| 2547 |                                                                                                                                    |                    |
| 2548 |                                                                                                                                    |                    |
| 2549 |                                                                                                                                    |                    |
| 2550 |                                                                                                                                    |                    |
| 2551 |                                                                                                                                    |                    |
| 2552 |                                                                                                                                    |                    |
| 2553 |                                                                                                                                    |                    |
| 2554 |                                                                                                                                    |                    |
| 2555 |                                                                                                                                    |                    |
| 2556 |                                                                                                                                    |                    |
| 2557 |                                                                                                                                    |                    |
| 2558 |                                                                                                                                    |                    |
| 2559 |                                                                                                                                    |                    |
| 2560 |                                                                                                                                    |                    |
| 2561 |                                                                                                                                    |                    |
| 2562 |                                                                                                                                    |                    |
| 2563 |                                                                                                                                    |                    |
| 2564 |                                                                                                                                    |                    |
| 2565 |                                                                                                                                    |                    |
| 2566 |                                                                                                                                    |                    |
| 2567 |                                                                                                                                    |                    |
| 2568 |                                                                                                                                    |                    |
| 2569 |                                                                                                                                    |                    |
| 2570 |                                                                                                                                    |                    |
| 2571 |                                                                                                                                    |                    |
| 2572 |                                                                                                                                    |                    |
| 2573 |                                                                                                                                    |                    |
| 2574 |                                                                                                                                    |                    |
| 2575 |                                                                                                                                    |                    |
| 2576 |                                                                                                                                    |                    |
| 2577 |                                                                                                                                    |                    |
| 2578 |                                                                                                                                    |                    |
| 2579 |                                                                                                                                    |                    |
| 2580 |                                                                                                                                    |                    |
| 2581 |                                                                                                                                    |                    |
| 2582 |                                                                                                                                    |                    |
| 2583 |                                                                                                                                    |                    |
| 2584 |                                                                                                                                    |                    |
| 2585 |                                                                                                                                    |                    |
| 2586 |                                                                                                                                    |                    |
| 2587 |                                                                                                                                    |                    |
| 2588 |                                                                                                                                    |                    |
| 2589 |                                                                                                                                    |                    |
| 2590 |                                                                                                                                    |                    |
| 2591 |                                                                                                                                    |                    |
| 2592 |                                                                                                                                    |                    |
| 2593 |                                                                                                                                    |                    |
| 2594 |                                                                                                                                    |                    |
| 2595 |                                                                                                                                    |                    |
| 2596 |                                                                                                                                    |                    |
| 2597 |                                                                                                                                    |                    |
| 2598 |                                                                                                                                    |                    |
| 2599 |                                                                                                                                    |                    |
| 2600 |                                                                                                                                    |                    |
| 2601 |                                                                                                                                    |                    |
| 2602 |                                                                                                                                    |                    |
| 2603 |                                                                                                                                    |                    |
| 2604 |                                                                                                                                    |                    |
| 2605 |                                                                                                                                    |                    |
| 2606 |                                                                                                                                    |                    |
| 2607 |                                                                                                                                    |                    |
| 2608 |                                                                                                                                    |                    |
| 2609 |                                                                                                                                    |                    |
| 2610 |                                                                                                                                    |                    |
| 2611 |                                                                                                                                    |                    |
| 2612 |                                                                                                                                    |                    |
| 2613 |                                                                                                                                    |                    |
| 2614 |                                                                                                                                    |                    |
| 2615 |                                                                                                                                    |                    |
| 2616 |                                                                                                                                    |                    |
| 2617 |                                                                                                                                    |                    |
| 2618 |                                                                                                                                    |                    |
| 2619 |                                                                                                                                    |                    |
| 2620 |                                                                                                                                    |                    |
| 2621 |                                                                                                                                    |                    |
| 2622 |                                                                                                                                    |                    |
| 2623 |                                                                                                                                    |                    |
| 2624 |                                                                                                                                    |                    |
| 2625 |                                                                                                                                    |                    |
| 2626 |                                                                                                                                    |                    |
| 2627 |                                                                                                                                    |                    |
| 2628 |                                                                                                                                    |                    |
| 2629 |                                                                                                                                    |                    |
| 2630 |                                                                                                                                    |                    |
| 2631 |                                                                                                                                    |                    |
| 2632 |                                                                                                                                    |                    |
| 2633 |                                                                                                                                    |                    |
| 2634 |                                                                                                                                    |                    |
| 2635 |                                                                                                                                    |                    |
| 2636 |                                                                                                                                    |                    |
| 2637 |                                                                                                                                    |                    |
| 2638 |                                                                                                                                    |                    |
| 2639 |                                                                                                                                    |                    |
| 2640 |                                                                                                                                    |                    |
| 2641 |                                                                                                                                    |                    |
| 2642 |                                                                                                                                    |                    |
| 2643 |                                                                                                                                    |                    |
| 2644 |                                                                                                                                    |                    |
| 2645 |                                                                                                                                    |                    |
| 2646 |                                                                                                                                    |                    |
| 2647 |                                                                                                                                    |                    |
| 2648 |                                                                                                                                    |                    |
| 2649 |                                                                                                                                    |                    |
| 2650 |                                                                                                                                    |                    |
| 2651 |                                                                                                                                    |                    |
| 2652 |                                                                                                                                    |                    |
| 2653 |                                                                                                                                    |                    |
| 2654 |                                                                                                                                    |                    |
| 2655 |                                                                                                                                    |                    |
| 2656 |                                                                                                                                    |                    |
| 2657 |                                                                                                                                    |                    |
| 2658 |                                                                                                                                    |                    |
| 2659 |                                                                                                                                    |                    |
| 2660 |                                                                                                                                    |                    |
| 2661 |                                                                                                                                    |                    |
| 2662 |                                                                                                                                    |                    |
| 2663 |                                                                                                                                    |                    |
| 2664 |                                                                                                                                    |                    |
| 2665 |                                                                                                                                    |                    |
| 2666 |                                                                                                                                    |                    |
| 2667 |                                                                                                                                    |                    |
| 2668 |                                                                                                                                    |                    |

Table 4: Hyper-parameters used in our experiments.

| Parameter                | Value               |
|--------------------------|---------------------|
| MEERKAT learning rate    | [2e-4, 2e-8]        |
| MEERKAT-VP learning rate | [2e-4, 2e-8]        |
| LoRA-FedZO learning rate | [2e-4, 2e-8]        |
| Full-FedZO learning rate | [2e-4, 2e-8]        |
| Batch size               | 16                  |
| Dirichlet alpha          | 0.5, 0.3, 0.1       |
| LoRA rank                | 16                  |
| LoRA alpha               | 16                  |
| initial phase steps      | 20                  |
| later phase steps        | 20                  |
| convergence threshold    | 1                   |
| quiescent step ratio     | [0.4, 0.5, 0.7]     |
| Initial to later ratio   | [1.5, 2, 5, 10, 15] |
| calibration steps        | 100                 |
| Total clients            | 10                  |

Table 5: Default MEERKAT-VP Hyperparameter Values

| initial phase steps | later phase steps | convergence threshold | quiescent step ratio | Initial to later ratio |
|---------------------|-------------------|-----------------------|----------------------|------------------------|
| 20                  | 20                | 1                     | 0.5                  | 5                      |

Table 6: Task-Specific VPCS Hyperparameters for RTE Task

| Model        | initial phase steps | later phase steps | convergence threshold | quiescent step ratio | Initial to later ratio |
|--------------|---------------------|-------------------|-----------------------|----------------------|------------------------|
| Gemma2-2B    | 20                  | 20                | 1                     | 0.7                  | 5                      |
| LLaMA-3.2-1B | 20                  | 20                | 0.5                   | 0.7                  | 5                      |
| Qwen2-1.5B   | 20                  | 20                | 0.5                   | 0.5                  | 5                      |

Table 7: Parameter Sensitivity Analysis for LLaMA-3.2-1B on SST-2 Task

| initial phase steps | later phase steps | convergence threshold | quiescent step ratio | Initial to later ratio | Performance |
|---------------------|-------------------|-----------------------|----------------------|------------------------|-------------|
| 20                  | 20                | 1                     | 0.5                  | 3                      | 0.922       |
| 20                  | 20                | 1                     | 0.5                  | 5                      | 0.922       |
| 20                  | 20                | 1                     | 0.5                  | 7                      | 0.922       |
| 20                  | 20                | 1                     | 0.5                  | 10                     | 0.922       |
| 20                  | 20                | 1                     | 0.5                  | 12                     | 0.922       |

Table 8: Parameter Sensitivity Analysis for RTE Task

| Model        | initial phase steps | later phase steps | convergence threshold | quiescent step ratio | Initial to later ratio | Performance |
|--------------|---------------------|-------------------|-----------------------|----------------------|------------------------|-------------|
| LLaMA-3.2-1B | 20                  | 20                | 0.8                   | 0.5                  | 7                      | 0.617       |
| LLaMA-3.2-1B | 20                  | 20                | 0.7                   | 0.5                  | 5                      | 0.617       |
| Gemma2-2B    | 20                  | 20                | 1                     | 0.5                  | 3                      | 0.657       |
| Gemma2-2B    | 20                  | 20                | 1                     | 0.5                  | 5                      | 0.657       |

## D.2 ADDITIONAL EXPERIMENT RESULTS

In this section, we present additional experimental results to compare MEERKAT, MEERKAT-VP, Full-FedZO, and LoRA-FedZO under various settings. The results include five tables and three figures, providing a detailed evaluation of performance across different models, datasets and experiment settings. Table 3 provides a description of the parameters used in MEERKAT-VP, and Table 4 lists the experiment parameters used in this experiment. Tables 5 and 6 list the hyperparameter values for MEERKAT-VP. Tables 7 and 8 demonstrate the robustness of the MEERKAT-VP parameter selection. Table 9 provides a quantitative analysis that demonstrates the significant disparity in gradient sensitivity across different parameter groups, thereby justifying our selection criteria. Table 11 shows that a domain-shifted calibration dataset can be used effectively to select sensitive model parameters. Furthermore, we designed an experiment where each client builds a local parameter

mask from its own dataset. The results demonstrate that aggregating these local masks into a union mask does not achieve better performance than using a single, globally unified mask. Table 13 compares MEERKAT and Full-FedZO on multiple tasks at the same communication frequency for Llama-3.2-1B, Qwen2-1.5B, and Gemma-2-2b models. Table 14 presents results in a Non-IID client data scenario, comparing MEERKAT-VP and MEERKAT under the same communication frequency and sparsity density, and demonstrating MEERKAT-VP improved performance. Table 15 investigates the robustness of MEERKAT by evaluating test accuracy with local step 1 across different sparsity densities. Table 16 compares MEERKAT, Full-FedZO and LoRA-FedZO under high communication frequency across IID and Non-IID client data settings. Table 23 details the number of training rounds required for convergence across different models and tasks. Table 24 benchmarks computational and communication efficiency, demonstrating that MEERKAT significantly reduces peak RAM usage and client download bandwidth compared to the Full-FedZO and LoRA-FedZO baselines. Table 25 shows that our MEERKAT-VP method achieves competitive performance against the back-propagation upper bound and substantially outperforms FedDYN Acar et al. (2021). Figure 7 and Figure 9 further illustrate the phenomenon of GradIP under IID and Non-IID client data settings.

Table 9: Gradient Sensitivity Analysis for Qwen2-1.5B Model on C4 Dataset (Top 0.1% Parameters).

To quantitatively analyze gradient sensitivity, we ranked all parameters by their average squared gradients from pre-training and divided them into four disjoint (non-overlapping) buckets: 0-0.1%, 0.1-1%, 1-10% and 10%-100%.

| Bucket / Metric     | Top 0.1%               | 0.1%-1%                | 1%-10%                 | 10%-100%                |
|---------------------|------------------------|------------------------|------------------------|-------------------------|
| Avg Gradient Square | $4.403 \times 10^{-3}$ | $8.536 \times 10^{-5}$ | $1.075 \times 10^{-5}$ | $1.764 \times 10^{-6}$  |
| Std Gradient Square | $8.094 \times 10^{-2}$ | $5.858 \times 10^{-5}$ | $6.255 \times 10^{-6}$ | $1.099 \times 10^{-6}$  |
| Max Gradient Square | $1.413 \times 10^1$    | $3.147 \times 10^{-4}$ | $3.505 \times 10^{-5}$ | $5.245 \times 10^{-6}$  |
| Min Gradient Square | $3.166 \times 10^{-4}$ | $3.529 \times 10^{-5}$ | $5.245 \times 10^{-6}$ | $1.025 \times 10^{-19}$ |

Table 10: Accuracy of MEERKAT vs. Random-Select (Qwen2-1.5B, 0.1% mask). Directly addressing the comparison with random selection, we ran a control experiment that shows our method is significantly better across all tasks. The local step is 10.

| Method             | SST-2        | AGNews       | Yelp         | BoolQ        | RTE          | WSC          | WIC          | Avg          |
|--------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| MEERKAT            | <b>0.949</b> | <b>0.881</b> | <b>0.934</b> | <b>0.752</b> | <b>0.813</b> | <b>0.682</b> | <b>0.628</b> | <b>0.806</b> |
| Random Select      | 0.821        | 0.543        | 0.852        | 0.667        | 0.711        | 0.663        | 0.539        | 0.685        |
| <b>Improvement</b> | +12.8%       | +33.8%       | +8.2%        | +8.5%        | +10.2%       | +1.9%        | +8.9%        | +12.1%       |

2322  
 2323 Table 11: Performance Comparison with Different Calibration Datasets and Methods. Our method  
 2324 does not require the original pre-training data. It uses a small sample (128 sequences) from any  
 2325 public, high-quality text corpus to create a transferable parameter mask. This table confirms  
 2326 MEERKAT’s flexibility and transferability across different domains, including web-text, code, and  
 2327 medical data, consistently outperforming the Full-FedZO baseline. We also explore UnionMask, a  
 2328 client-specific mask aggregation approach: (1) Each client computes its own mask based on local  
 2329 data distribution; (2) Clients send masks to the server for aggregation into a union mask; (3) All  
 2330 clients use this union mask for ZO training; (4) The server uses the union mask for parameter  
 2331 updates. Results show that the specialized UnionMask performs similarly to our transferable mask,  
 2332 validating our universality approach. The local step is 10. Code data: microsoft/rStar-Coder. Medical  
 2333 data: FreedomIntelligence/medical-01-reasoning-SFT.  
 2334

| Method                                  | SST-2        | AGNews       | Yelp         | BoolQ        | RTE          | WSC          | WIC          | Avg          |
|-----------------------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Full-FedZO                              | 0.909        | 0.705        | 0.940        | 0.641        | 0.542        | 0.634        | 0.523        | 0.699        |
| <i>Web-Text Domain Calibration Data</i> |              |              |              |              |              |              |              |              |
| MEERKAT (C4, 0.1%)                      | <b>0.916</b> | <b>0.872</b> | <b>0.964</b> | 0.695        | <b>0.600</b> | <b>0.653</b> | <b>0.614</b> | <b>0.759</b> |
| MEERKAT (Wiki, 0.1%)                    | 0.913        | 0.855        | 0.952        | 0.646        | 0.582        | 0.634        | 0.567        | 0.736        |
| MEERKAT (ArXiv, 0.1%)                   | 0.901        | 0.851        | 0.949        | <b>0.714</b> | 0.573        | 0.644        | 0.562        | 0.742        |
| MEERKAT (FineWeb, 0.1%)                 | 0.902        | 0.846        | 0.958        | 0.695        | 0.584        | 0.634        | 0.561        | 0.740        |
| <i>Domain-Shifted Calibration Data</i>  |              |              |              |              |              |              |              |              |
| MEERKAT (Code, 0.1%)                    | 0.915        | 0.843        | 0.956        | 0.695        | 0.551        | 0.612        | 0.602        | 0.739        |
| MEERKAT (Bio, 0.1%)                     | 0.912        | 0.850        | 0.956        | 0.694        | 0.560        | 0.625        | 0.595        | 0.742        |
| <i>Client-Specific Mask Aggregation</i> |              |              |              |              |              |              |              |              |
| UnionMask (per-client, C4, 0.1%)        | 0.902        | 0.845        | 0.950        | 0.669        | 0.582        | 0.634        | 0.569        | 0.736        |

2346  
 2347  
 2348 Table 12: Transferability of the sparse mask between legal-domain (MultiEURLEX) Chalkidis et al.  
 2349 (2021) calibration datasets on LLaMA-3.2-1B.  
 2350

| Mask domain                | SST2  | AgNews | Yelp  | BoolQ |
|----------------------------|-------|--------|-------|-------|
| Legal-domain (MultiEURLEX) | 0.912 | 0.845  | 0.948 | 0.703 |

2354  
 2355  
 2356 Table 13: Performance comparison of MEERKAT and Full-FedZO on tasks SST-2, AgNews, Yelp,  
 2357 BoolQ, RTE, WSC, WIC under an IID client data setting. “Acc” is the average test accuracy across  
 2358 tasks. Bold numbers indicate the highest value in each row.  
 2359

|              | Methods    | Local Step | SST-2        | AgNews       | Yelp         | BoolQ        | RTE          | WSC          | WIC          | Acc          |
|--------------|------------|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| LLaMA-3.2-1B | Full-FedZO | 10         | 0.913        | 0.700        | 0.938        | 0.646        | 0.537        | 0.634        | 0.540        | 0.701        |
|              | MEERKAT    | 10         | <b>0.925</b> | <b>0.881</b> | <b>0.964</b> | <b>0.751</b> | <b>0.684</b> | 0.634        | <b>0.648</b> | <b>0.784</b> |
|              | Full-FedZO | 30         | 0.913        | 0.700        | 0.935        | 0.643        | 0.542        | 0.634        | 0.528        | 0.699        |
|              | MEERKAT    | 30         | <b>0.919</b> | <b>0.865</b> | <b>0.967</b> | <b>0.729</b> | <b>0.644</b> | <b>0.663</b> | <b>0.617</b> | <b>0.772</b> |
|              | Full-FedZO | 50         | 0.913        | 0.698        | 0.939        | 0.641        | 0.520        | 0.634        | 0.539        | 0.698        |
|              | MEERKAT    | 50         | <b>0.920</b> | <b>0.871</b> | <b>0.966</b> | <b>0.734</b> | <b>0.648</b> | <b>0.653</b> | <b>0.614</b> | <b>0.772</b> |
| Qwen2-1.5b   | Full-FedZO | 100        | 0.903        | 0.705        | 0.934        | 0.656        | 0.537        | 0.634        | 0.537        | 0.701        |
|              | MEERKAT    | 100        | <b>0.913</b> | <b>0.842</b> | <b>0.945</b> | <b>0.722</b> | <b>0.573</b> | 0.634        | <b>0.595</b> | <b>0.746</b> |
|              | Full-FedZO | 10         | 0.891        | 0.701        | 0.931        | 0.696        | 0.800        | 0.682        | 0.579        | 0.754        |
|              | MEERKAT    | 10         | <b>0.944</b> | <b>0.889</b> | <b>0.942</b> | <b>0.788</b> | <b>0.817</b> | <b>0.700</b> | <b>0.656</b> | <b>0.819</b> |
|              | Full-FedZO | 30         | 0.902        | 0.702        | 0.930        | 0.709        | 0.817        | 0.663        | 0.583        | 0.758        |
|              | MEERKAT    | 30         | <b>0.942</b> | <b>0.895</b> | <b>0.940</b> | <b>0.786</b> | <b>0.840</b> | <b>0.710</b> | <b>0.659</b> | <b>0.825</b> |
| Gemma2-2b    | Full-FedZO | 50         | 0.902        | 0.705        | 0.929        | 0.701        | 0.808        | <b>0.663</b> | 0.590        | 0.757        |
|              | MEERKAT    | 50         | <b>0.942</b> | <b>0.885</b> | <b>0.934</b> | <b>0.784</b> | <b>0.840</b> | 0.634        | <b>0.637</b> | <b>0.808</b> |
|              | Full-FedZO | 100        | 0.899        | 0.714        | 0.928        | 0.705        | <b>0.831</b> | <b>0.682</b> | 0.594        | 0.765        |
|              | MEERKAT    | 100        | <b>0.946</b> | <b>0.886</b> | <b>0.930</b> | <b>0.776</b> | 0.804        | 0.653        | <b>0.653</b> | <b>0.807</b> |
|              | Full-FedZO | 10         | 0.87         | 0.732        | 0.944        | 0.717        | 0.564        | 0.634        | 0.592        | 0.723        |
|              | MEERKAT    | 10         | <b>0.943</b> | <b>0.892</b> | <b>0.97</b>  | <b>0.817</b> | <b>0.724</b> | <b>0.653</b> | <b>0.636</b> | <b>0.805</b> |
|              | Full-FedZO | 30         | 0.91         | 0.81         | 0.942        | 0.73         | 0.56         | 0.644        | 0.578        | 0.739        |
|              | MEERKAT    | 30         | <b>0.943</b> | <b>0.887</b> | <b>0.973</b> | <b>0.812</b> | <b>0.617</b> | <b>0.663</b> | <b>0.608</b> | <b>0.786</b> |
|              | Full-FedZO | 50         | 0.911        | 0.812        | 0.942        | 0.735        | 0.551        | 0.634        | 0.572        | 0.737        |
|              | MEERKAT    | 50         | <b>0.94</b>  | <b>0.873</b> | <b>0.964</b> | <b>0.812</b> | <b>0.604</b> | 0.634        | <b>0.617</b> | <b>0.778</b> |
|              | Full-FedZO | 100        | 0.917        | 0.83         | 0.936        | 0.728        | 0.56         | <b>0.644</b> | 0.59         | 0.744        |
|              | MEERKAT    | 100        | <b>0.949</b> | <b>0.87</b>  | <b>0.954</b> | <b>0.815</b> | <b>0.568</b> | 0.634        | <b>0.592</b> | <b>0.769</b> |

Table 14: Comparison of MEERKAT-VP and MEERKAT under Non-IID client data setting, with the same local step and sparsity. Tasks include SST-2, AgNews, Yelp, BoolQ, RTE, WSC, and WIC. “Acc” indicates the average test accuracy across all tasks. Bold numbers highlight the best result in each row.

|              | Methods    | Local Step | SST-2        | AgNews       | Yelp         | BoolQ        | RTE          | WSC          | WIC          | Acc          |
|--------------|------------|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| LLaMA-3.2-1B | MEERKAT-VP | 10         | <b>0.922</b> | 0.864        | 0.962        | <b>0.713</b> | <b>0.617</b> | 0.644        | 0.625        | <b>0.764</b> |
|              | MEERKAT    | 10         | 0.916        | <b>0.872</b> | <b>0.964</b> | 0.695        | 0.600        | <b>0.653</b> | <b>0.614</b> | 0.759        |
|              | MEERKAT-VP | 30         | <b>0.919</b> | 0.825        | 0.963        | <b>0.685</b> | <b>0.595</b> | 0.634        | <b>0.631</b> | <b>0.750</b> |
|              | MEERKAT    | 30         | 0.897        | <b>0.862</b> | <b>0.965</b> | 0.646        | 0.577        | <b>0.644</b> | 0.583        | 0.739        |
|              | MEERKAT-VP | 50         | 0.909        | <b>0.836</b> | 0.959        | <b>0.691</b> | 0.577        | 0.615        | <b>0.615</b> | <b>0.743</b> |
|              | MEERKAT    | 50         | 0.909        | 0.827        | <b>0.965</b> | 0.647        | <b>0.595</b> | <b>0.634</b> | 0.567        | 0.734        |
|              | MEERKAT-VP | 100        | <b>0.904</b> | <b>0.824</b> | <b>0.962</b> | <b>0.684</b> | 0.577        | <b>0.653</b> | <b>0.630</b> | <b>0.747</b> |
|              | MEERKAT    | 100        | 0.896        | 0.777        | 0.961        | 0.658        | 0.577        | 0.644        | 0.573        | 0.726        |
|              | MEERKAT-VP | 10         | 0.941        | <b>0.886</b> | <b>0.947</b> | <b>0.76</b>  | <b>0.822</b> | 0.653        | <b>0.636</b> | <b>0.806</b> |
|              | MEERKAT    | 10         | <b>0.949</b> | 0.881        | 0.934        | 0.752        | 0.813        | <b>0.682</b> | 0.628        | 0.805        |
| Qwen2-1.5b   | MEERKAT-VP | 30         | 0.935        | 0.876        | <b>0.953</b> | <b>0.759</b> | <b>0.822</b> | 0.653        | <b>0.626</b> | <b>0.803</b> |
|              | MEERKAT    | 30         | <b>0.944</b> | <b>0.878</b> | 0.928        | 0.734        | 0.800        | <b>0.663</b> | 0.624        | 0.795        |
|              | MEERKAT-VP | 50         | 0.931        | <b>0.882</b> | <b>0.946</b> | <b>0.754</b> | <b>0.804</b> | 0.644        | <b>0.63</b>  | <b>0.798</b> |
|              | MEERKAT    | 50         | <b>0.948</b> | 0.872        | 0.926        | 0.746        | 0.795        | <b>0.663</b> | 0.594        | 0.792        |
|              | MEERKAT-VP | 100        | 0.935        | 0.874        | <b>0.947</b> | <b>0.751</b> | <b>0.817</b> | 0.653        | <b>0.644</b> | <b>0.803</b> |
|              | MEERKAT    | 100        | <b>0.936</b> | <b>0.878</b> | 0.925        | 0.741        | 0.795        | <b>0.663</b> | 0.61         | 0.792        |
|              | MEERKAT-VP | 10         | <b>0.948</b> | <b>0.873</b> | <b>0.971</b> | 0.802        | <b>0.657</b> | <b>0.663</b> | 0.609        | <b>0.789</b> |
|              | MEERKAT    | 10         | 0.939        | 0.869        | 0.96         | <b>0.804</b> | 0.591        | 0.634        | 0.609        | 0.772        |
|              | MEERKAT-VP | 30         | <b>0.948</b> | <b>0.86</b>  | <b>0.974</b> | <b>0.799</b> | <b>0.6</b>   | 0.634        | <b>0.619</b> | <b>0.776</b> |
|              | MEERKAT    | 30         | 0.94         | 0.855        | 0.947        | 0.734        | 0.568        | <b>0.644</b> | 0.601        | 0.755        |
| Gemma2-2b    | MEERKAT-VP | 50         | <b>0.949</b> | 0.853        | <b>0.969</b> | <b>0.782</b> | 0.551        | 0.615        | 0.620        | 0.762        |
|              | MEERKAT    | 50         | 0.945        | <b>0.857</b> | 0.966        | 0.767        | <b>0.613</b> | <b>0.634</b> | <b>0.623</b> | <b>0.772</b> |
|              | MEERKAT-VP | 100        | <b>0.944</b> | 0.812        | <b>0.97</b>  | 0.733        | 0.551        | 0.634        | <b>0.634</b> | <b>0.754</b> |
|              | MEERKAT    | 100        | 0.94         | <b>0.851</b> | 0.951        | <b>0.745</b> | 0.551        | 0.634        | 0.574        | 0.749        |

Table 15: MEERKAT performance at local step = 1 with varying outlier percentages across the LLaMA-3.2-1B, Qwen2-1.5b, and Gemma2-2b models. We report test accuracy on SST-2, AgNews, Yelp, BoolQ, RTE, WSC, and WIC under both IID and Non-IID client data settings. Bold numbers indicate the highest value in each row.

| Model        | Outlier Percentage | IID          |              |              |              |              |              | Non-IID      |              |              |             |              |              |              |              |
|--------------|--------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|-------------|--------------|--------------|--------------|--------------|
|              |                    | SST-2        | AgNews       | Yelp         | BoolQ        | RTE          | WSC          | WIC          | SST-2        | AgNews       | Yelp        | BoolQ        | RTE          | WSC          | WIC          |
| LLaMA-3.2-1B | 5e-1               | 0.917        | 0.72         | 0.965        | 0.725        | 0.653        | 0.644        | 0.634        | 0.895        | 0.669        | 0.964       | 0.684        | 0.644        | 0.653        | 0.594        |
|              | 5e-2               | 0.913        | 0.861        | 0.966        | 0.749        | 0.653        | 0.644        | 0.633        | 0.915        | 0.87         | <b>0.97</b> | 0.722        | 0.653        | 0.644        | 0.619        |
|              | 5e-3               | 0.900        | <b>0.885</b> | <b>0.971</b> | 0.769        | 0.702        | 0.653        | 0.614        | <b>0.930</b> | 0.874        | 0.963       | <b>0.753</b> | 0.620        | 0.66         | 0.62         |
|              | 5e-4               | 0.910        | 0.877        | 0.954        | <b>0.773</b> | <b>0.720</b> | <b>0.663</b> | 0.641        | 0.911        | <b>0.888</b> | 0.956       | 0.700        | <b>0.693</b> | <b>0.663</b> | <b>0.628</b> |
|              | 5e-5               | <b>0.922</b> | 0.879        | 0.964        | 0.724        | 0.631        | 0.625        | <b>0.648</b> | 0.92         | 0.876        | 0.940       | 0.725        | 0.613        | <b>0.663</b> | 0.626        |
| Qwen2-1.5b   | 5e-1               | 0.854        | 0.856        | 0.947        | 0.766        | 0.82         | 0.663        | 0.644        | 0.845        | 0.854        | 0.946       | 0.753        | 0.826        | 0.682        | 0.631        |
|              | 5e-2               | 0.925        | 0.868        | 0.949        | 0.778        | 0.826        | 0.692        | 0.647        | 0.93         | 0.853        | 0.943       | 0.759        | 0.822        | 0.663        | 0.663        |
|              | 5e-3               | <b>0.926</b> | 0.851        | 0.945        | 0.765        | <b>0.813</b> | <b>0.692</b> | <b>0.658</b> | 0.924        | <b>0.866</b> | 0.94        | 0.759        | 0.822        | <b>0.692</b> | 0.661        |
|              | 5e-4               | 0.92         | 0.764        | 0.943        | 0.774        | 0.813        | 0.682        | 0.645        | 0.918        | 0.848        | 0.943       | <b>0.762</b> | 0.813        | 0.682        | 0.647        |
|              | 5e-5               | 0.903        | 0.78         | 0.941        | 0.748        | 0.80         | 0.673        | 0.625        | 0.896        | 0.799        | 0.937       | 0.739        | 0.80         | 0.673        | 0.633        |
| Gemma2-2b    | 5e-1               | 0.842        | 0.867        | 0.963        | 0.751        | 0.657        | <b>0.673</b> | 0.626        | 0.871        | 0.855        | 0.952       | 0.695        | <b>0.663</b> | <b>0.663</b> | 0.619        |
|              | 5e-2               | 0.932        | <b>0.878</b> | <b>0.977</b> | 0.809        | 0.791        | 0.663        | 0.623        | 0.92         | <b>0.863</b> | 0.968       | 0.786        | 0.706        | 0.653        | 0.634        |
|              | 5e-3               | <b>0.952</b> | 0.871        | 0.971        | <b>0.837</b> | <b>0.800</b> | 0.663        | <b>0.639</b> | <b>0.942</b> | 0.853        | <b>0.97</b> | 0.807        | <b>0.751</b> | 0.653        | <b>0.645</b> |
|              | 5e-4               | 0.941        | 0.824        | 0.967        | 0.83         | 0.764        | 0.663        | 0.612        | 0.941        | 0.83         | 0.962       | <b>0.831</b> | 0.746        | 0.634        | 0.63         |
|              | 5e-5               | 0.92         | 0.828        | 0.952        | 0.797        | 0.6          | 0.634        | 0.606        | 0.922        | 0.764        | 0.949       | 0.774        | 0.56         | 0.634        | 0.601        |

Table 16: Performance comparison of Full-FedZO, LoRA-FedZO, and MEERKAT under synchronous updates with  $localstep = 1$ , evaluated on both IID and Non-IID client data settings(**Dirichlet**  $\alpha = 0.5$ ) across LLaMA-3.2-1B, Qwen2-1.5b, and Gemma2-2b. We report test accuracy on SST-2, AgNews, Yelp, BoolQ, RTE, WSC, and WIC. Bold numbers indicate the highest value in each row.

| Model                  | Method     | SST-2        | AgNews       | Yelp         | BoolQ        | RTE          | WSC          | WIC          | Acc          |
|------------------------|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| LLaMA-3.2-1B (IID)     | Full-FedZO | <b>0.918</b> | 0.801        | 0.937        | 0.686        | 0.54         | 0.625        | 0.58         | 0.726        |
|                        | LoRA-FedZO | 0.915        | 0.855        | 0.944        | 0.672        | 0.599        | <b>0.663</b> | 0.599        | 0.749        |
|                        | MEERKAT    | 0.900        | <b>0.885</b> | <b>0.971</b> | <b>0.773</b> | <b>0.702</b> | 0.653        | <b>0.614</b> | <b>0.785</b> |
| LLaMA-3.2-1B (Non-IID) | Full-FedZO | 0.911        | 0.831        | 0.937        | 0.672        | 0.528        | 0.587        | 0.567        | 0.719        |
|                        | LoRA-FedZO | 0.8669       | 0.842        | 0.944        | 0.659        | 0.53         | 0.567        | 0.578        | 0.712        |
|                        | MEERKAT    | <b>0.93</b>  | <b>0.888</b> | <b>0.963</b> | <b>0.753</b> | <b>0.67</b>  | <b>0.66</b>  | <b>0.62</b>  | <b>0.783</b> |
| Qwen2-1.5b (IID)       | Full-FedZO | 0.9013       | 0.726        | 0.918        | 0.700        | 0.797        | <b>0.710</b> | 0.579        | 0.761        |
|                        | LoRA-FedZO | <b>0.935</b> | 0.752        | 0.925        | 0.686        | 0.794        | 0.673        | 0.606        | 0.767        |
|                        | MEERKAT    | 0.926        | <b>0.851</b> | <b>0.945</b> | <b>0.778</b> | <b>0.813</b> | 0.692        | <b>0.658</b> | <b>0.809</b> |
| Qwen2-1.5b (Non-IID)   | Full-FedZO | 0.844        | 0.725        | 0.937        | 0.688        | 0.769        | 0.663        | 0.565        | 0.741        |
|                        | LoRA-FedZO | <b>0.932</b> | 0.76         | <b>0.944</b> | 0.682        | 0.773        | 0.682        | 0.565        | 0.763        |
|                        | MEERKAT    | 0.924        | <b>0.866</b> | 0.94         | <b>0.762</b> | <b>0.822</b> | <b>0.692</b> | <b>0.661</b> | <b>0.809</b> |
| Gemma2-2b (IID)        | Full-FedZO | 0.934        | 0.84         | 0.953        | 0.774        | 0.542        | 0.644        | 0.606        | 0.756        |
|                        | LoRA-FedZO | 0.942        | 0.856        | 0.94         | 0.735        | 0.52         | 0.644        | 0.606        | 0.749        |
|                        | MEERKAT    | <b>0.952</b> | <b>0.871</b> | <b>0.971</b> | <b>0.837</b> | <b>0.8</b>   | <b>0.663</b> | <b>0.639</b> | <b>0.819</b> |
| Gemma2-2b (Non-IID)    | Full-FedZO | 0.93         | 0.824        | 0.95         | 0.744        | 0.56         | 0.625        | 0.575        | 0.744        |
|                        | LoRA-FedZO | 0.9415       | 0.825        | 0.954        | 0.711        | 0.528        | 0.625        | 0.578        | 0.737        |
|                        | MEERKAT    | <b>0.942</b> | <b>0.853</b> | <b>0.97</b>  | <b>0.807</b> | <b>0.751</b> | <b>0.653</b> | <b>0.645</b> | <b>0.803</b> |

2430  
 2431 Table 17: Performance comparison of LoRA-FedZO, and MEERKAT under synchronous updates  
 2432 with  $localstep = 1$ , evaluated on Non-IID client data settings (**Dirichlet**  $\alpha = 0.3$ ) across  
 2433 LLaMA-3.2-1B, Qwen2-1.5b, and Gemma2-2b. We report test accuracy on SST-2, AgNews, Yelp,  
 2434 BoolQ, RTE, WSC, and WIC. Bold numbers indicate the highest value in each row.  
 2435

| Model                  | Method     | SST-2        | AgNews       | Yelp         | BoolQ        | RTE          | WSC          | WIC          | Acc          |
|------------------------|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| LLaMA-3.2-1B (Non-IID) | Full-FedZO | 0.891        | 0.759        | 0.94         | 0.623        | 0.528        | 0.644        | 0.551        | 0.705        |
|                        | LoRA-FedZO | 0.915        | 0.866        | 0.952        | 0.646        | 0.586        | 0.653        | 0.554        | 0.739        |
|                        | MEERKAT    | <b>0.918</b> | <b>0.843</b> | <b>0.97</b>  | <b>0.761</b> | <b>0.626</b> | <b>0.653</b> | <b>0.609</b> | <b>0.769</b> |
| Qwen2-1.5b (Non-IID)   | Full-FedZO | 0.52         | 0.347        | 0.45         | 0.62         | 0.532        | 0.632        | 0.51         | 0.516        |
|                        | LoRA-FedZO | 0.855        | 0.732        | 0.907        | 0.674        | 0.72         | 0.634        | 0.603        | 0.732        |
|                        | MEERKAT    | <b>0.91</b>  | <b>0.809</b> | <b>0.954</b> | <b>0.772</b> | <b>0.822</b> | <b>0.682</b> | <b>0.661</b> | <b>0.801</b> |
| Gemma2-2b (Non-IID)    | Full-FedZO | 0.881        | 0.761        | 0.94         | 0.688        | 0.552        | 0.613        | 0.603        | 0.720        |
|                        | LoRA-FedZO | 0.922        | 0.826        | 0.921        | 0.681        | 0.52         | 0.625        | 0.606        | 0.729        |
|                        | MEERKAT    | <b>0.942</b> | <b>0.873</b> | <b>0.97</b>  | <b>0.806</b> | <b>0.688</b> | <b>0.634</b> | <b>0.615</b> | <b>0.79</b>  |

2443  
 2444 Table 18: Performance comparison of LoRA-FedZO, and MEERKAT under synchronous updates  
 2445 with  $localstep = 1$ , evaluated on Non-IID client data settings (**Dirichlet**  $\alpha = 0.1$ ) across  
 2446 LLaMA-3.2-1B, Qwen2-1.5b, and Gemma2-2b. We report test accuracy on SST-2, AgNews, Yelp,  
 2447 BoolQ, RTE, WSC, and WIC. Bold numbers indicate the highest value in each row.  
 2448

| Model                  | Method     | SST-2        | AgNews       | Yelp         | BoolQ        | RTE          | WSC          | WIC          | Acc          |
|------------------------|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| LLaMA-3.2-1B (Non-IID) | Full-FedZO | 0.891        | 0.754        | 0.933        | 0.626        | 0.522        | 0.365        | 0.512        | 0.658        |
|                        | LoRA-FedZO | 0.902        | 0.845        | 0.942        | 0.643        | 0.533        | 0.365        | 0.559        | 0.684        |
|                        | MEERKAT    | <b>0.92</b>  | <b>0.794</b> | <b>0.965</b> | <b>0.745</b> | <b>0.582</b> | <b>0.644</b> | <b>0.603</b> | <b>0.750</b> |
| Qwen2-1.5b (Non-IID)   | Full-FedZO | 0.49         | 0.247        | 0.44         | 0.62         | 0.528        | 0.634        | 0.5          | 0.494        |
|                        | LoRA-FedZO | 0.848        | 0.735        | 0.92         | 0.67         | 0.746        | 0.548        | 0.601        | 0.724        |
|                        | MEERKAT    | <b>0.889</b> | <b>0.78</b>  | <b>0.944</b> | <b>0.732</b> | <b>0.822</b> | <b>0.634</b> | <b>0.637</b> | <b>0.777</b> |
| Gemma2-2b (Non-IID)    | Full-FedZO | 0.879        | 0.741        | 0.937        | 0.681        | 0.48         | 0.634        | 0.601        | 0.708        |
|                        | LoRA-FedZO | 0.91         | 0.78         | 0.914        | 0.682        | 0.551        | 0.567        | 0.608        | 0.716        |
|                        | MEERKAT    | <b>0.944</b> | <b>0.866</b> | <b>0.971</b> | <b>0.805</b> | <b>0.728</b> | <b>0.605</b> | <b>0.628</b> | <b>0.792</b> |

2459  
 2460 Table 19: Test accuracy of MEERKAT versus DecomFL on Qwen2-1.5b with a single local step under  
 2461 Non-IID data settings (Dirichlet  $\alpha = 1$ ). Results are shown for SST-2, BoolQ, RTE, and WSC; bold  
 2462 indicates the best score in each row. Experiments use 8 clients in total, with 2 clients participating in  
 2463 each round, following the DecomFL configuration.  
 2464

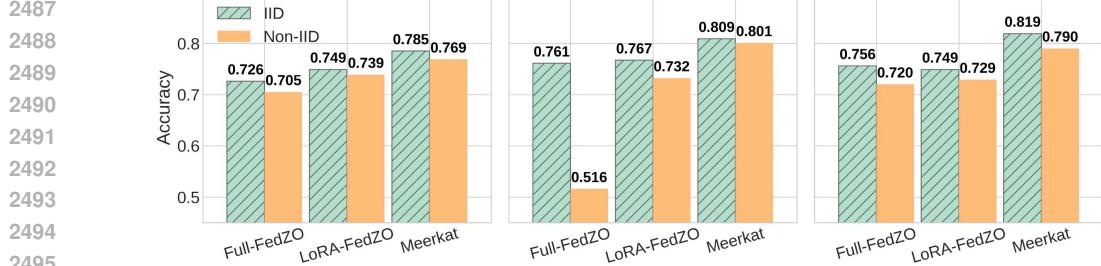
| Model      | Method  | SST-2        | BoolQ        | RTE          | WSC          |
|------------|---------|--------------|--------------|--------------|--------------|
| Qwen2-1.5b | DecomFL | 0.868        | 0.674        | 0.773        | 0.653        |
|            | MEERKAT | <b>0.918</b> | <b>0.734</b> | <b>0.817</b> | <b>0.682</b> |

2472 Table 20: Performance comparison of Task-Mask, and MEERKAT under synchronous updates with  
 2473  $localstep = 1$ , evaluated on IID client data settings across LLaMA-3.2-1B, Qwen2-1.5b, and  
 2474 Gemma2-2b. We report test accuracy on SST-2, AgNews, Yelp, BoolQ, RTE, WSC, and WIC. Bold  
 2475 numbers indicate the highest value in each row.  
 2476

| Model              | Method  | SST-2        | AgNews       | Yelp         | BoolQ        | RTE          | WSC          | WIC          |
|--------------------|---------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| LLaMA-3.2-1B (IID) | Task    | <b>0.910</b> | 0.847        | 0.957        | 0.718        | 0.661        | 0.644        | <b>0.661</b> |
|                    | MEERKAT | 0.90         | <b>0.885</b> | <b>0.971</b> | <b>0.773</b> | <b>0.702</b> | <b>0.653</b> | 0.614        |
| Qwen2-1.5b (IID)   | Task    | 0.936        | 0.827        | 0.954        | 0.765        | 0.83         | <b>0.711</b> | <b>0.664</b> |
|                    | MEERKAT | 0.926        | <b>0.851</b> | 0.945        | <b>0.778</b> | <b>0.813</b> | 0.692        | 0.658        |
| Gemma2-2b (IID)    | Task    | 0.942        | 0.868        | <b>0.972</b> | 0.78         | 0.728        | 0.644        | 0.6          |
|                    | MEERKAT | <b>0.952</b> | <b>0.871</b> | 0.971        | <b>0.837</b> | <b>0.8</b>   | <b>0.663</b> | <b>0.639</b> |

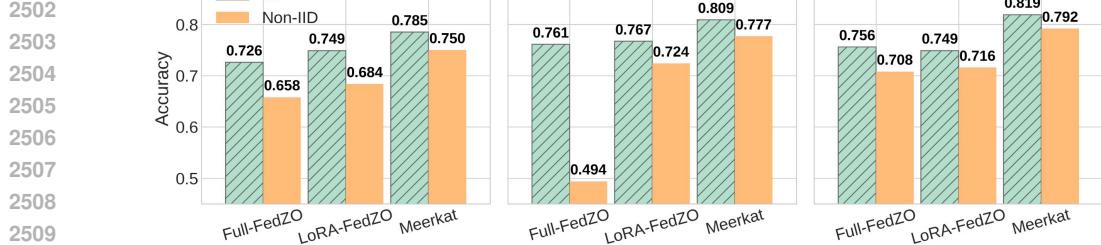
2484

2485

2486 **LLaMA-3.2-1B**

2487 (a) This figure compares three methods—Full-FedZO, LoRA-FedZO, and MEERKAT—on three  
 2488 LLMs: LLaMA-3.2-1B, Qwen2-1.5b, and Gemma2-2b. The x-axis shows the different methods, and  
 2489 each method has two bars indicating performance under IID and Non-IID settings. The Non-IID  
 2490 results are obtained under a Dirichlet  $\alpha = 0.3$ . The y-axis represents the average test accuracy across  
 2491 multiple downstream tasks—SST2, AgNews, Yelp, BoolQ, RTE, WSC, and WiC.

2492

2493 **Qwen2-1.5b**

2494

2495

2496 (b) This figure compares three methods—Full-FedZO, LoRA-FedZO, and MEERKAT—on three  
 2497 LLMs: LLaMA-3.2-1B, Qwen2-1.5b, and Gemma2-2b. The x-axis shows the different methods, and  
 2498 each method has two bars indicating performance under IID and Non-IID settings. The Non-IID  
 2499 results are obtained under a Dirichlet  $\alpha = 0.1$ . The y-axis represents the average test accuracy across  
 2500 multiple downstream tasks—SST2, AgNews, Yelp, BoolQ, RTE, WSC, and WiC.

2501

2502

2503 Figure 6: Comparison of Full-FedZO, LoRA-FedZO, and MEERKAT on LLaMA-3.2-1B, Qwen2-  
 2504 1.5b, and Gemma2-2b under IID and Non-IID settings with varying Dirichlet  $\alpha$ . Subfigure(a) presents  
 2505 results for Non-IID data generated with  $\alpha = 0.3$ , while Subfigure(b) shows results for Non-IID data  
 2506 with  $\alpha = 0.1$ .

2507

2508

2509

2510

2511

2512

2513

2514

2515

2516

2517

2518

2519

2520

2521

2522

2523

2524

2525

2526

2527

2528

2529

2530

2531

2532

2533

2534

2535

2536

2537

2538 Table 21: Performance comparison of Task-Mask, which uses downstream task data to select  
 2539 sensitive model parameters, and MEERKAT under synchronous updates with  $localstep = 1$ ,  
 2540 evaluated on Non-IID client data settings (**Dirichlet**  $\alpha = 0.5$ ) across LLaMA-3.2-1B, Qwen2-1.5b,  
 2541 and Gemma2-2b. We report test accuracy on SST-2, AgNews, Yelp, BoolQ, RTE, WSC, and WIC.  
 2542 Bold numbers indicate the highest value in each row.

| Model                         | Method  | SST-2        | AgNews       | Yelp         | BoolQ        | RTE          | WSC          | WIC          |
|-------------------------------|---------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| <b>LLaMA-3.2-1B (Non-IID)</b> | Task    | 0.904        | 0.874        | 0.956        | 0.744        | 0.591        | 0.615        | 0.622        |
|                               | MEERKAT | <b>0.93</b>  | <b>0.888</b> | <b>0.963</b> | <b>0.753</b> | <b>0.62</b>  | <b>0.66</b>  | <b>0.62</b>  |
| <b>Qwen2-1.5b (Non-IID)</b>   | Task    | 0.938        | 0.863        | 0.956        | 0.779        | 0.817        | 0.692        | 0.65         |
|                               | MEERKAT | <b>0.924</b> | <b>0.866</b> | <b>0.94</b>  | <b>0.762</b> | <b>0.822</b> | <b>0.692</b> | <b>0.661</b> |
| <b>Gemma2-2b (Non-IID)</b>    | Task    | 0.91         | 0.834        | 0.966        | 0.822        | 0.72         | 0.644        | 0.578        |
|                               | MEERKAT | <b>0.942</b> | <b>0.853</b> | <b>0.97</b>  | <b>0.807</b> | <b>0.751</b> | <b>0.653</b> | <b>0.645</b> |

Table 22: Test accuracy of MEERKAT versus Task-Mask on Qwen2-1.5b with a 10 local step under Non-IID data settings (Dirichlet  $\alpha = 0.5$ ). Results are shown for SST-2, BoolQ, RTE, and WSC; bold indicates the best score in each row. Experiments use 8 clients in total, with 2 clients participating in each round, following the DecomFL configuration.

| Model      | Method  | SST-2        | BoolQ        | RTE          | WSC          |
|------------|---------|--------------|--------------|--------------|--------------|
| Qwen2-1.5b | Task    | 0.932        | 0.784        | 0.823        | 0.681        |
|            | MEERKAT | <b>0.944</b> | <b>0.752</b> | <b>0.813</b> | <b>0.682</b> |

Table 23: MEERKAT Convergence Rounds for the LLaMA-3.2-1B, Gemma2-2B, and Qwen2-1.5B models on the SST-2, AgNews, Yelp, and BoolQ tasks, with 10 local steps.

| Model        | SST-2 | AgNews | Yelp | BoolQ |
|--------------|-------|--------|------|-------|
| Gemma2-2B    | 39    | 61     | 29   | 43    |
| Qwen2-1.5B   | 51    | 75     | 36   | 70    |
| LLaMA-3.2-1B | 85    | 77     | 52   | 97    |

Table 24: Computation and Communication Efficiency Benchmark Shows MEERKAT’s Superior Resource Usage over Baselines. We benchmarked resource usage on Qwen2-1.5B with 10 clients (FP16). Setting: Full-FedZO vs Meerkat vs LoRA-FedZO, where LoRA is configured with rank = 16,  $\alpha = 16$ —the same setting used in Table 1.

| Method/Metrics      | RAM (Peak)       | Upload/Client | Download/Client |
|---------------------|------------------|---------------|-----------------|
| Full-FedZO          | 12,600 MiB       | 0.078 KB      | 2.875 GB        |
| LoRA-FedZO          | 10,741 MiB       | 0.078 KB      | 35.22 MB        |
| MEERKAT (0.1% mask) | <b>7,850 MiB</b> | 0.078 KB      | <b>2.50 MB</b>  |

Table 25: Performance comparison on LLaMA-3.2-1B under Non-IID Dirichlet partition ( $\alpha = 0.5$ ) with  $T = 10$  local steps. While ZO methods cannot match back-propagation due to gradient noise from limited sampling, MEERKAT-VP achieves competitive accuracy (0.764 avg) with significantly lower memory consumption, and outperforms several Non-IID FL baselines (FedDYN, FedAvgM, FedSA-LoRA, and stochastic controlled averaging) under the same training setup. We adapt FedDYN following the original paper with  $\alpha = 0.01$ .

| Method                | SST-2 | AGNews | Yelp  | BoolQ | RTE   | WSC   | WIC   | Avg   |
|-----------------------|-------|--------|-------|-------|-------|-------|-------|-------|
| Back-propagation      | 0.925 | 0.893  | 0.968 | 0.751 | 0.644 | 0.660 | 0.630 | 0.782 |
| Stochastic Controlled | 0.880 | 0.720  | 0.901 | 0.612 | 0.523 | 0.612 | 0.580 | 0.690 |
| FedAvgM               | 0.901 | 0.821  | 0.941 | 0.629 | 0.580 | 0.613 | 0.600 | 0.726 |
| FedSA-LoRA            | 0.905 | 0.832  | 0.920 | 0.630 | 0.570 | 0.622 | 0.570 | 0.721 |
| MEERKAT+FedDYN        | 0.917 | 0.841  | 0.954 | 0.638 | 0.564 | 0.615 | 0.570 | 0.728 |
| MEERKAT-VP            | 0.922 | 0.864  | 0.962 | 0.713 | 0.617 | 0.644 | 0.625 | 0.764 |

Table 26: Effect of different sparsity density ratios on LLaMA-3-8B under Non-IID Dirichlet  $\alpha = 0.5$ . We compare density ratios  $10^{-3}$  and  $10^{-4}$  using the same transferable mask construction pipeline. Both settings achieve strong performance.

| Density ratio      | SST-2 | AGNews | Yelp  | BoolQ | RTE   | WSC   | WIC   | Avg   |
|--------------------|-------|--------|-------|-------|-------|-------|-------|-------|
| $1 \times 10^{-3}$ | 0.950 | 0.851  | 0.954 | 0.831 | 0.755 | 0.674 | 0.660 | 0.811 |
| $1 \times 10^{-4}$ | 0.941 | 0.862  | 0.956 | 0.861 | 0.783 | 0.664 | 0.640 | 0.815 |

Table 27: Scalability of MEERKAT and MEERKAT-VP with respect to the number of clients on Qwen2-1.5B. Increasing the number of clients from 10 to 20 does not degrade performance; MEERKAT-VP-20 even slightly improves the average accuracy.

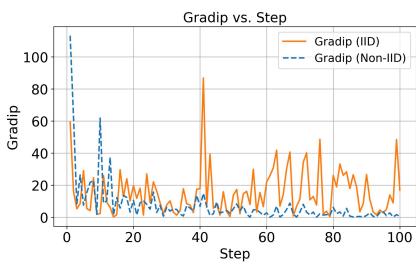
| Model                   | SST2  | AgNews | Yelp  | BoolQ |
|-------------------------|-------|--------|-------|-------|
| MEERKAT-VP (20 clients) | 0.951 | 0.885  | 0.936 | 0.756 |
| MEERKAT (20 clients)    | 0.929 | 0.869  | 0.922 | 0.719 |
| MEERKAT (10 clients)    | 0.949 | 0.881  | 0.934 | 0.752 |

2592

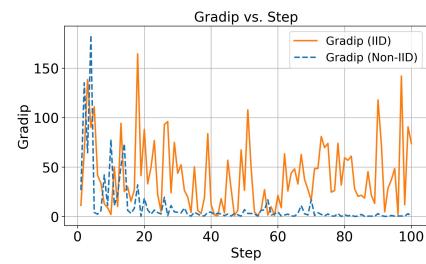
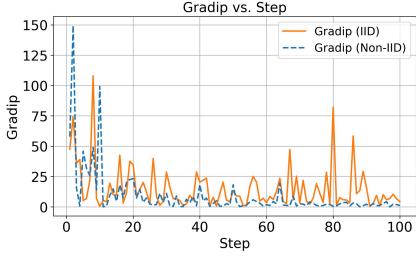
2593

2594

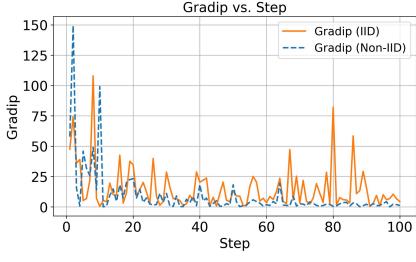
2595



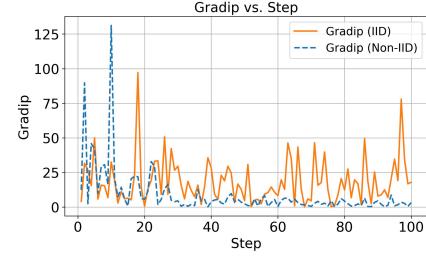
(a) The GradIP measured for IID and Non-IID clients data under the WIC task using the Llama-3.2-1B model.



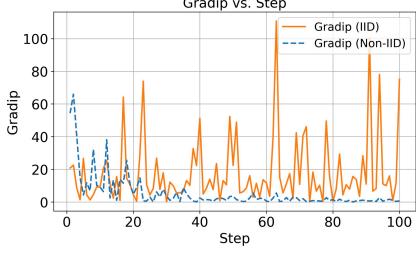
(b) The GradIP measured for IID and Non-IID clients data under the AgNews task using the Llama-3.2-1B model.



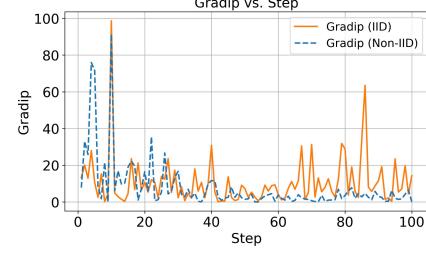
(c) The GradIP measured for IID and Non-IID clients data under the Yelp task using the Llama-3.2-1B model.



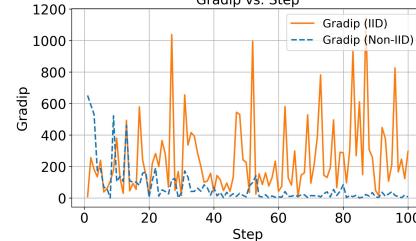
(d) The GradIP measured for IID and Non-IID clients data under the BoolQ task using the Llama-3.2-1B model.



(e) The GradIP measured for IID and Non-IID clients data under the RTE task using the Llama-3.2-1B model.



(f) The GradIP measured for IID and Non-IID clients data under the WSC task using the Llama-3.2-1B model.



(g) The GradIP measured for IID and Non-IID clients data under the BoolQ task using the Gemma-2-2b model.

Figure 7: These figures show GradIP (Definition 2.3) curves under IID and Non-IID settings, computed over 100 local training steps on six datasets (WSC, BoolQ, RTE, WIC, AgNews, Yelp) using the Llama-3.2-1B model with density level  $5 \times 10^{-3}$ . An extra BoolQ result is shown for the Gemma-2-2B model.

2639

2640

2641

2642

2643

2644

2645

2646  
2647  
2648  
2649  
2650

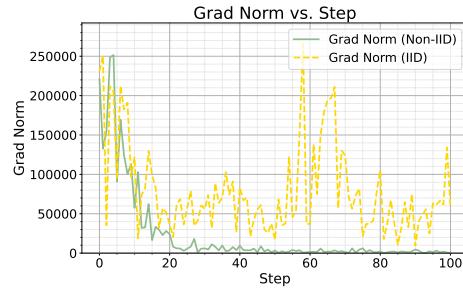
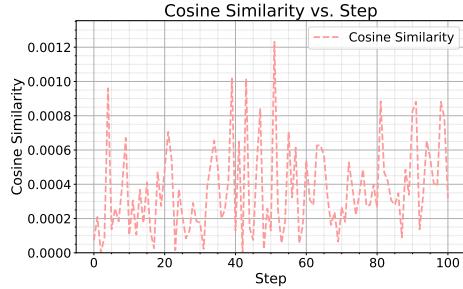


Figure 8: The left panel shows the cosine similarity between locally computed ZO gradients and gradients from the C4-pre-trained data, illustrating that the two gradient vectors remain nearly orthogonal throughout training. The right panel presents the norm of local ZO gradients over training steps, showing a consistent decay and convergence in magnitude under Non-IID and IID data distribution. These observations are obtained under density level of  $5 \times 10^{-3}$ .

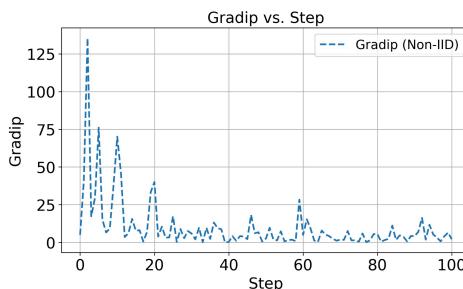
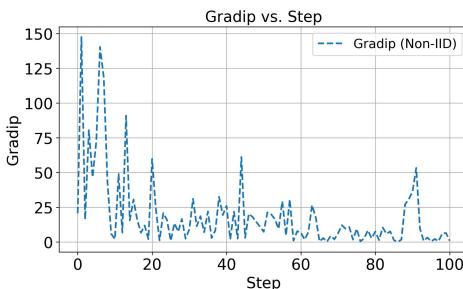
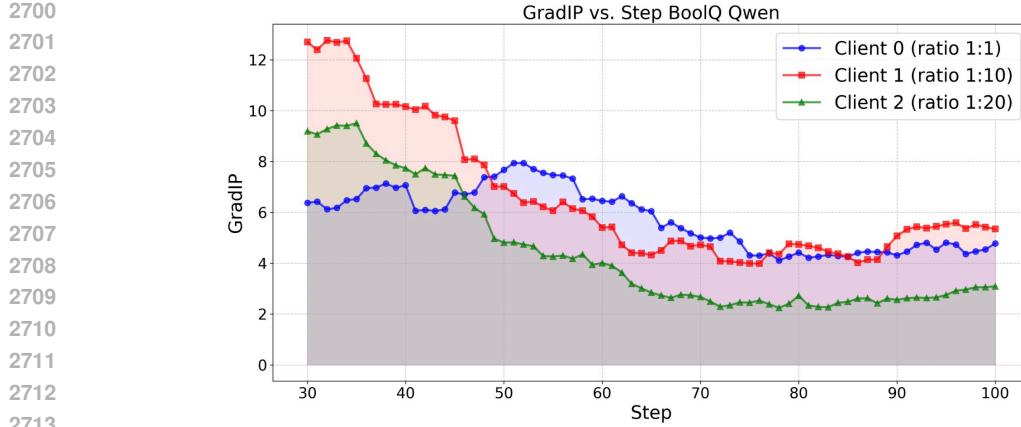
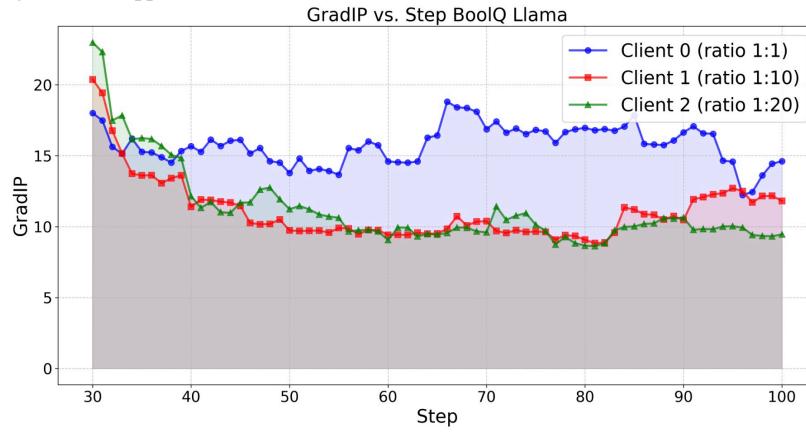


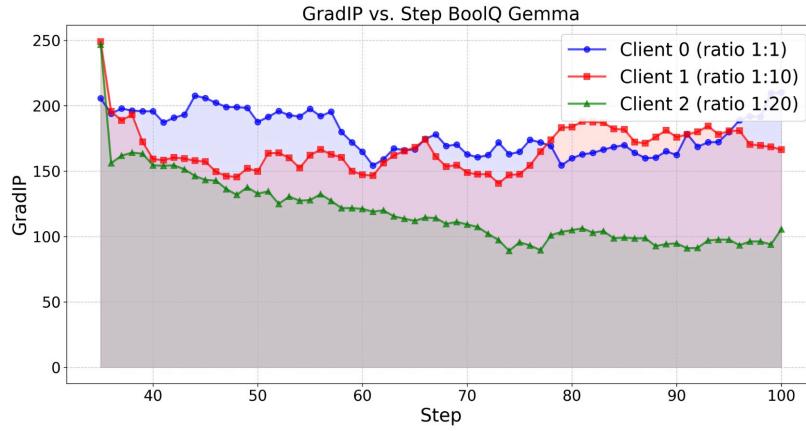
Figure 9: These subfigures show GradIP (see Definition 2.3) for LLaMA-3.2-1B under Non-IID client data with 100 local training steps. Subfigure (a) uses AgNews (5 vs. 89), while Subfigure (b) uses BoolQ (6 vs. 190).



2714 (a) The experiments, conducted using the Qwen2-1.5B model on the BoolQ dataset, reveal that under  
2715 Non-IID settings—especially with a 1:20 class imbalance—there is a pronounced decline in GradIP  
2716 between the early and later stages of training. In the extreme Non-IID case, the GradIP values in the  
2717 later stages tend to approach zero.



2731 (b) The experiments, conducted using the Llama-3.2-1B model on the BoolQ dataset, reveal that  
2732 under Non-IID settings—especially with a 1:20 class imbalance—there is a pronounced decline in  
2733 GradIP between the early and later stages of training. In the extreme Non-IID case, the GradIP values  
2734 in the later stages tend to approach zero.



2748 (c) The experiments, conducted using the Gemma-2-2B model on the BoolQ dataset, reveal that  
2749 under Non-IID settings—especially with a 1:20 class imbalance—there is a pronounced decline in  
2750 GradIP between the early and later stages of training. In the extreme Non-IID case, the GradIP values  
2751 in the later stages tend to approach zero.

2752 Figure 10: GradIP analysis for different models on the BoolQ dataset under Non-IID and IID  
2753 conditions: As the class imbalance ratio increases, GradIP in the later training stages tends to  
approach zero. This decline is more pronounced under Non-IID settings, where the gap between  
initial and final GradIP values is larger than in the IID case. All trends are visualized using a moving  
average for clarity.

2754

2755

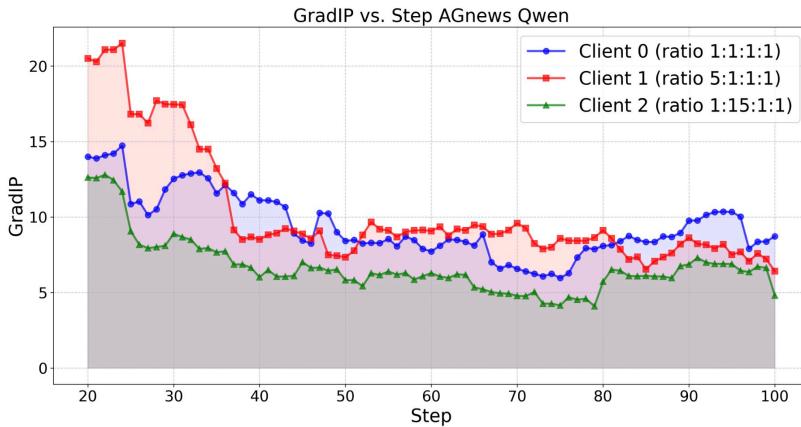
2756

2757

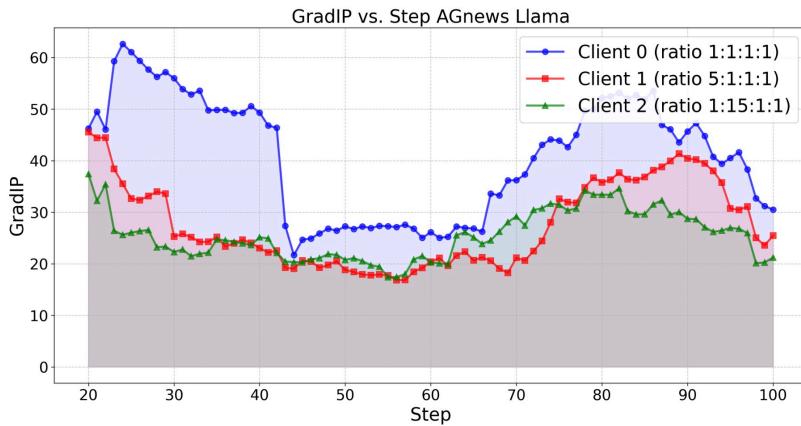
2758

2759

2760



2774 (a) The experiments, conducted using the Qwen2-1.5B model on the AGNews dataset, reveal that  
 2775 under Non-IID settings—especially with a 1:15:1:1 class imbalance—there is a pronounced decline  
 2776 in GradIP between the early and later stages of training. In the extreme Non-IID case, the GradIP  
 2777 values in the later stages tend to approach zero.



2791 (b) The experiments, conducted using the Llama-3.2-1B model on the AGNews dataset, reveal that  
 2792 under Non-IID settings—especially with a 1:15:1:1 class imbalance—there is a pronounced decline  
 2793 in GradIP between the early and later stages of training. In the extreme Non-IID case, the GradIP  
 2794 values in the later stages tend to approach zero.

2795

2796

2797

2798

2799

2800

2801

2802

2803

2804

2805

2806

2807

Figure 11: GradIP analysis for different models on the AGNews dataset under Non-IID and IID conditions: As the class imbalance ratio increases, GradIP in the later training stages tends to approach zero. This decline is more pronounced under Non-IID settings, where the gap between initial and final GradIP values is larger than in the IID case. All trends are visualized using a moving average for clarity; consequently, the plotted lines do not begin at step zero, as the initial data points are used to compute the first averaged value. This is an intentional effect of the visualization, not an error or a result of missing data.

2808

2809

2810

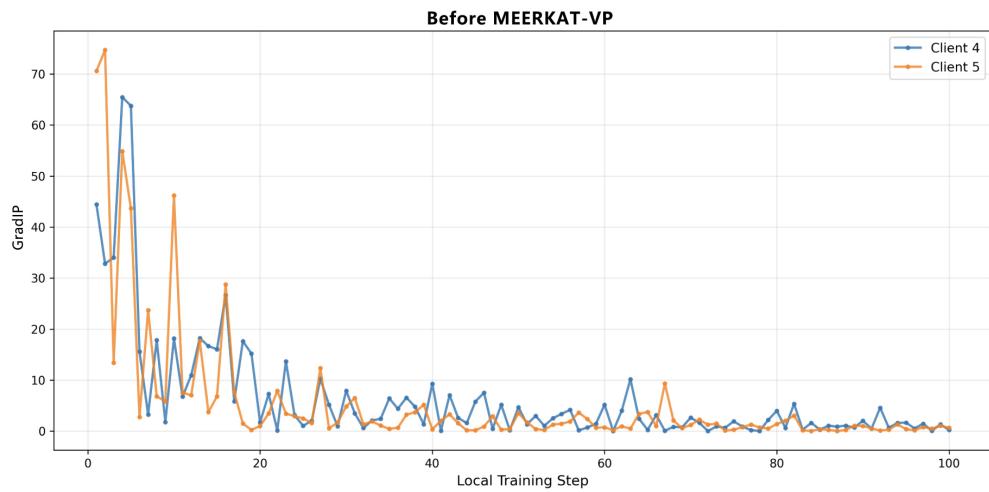
2811

2812

2813

2814

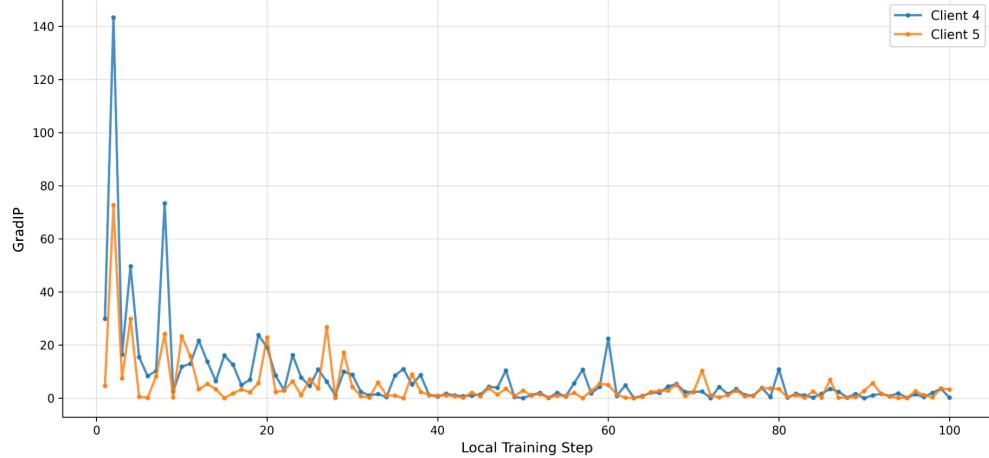
2815



(a) GradIP trajectories for the two extreme Non-IID clients on the SST2 task with Qwen2-1.5B, before MEERKAT-VP training.

2832

**MEERKAT-VP Managed**



(b) GradIP trajectories for the same two extreme Non-IID clients on SST2 *after MEERKAT-VP training has converged* (global validation accuracy  $\approx 90\%$ ).

2850

Figure 12: GradIP trajectories (Definition 2.3) for Qwen2-1.5B on SST2 with 6 clients (2 extreme Non-IID, 4 IID). Subfigure (a) shows the GradIP trajectories of the two extreme Non-IID clients at initialization, while Subfigure (b) shows the trajectories for the same clients after MEERKAT-VP training has converged. The shape of the trajectories remains similar before and after training, supporting our claim that GradIP is primarily a data-distribution-driven signal rather than a direct reflection of the global model state.

2856

2857

2858

2859

2860

2861