ENERGY AND MEMORY-EFFICIENT FEDERATED LEARNING WITH ORDERED LAYER FREEZING AND TENSOR OPERATION APPROXIMATION

Anonymous authors

006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026 027 028

029

Paper under double-blind review

ABSTRACT

The effectiveness of Federated Learning (FL) in the context of the Internet of Things (IoT) is hindered by the resource constraints of IoT devices, such as limited computing capability, memory space and bandwidth support. These constraints create significant computation and communication bottlenecks for training and transmitting deep neural networks. Various FL frameworks have been proposed to reduce computation and communication overheads through dropout or layer freezing. However, these approaches often sacrifice accuracy or neglect memory constraints. In this work, we introduce Federated Learning with Ordered Layer Freezing (FedOLF) to improve energy efficiency and reduce memory footprint while maintaining accuracy. Additionally, we employ the Tensor Operation Approximation technique to reduce the communication (and accordingly energy) cost, which can better preserve accuracy compared to traditional quantization methods. Experimental results demonstrate that FedOLF achieves higher accuracy and energy efficiency as well as lower memory footprint across EMNIST, CIFAR-10, CIFAR-100, and CINIC-10 benchmarks compared to existing methods.

1 INTRODUCTION

Federated Learning (FL) (McMahan et al., 2017) has gained significant traction in the Internet of 031 Things (IoT) for processing decentralized data and providing privacy-preserving intelligent services to clients (Jin et al., 2024; Zheng et al., 2023; Nguyen et al., 2022). However, the heterogeneous 033 nature of client devices poses a challenge due to varying system capacities. In real-world IoT envi-034 ronments, clients, often edge devices, exhibit diverse configurations in terms of processor, battery, bandwidth, and memory. Resource-constrained devices with limited hardware and bandwidth face difficulties in training and transmitting large neural networks, leading to straggling, low quality-of-037 service, and excessive computation and communication costs. Moreover, devices with insufficient 038 memory may be unable to handle memory-intensive neural networks, thus being excluded from FL with severe information loss. Therefore, addressing the issue of resource constraints is crucial for the successful application of FL in IoT systems (Imteaj et al., 2022; Pfeiffer et al., 2023a). 040

041 Several studies have been proposed to address resource constraints through techniques such as 042 dropout (Caldas et al., 2018; Horváth et al., 2021; Diao et al., 2021; Kim et al., 2023) or layer 043 freezing (Pfeiffer et al., 2023a;b). These methods involve training a subset of the global model with 044 reduced requirements on hardware, bandwidth, and memory on edge devices. Specifically, dropout involves pruning a fraction of the global model and sending the remaining sub-model to clients for training. However, it may significantly degrade accuracy in non-independent identical (non-iid) lo-046 cal data distributions. In such settings, data importance among clients may vary, and training an 047 underparameterized sub-model for an important client with data resembling the global distribution 048 may not sufficiently capture knowledge from local data, leading to decreased accuracy of the global 049 model (Pfeiffer et al., 2023b; Acar et al., 2021). 050

Instead of sub-models, layer freezing involves sending the full global model to all devices and al lowing resource-constrained devices to freeze some layers during training. For example, CoCoFL
 (Pfeiffer et al., 2023b) allows clients to randomly train certain layers while freezing the remaining, while SLT (Pfeiffer et al., 2023a) enables clients to sequentially train each layer in a bottom-up man-

054 ner, with other layers partially frozen. Compared to dropout, layer freezing is more resilient to noniid data by preserving the full model architecture on each client (Pfeiffer et al., 2023b). However, 056 layer freezing introduces heavy communication overhead since the global model must be transmitted to clients. Additionally, these methods overlook the fact that top-level layers, even though frozen, still need to store and pass gradient information back to lower-level active layers during backpropagation, resulting in heavy memory consumption. For example, Figure 1 illustrates a comparison between two training modules: (a) random layer freezing requiring more memory than (b) ordered 060 layer freezing due to a longer path for backpropagation of gradients (a longer red arrow in Figure 061 1(a)). To validate this analysis, we implement these two layer-freezing strategies using ResNet20 062 (He et al., 2016) with the CIFAR-100 dataset (Krizhevsky et al., 2009), and measure their maximum 063 memory usage using the TORCH. CUDA. MAX MEMORY ALLOCATED function (PyTorch, 2023) in 064 PyTorch. As depicted in Figure 1(c), random layer freezing consumes more memory compared to 065 ordered layer freezing, even when the same number of layers are frozen. 066



Figure 1: A comparison between (a) Random Layer Freezing and (b) Ordered Layer Freezing. The
 former requires more memory space to pass the gradient information back towards low-level active
 layers. (c) shows the maximum memory usage of these two modules in practice.

To address the shortcomings of existing methods, we introduce a new FL framework named Feder-087 ated Learning with Ordered Layer Freezing (FedOLF). In FedOLF, resource-constrained devices selectively freeze some low-level layers while training the remaining top-level layers. This approach substantially reduces the computation overhead and memory requirements of training, by shortening 090 the gradient backpropagation path as illustrated in Figure 1(b). Additionally, we empirically observe 091 that the gradient loss resulting from low-level frozen layers tends to diminish as training moves for-092 ward to top-level layers, which helps FedOLF maintain accuracy. Furthermore, we adopt an adapted Tensor Operation Approximation (TOA) scheme (Adelman et al., 2021) to reduce the communi-094 cation cost in FedOLF. Instead of the full global model, clients receive a low-rank approximation 095 of the frozen layers along with all active layers from the server during communication. Unlike con-096 ventional quantization methods, TOA minimally impacts training and significantly preserves model accuracy. The contributions of this paper are summarized as follows: 098

- We introduce FedOLF, an efficient FL framework addressing the memory shortage problem by allowing resource-constrained devices to train partial top-level layers of the global model. We also provide convergence analysis of FedOLF in non-convex settings.
- We propose an adapted TOA framework to reduce communication costs and memory footprint of FedOLF. Unlike the initial method that works on all layers, the adjusted TOA framework only works on frozen layers to ensure the active layers get fully trained.
- We evaluate FedOLF on EMNIST (with CNN), CIFAR-10 (with AlexNet), CIFAR-100 and CINIC-10 (with ResNet20 and ResNet44). Experimental results demonstrate that FedOLF outperforms the state-of-the-art by improving the accuracy by at least 0.3%, 6.4%, 12.8%, 4.4%, 6.6% and 1.29%, with higher energy efficiency and lower memory footprint.

099

102

108 2 LITERATURE REVIEW

110 Efficient Federated Learning: This stream of research aims to alleviate the computational and 111 communication costs associated with FL. Various approaches have been proposed to enhance com-112 putation efficiency, such as FedProx (Li et al., 2020), FedParl (Imteaj & Amini, 2021), and Pyra-113 midFL (Li et al., 2022), which reduce client training epochs to mitigate computation costs. To improve communication efficiency, methods like FedCOM (Haddadpour et al., 2021), FetchSGD 114 (Rothchild et al., 2020), and STC (Sattler et al., 2020) reduce the size of transmitted parameters 115 through message compression. Additionally, approaches like FedSL (Zhang et al., 2024), FedOBD 116 (Chen et al., 2022a), FedNew (Elgabli et al., 2022), Fedproto (Tan et al., 2022), and DS-FL (Itahara 117 et al., 2023) advocate for transmitting lightweight replacement messages, such as logits and pro-118 totypes, instead of the full global model. However, these methods often focus on singular aspects 119 of efficiency and fail to simultaneously address both computation and communication challenges. 120 Moreover, they do not adequately account for memory constraints on devices, as they typically in-121 volve full-model training on all clients. Adaptive dropout (Li et al., 2021a; Jiang et al., 2022; 2023; 122 Li et al., 2021b) offers a more comprehensive approach by enabling clients to train and transmit 123 lightweight sub-models, thereby achieving both computation and communication efficiency. Never-124 theless, adaptive dropout overlooks memory constraints, as clients must prune unimportant neurons 125 to generate sub-models, a process that requires pre-training the full model locally. FLrce (Niu et al., 2024) mitigates overall computation and communication costs by reducing FL iterations with an 126 early-stopping mechanism. However, it still entails full-model training on all devices irrespective of 127 memory constraints. 128

129 Federated Learning on Resource-Constrained Devices: The primary distinction between ef-130 ficient FL and resource-constrained FL lies in the latter's consideration of devices with limited 131 resources, such as memory space or bandwidth support, which are unable to train or transmit the entire model. To tackle this challenge, (Caldas et al., 2018; Horváth et al., 2021; Kim et al., 2023; 132 Diao et al., 2021) introduce the concept of sub-models, which contain fewer parameters and can be 133 trained and transmitted by resource-constrained clients. Specifically, Feddrop (Caldas et al., 2018) 134 employs random neuron pruning, FjORD (Horváth et al., 2021) and HeteroFL (Diao et al., 2021) 135 adopt a right-to-left approach for neuron pruning, and DepthFL (Kim et al., 2023) employs top-first 136 layer pruning. Unlike adaptive dropout, these works execute dropout at the server side, eliminating 137 the need for clients to pre-train a full model. However, these methods are susceptible to non-iid data 138 among clients, as training small sub-models on crucial clients may not capture sufficient knowledge 139 to construct an accurate global model. In contrast, CoCoFL (Pfeiffer et al., 2023b) and SLT (Pfeif-140 fer et al., 2023a) advocate for maintaining the full model architecture on all clients while freezing 141 certain layers on resource-constrained devices. CoCoFL randomly freezes layers within the local 142 model, whereas SLT partially freezes top-level layers and sequentially trains all layers from the bottom. The frozen layers remain untrained and untransmitted to enhance computation and com-143 munication efficiency. However, these approaches lead to increased memory usage, particularly in 144 the case of frozen top-level layers, which consume significant memory space to transmit gradient 145 information backward, as illustrated in Figure 1. 146

147 148

149

3 METHODOLOGY

3.1 PROBLEM SETUP 150

151 Given a network with one server and K devices (clients), and a global model w stored on the server 152 side, the goal of FL is to optimize the following problem: 153

$$\min_{w} f(w) := \mathbb{E}[f_{k}(w)] := \sum_{k=1}^{K} \frac{n_{k}}{n} (f_{k}(w)),$$

$$f_{k}(w) := \frac{1}{n_{k}} \sum_{n_{k}}^{i=1} \mathcal{L}(w, (\boldsymbol{x}_{i}, y_{i})).$$
(1)

156 157 158

154

159 f, the global objective function, is a weighted average of all local objective functions f_k ($1 \le k \le$ 160 K). For a client k, the local objective function f_k is equivalent to the empirical risk over its personal 161 dataset D_k , $n_k = |D_k|$ is the size of the local dataset and $\mathcal{L}(w, (x_i, y_i))$ is the prediction loss of 162 w over the i-th sample (x_i, y_i) in D_k . $n = \sum_{k=1}^{K} n_k$ is the total number of samples across all 163 local datasets. Moreover, let N denote the total number of layers in the global model w, and W_l 164 represent the l-th layer with parameter θ_l $(1 \le l \le N)$. The layer W_l can be viewed as a function 165 that takes the input feature representation x_{l-1} from the previous layer, and outputs a new feature 166 representation x_l , i.e. $x_l = W_l(x_{l-1}, \theta_l)$. Specially, $x_0 = x$ is the initial data sample, and $x_N = \hat{y}$ 167 is the model's final prediction.

169 Algorithm 1 FedOLF

170 **Require:** maximum global iteration T, clients $C = \{1, ..., K\}$ with numbers of frozen layers 171 $\{l_1, ..., l_K\}$, and initial global model w^0 , scale factor s. 1: for t = 1, 2, ..., T do 172 Server randomly samples a set of participating clients $C_t \subset C$. 173 2: 3: for every client $k \in C_t$ the server does: 174 Decompose w^t into $w_{F,k}^t$ and $w_{A,k}^t$ based on l_k . 4: 175 5: $\hat{w}_{F,k}^t \leftarrow TOA(w_{F,k}^t, s, l_k).$ ▷ Algorithm 2 176 Send $\hat{w}_{F,k}^t$ and $w_{A,k}^t$ to k. 6: 177 each $k \in C_t$ in parallel does: 7: 178
$$\begin{split} & w_k^t \leftarrow \hat{w}_{F,k}^t \circ w_{A,k}^t. \\ & \text{For local epochs } 1, ..., E: \\ & w_{A,k}^{t+1} = w_{A,k}^t - \eta \nabla f_k'(w_{A,k}^t). \end{split}$$
179 8: 9: ⊳ SGD 10: 181 Upload $w_{A,k}^{t+1}$ to the server. 11: for each layer $W_l \in w^t$, the server does: 183 12: 13: $C_{t,l} \leftarrow \{k : k \in C_t \land W_l \in w_{A,k}^{t+1}\}.$ \triangleright Obtain all clients that include W_l $\begin{array}{l} n_l \leftarrow \sum_{k \in C_{t,l}} n_k. \\ W_l \leftarrow \mathbb{E}(W_{k,l}) := \sum_{k \in C_{t,l}} \frac{n_k}{n_l} W_{k,l}. \end{array}$ 185 14: 186 15: Layer-wise aggregation 187 16: end for 188 17: return w^t

189 190 191

192

201

202

168

3.2 FEDOLF: FEDERATED LEARNING WITH ORDERED LAYER FREEZING

1

For a client k, the architecture of model w can be decomposed into two components $w_{F,k}$ and $w_{A,k}$ such that $w = w_{F,k} \circ w_{A,k}$. $w_{F,k} = \{W_1, ..., W_{l_k}\}$ and $w_{A,k} = \{W_{l_k+1}, ..., W_N\}$ are respectively the set of frozen and active layers. $l_k \in \{0, 1, ..., N-1\}$ is the number of frozen layers in training whose value depends on k's device capacity. For a powerful device that can train the entire model, we have $l_k = 0$ and $w_{F,k} = \emptyset$.

At global iteration t, client k downloads the global model w^t and decomposes w^t into $w^t_{F,k}$ and $w^t_{A,k}$ based on l_k . Afterwards, client k freezes $w^t_{F,k}$ and locally trains all parameters in $w^t_{A,k}$ by applying stochastic gradient descent (SGD) on dataset D_k according to Equation (2):

$$v_{A,k}^{t+1} = w_{A,k}^t - \eta \nabla f_k'(w_{A,k}^t)$$
⁽²⁾

 η is the learning rate and $\nabla f'_k$ is a low-error-rate approximation of the gradient ∇f_k in the case 203 of layer freezing. With layer freezing, the layers in $w_{F,k}^t$ will remain constant as training goes on, 204 and will subsequently generate a straggling feature representation $x'_{l_k} = x_{l_k} + \sigma_{l_k}$. x_{l_k} is the 205 true representation generated by $w_{F,k}^t$ if it is non-freezing, and σ_{l_k} is an error term representing the 206 divergence between x_{l_k} and x'_{l_k} . Feeding x_{l_k} and x'_{l_k} forward will respectively result in ∇f_k and 207 208 $\nabla f'_k$. Once local training is completed, client k only sends the updated layers $w^{t+1}_{A,k}$ to the server 209 for communication efficiency. After receiving the results from all participating clients, the server 210 updates the global model using a layer-wise aggregation strategy same as in (Pfeiffer et al., 2023b). 211 The details of FedOLF are outlined in Algorithm 1.

212

214

213 3.3 FEDOLF WITH TENSOR OPERATION APPROXIMATION

Furthermore, we propose an adapted Tensor Operation Approximation (TOA) framework (Adelman et al., 2021) dedicated to reducing the communication cost in FedOLF. Instead of the entire global

model w, a client k downloads $\hat{w}_{F,k}^t$ and $w_{A,k}^t$ from the server, where $\hat{w}_{F,k}^t$ is a low-rank approxima-tion of the frozen layers $w_{F,k}^{t}$ with fewer parameters. Unlike the initial TOA method which works on all layers, in this paper, the modified TOA works only on the frozen layers to ensure all active layers get fully trained. For illustration, let H_q denote the number of tensors in a frozen layer W_q , where a tensor is a filter or neuron if W_q is a convolution or fully-connected layer, respectively.



Figure 2: Within each frozen fully-connected layer W_q ($1 \le q < l$) containing H_q neurons, a subset W'_q (blue neurons) is derived by sampling $H'_q = \lfloor sH_q \rfloor$ neurons of the layer. Consequently, the approximation of $w^t_{F,k}$, represented as $\hat{w}^t_{F,k}$ is $\hat{w}^t_{F,k} = W'_1 \circ \dots \circ W'_{l-1} \circ W_l$.

For example, Figure 2 shows how TOA is applied on a fully-connected neural network with l frozen layers. For every layer W_q ($1 \le q < l$), except for the last frozen layer, the server samples $\lfloor sH_q \rfloor$ tensors from the layer and sends this subset of tensors to client k. $s (0 < s \le 1)$ is a scaling factor that determines the trade-off degree between accuracy and communication efficiency, with s = 1 representing that no TOA is applied. Moreover, TOA is not performed on the last frozen layer as shown in Figure 2, so that the dimensions of the output representation x'_1 and the following active layers remain unchanged. Based on the study of Adelman et al. (2021), we apply a weighted sampling strategy on TOA. With this strategy, TOA selects a tensor Z_j $(1 \le j \le H_q)$ within a frozen layer W_q with probabilities proportional to their *Frobenius* norms:

$$\mathbb{P}(\mathbf{Z}_j \in W'_q) = \frac{\|\mathbf{Z}_j\|_F}{\sum_{j=1}^{H_q} \|\mathbf{Z}_j\|_F}.$$
(3)

In this case, the approximation error $\mathbb{E}[\|x_l' - x_{l,TOA}'\|^2]$ will be minimized, where $x_{l,TOA}'$ and x_l' are respectively the output representations with and without TOA. The TOA technique significantly reduces the downstream communication cost in FedOLF by approximately $O(s^2)$. The procedure of TOA is shown in Algorithm 2.

Algorithm 2 TOA

Require: set of frozen layers w_F , scaling factor s, number of frozen layers l_k : 1: For every layer $W_q \in w_F$, $1 \le q \le l_k - 1$: 2: $H_q \leftarrow len(W_q)$. 3: $W'_q \leftarrow \text{sample(candidates=}\{Z_j\}_{j=1}^{H_q}, \text{weights=}\{\mathbb{P}(Z_j \in W'_q)\}_{j=1}^{H_q}, \text{number=}\lfloor sH_q \rfloor)$. 4: return $w'_F := W'_1 \circ \dots \circ W'_{l_k-1} \circ W_{l_k}$

3.4 DETERMINING THE NUMBER OF FROZEN LAYERS

Given a neural network w with N layers, the memory footprint m(w) can be computed as:

$$m(w) = \sum_{q=1}^{N} m_{\rm AM}(W_q) + m_{\rm G}(W_q) + m_{\rm W}(W_q) \approx \sum_{q=1}^{N} m_{\rm AM}(W_q)$$
(4)

That is, the overall memory footprint m is the accumulated memory footprint of three components, which are parameter weights (m_W) , gradients (m_G) and activation maps (m_{AM}) across all layers. Moreover, compared with weights and gradients, the size of activation maps is much more massive

and consumes a dominant memory space. Therefore, the overall memory footprint can be approximated as the total size of activation maps across all layers (Pfeiffer et al., 2023a).

In FedOLF, for a frozen layer W_q , $m_{AM}(W_q)$ becomes zero, as no activation maps have to be stored for training (Pfeiffer et al., 2023a). Accordingly, a client k can choose l_k to be the smallest value, given $\sum_{q=l_k+1}^{N} m_{AM}(W_q)$ (the size of activation maps in the remaining active layers) not exceeding its memory limit.

278 3.5 LOW-LEVEL LAYER SHARING AMONG CLIENTS

 According to the studies of (Zhang et al., 2024; Luo et al., 2021), low-level layers across various local models usually have higher degrees of Centered Kernal Alignment (CKA) similarity across different datasets (Kornblith et al., 2019), which means that these layers contain substantial redundant information and may generate similar feature representations. Motivated by this insight, in FedOLF, a resource-constrained device k can "borrow" the highly-generalized low-level layers from other clients by downloading $w_{F,k}^t$ from the server. Layers in $w_{F,k}^t$ have been trained by more powerful clients in previous rounds, and can be directly employed by k during the forward propagation phase of training without incurring significant errors.

3.6 VANISHING REPRESENTATION ERROR AND BOUNDED GRADIENT LOSS



Figure 3: During training, the *l* frozen layers will generate a feature representation x'_l that diverges from the true x_l . Affected by x'_l , the following active layers also generate inaccurate representations.

In addition to reducing memory usage, FedOLF preserves accuracy by mitigating representation errors induced by ordered layer freezing, with these errors diminishing as training advances through layers. For illustration, Figure 3 presents an exemplary model with l frozen layers and N - lactive layers. As described in subsection 3.2, owing to the staleness of frozen layers, all feature representations after layer W_l diverge from the true representations. However, the representation errors $\|\sigma_l\|, \|\sigma_{l+1}\|, ..., \|\sigma_N\|$ tends to decrease as the depth grows, where $\|\cdot\|$ represents l2-norm.

To verify our hypothesis, we first make the following assumption:

Assumption 1. The intrinsic function of each layer W_l is B_l -Lipschitz continuous with $B_l > 0$:

$$\forall \boldsymbol{x}_1, \boldsymbol{x}_2, \|W_l(\boldsymbol{x}_1, \boldsymbol{\theta}_l) - W_l(\boldsymbol{x}_2, \boldsymbol{\theta}_l)\| \le B_l \|\boldsymbol{x}_1 - \boldsymbol{x}_2\|.$$
(5)

Since $\mathbf{x}_{l+1} = W_l(\mathbf{x}_l, \boldsymbol{\theta}_l)$ and $\mathbf{x}'_{l+1} = W_l(\mathbf{x}'_l, \boldsymbol{\theta}_l)$, we can rewrite Equation (5) as:

$$\|\boldsymbol{x}_{l+1}' - \boldsymbol{x}_{l+1}\| \le B_l \|\boldsymbol{x}_l' - \boldsymbol{x}_l\|.$$
(6)

315 By induction, we have:

$$\|\boldsymbol{x}_{N}'-\boldsymbol{x}_{N}\| \leq \prod_{q=l}^{N-1} B_{q} \|\boldsymbol{x}_{l}'-\boldsymbol{x}_{l}\|, \text{ i.e. } \|\boldsymbol{\sigma}_{N}\| \leq \prod_{q=l}^{N-1} B_{q} \|\boldsymbol{\sigma}_{l}\|.$$
 (7)

In the experiment, we find that the term $\prod_{q=l}^{N-1} B_q$ is always shrinking (see Appendix A for evidence). Consequently, the representation error $\|\sigma_d\|$ $(l \le d \le N)$ caused by layer freezing tends to be vanishing as *d* increases, and the representation x'_d is gradually approaching the true representation x_d , thereby narrowing the gap between the computed gradient $\nabla \theta'_d$ and the true gradient $\nabla \theta_d$. As a result, the accumulated training error $\sum_{l}^{N} \|\nabla \theta'_l - \nabla \theta_l\|$ will be bounded (see Appendix B.1).

4 **CONVERGENCE ANALYSIS**

324

325

327

330 331 332

333

334 335

336

337 338

339 340 341

342

343

344 345

347 348

353

354 355 356

361

362

326 In this section, we analyze the convergence results for FedOLF on non-convex smooth objective functions. We do not require the objective function to be convex in the case of deep-learning neural 328 networks (Karimireddy et al., 2020). We make the following assumptions:

Assumption 2 (smoothness). The objective function f_k is L-smooth:

$$\forall w_1, w_2, \|\nabla f_k(w_1) - \nabla f_k(w_2)\| \le L \|w_1 - w_2\|.$$
(8)

Assumption 3 (Bounded variance). The variance of local gradients to the global gradient is bounded:

$$\forall k, w, \ \mathbb{E}(\|\nabla f_k(w) - \nabla f(w)\|^2) \le \gamma^2.$$
(9)

Furthermore, from Assumption 1 and Assumption 2, we can infer that the divergence of local gradient $\|\nabla f'_k - \nabla f_k\|$ resulting from layer freezing is bounded, which is defined in Corollary 1.

Corollary 1. For any client k, the divergence between the local gradient with and without layer freezing is bounded:

$$\forall k, w, \|\nabla f'_k(w) - \nabla f_k(w)\|^2 \le D^2.$$
 (10)

Based on Assumptions 1-3 and Corollary 1, we derive the following theorems:

Theorem 1. When the learning rate η satisfies $\frac{1}{L} < \eta < \frac{3}{2L}$, we have:

$$f(w^{t+1}) - f(w^{t}) \leq \frac{\eta}{2} (2\eta L - 3) (\mathbb{E}[\|\nabla f(w^{t})\|])^{2} + \eta D(\eta L - 1) \mathbb{E}[\|\nabla f(w^{t})\|] + \frac{\eta}{2} (2\eta L \gamma^{2} - \gamma^{2} + \eta L D^{2} + 2\eta L D \gamma).$$
(11)

Theorem 2. When the learning rate η satisfies $\eta \leq \frac{1}{L}$, we have:

$$f(w^{t+1}) - f(w^t) \le \frac{\eta}{2} \times (-\mathbb{E}[\|\nabla f(w^t)\|^2] + D^2 + \gamma^2 + 2D\gamma).$$
(12)

According to Theorem 1 and Theorem 2, when the learning rate is less than $\frac{3}{2L}$, the objective function f continues to decrease before w^t reaching a ϵ -critical point where $\|\nabla f(w^t)\| \leq \epsilon$. Specifically, when $\frac{1}{L} < \eta < \frac{3}{2L}$, we have $\epsilon = \epsilon_1 = \frac{D(\eta L - 1) + \sqrt{\eta D^2 L + 8\eta L \gamma^2 + 6\eta D L \gamma + D^2 - 3\gamma^2}}{3 - 2\eta L}$. When $\eta \leq \frac{1}{L}$, we have $\epsilon = \epsilon_2 = D + \gamma$. The proof can be found in Appendix B.

357 For FedOLF with TOA, the above theorems remain valid. The only difference is that the boundary 358 D in Corollary 1 is expected to become larger as TOA slightly increases the representation error. 359 Subsequently, the critical points ϵ_1 and ϵ_2 also increase, resulting in an earlier halt in the decay of f. 360

5 **EXPERIMENTS**

5.1 EXPERIMENT SETUP 364

365 Datasets and models. We evaluate the performance of FedOLF on the Extended MNIST (EMNIST) 366 (Cohen et al., 2017), CIFAR-10 (Krizhevsky et al., 2009), CIFAR-100 (Krizhevsky et al., 2009) and 367 CINIC-10 (Darlow et al., 2018) datasets. For EMNIST, we adopt a convolutional neural network 368 (CNN) consisting of two convolution layers and one fully-connected (FC) classifier (Horváth et al., 2021). For CIFAR-10, we employ AlexNet (Krizhevsky et al., 2012) (five convolution layers + two 369 FC layers). For CIFAR-100 and CINIC-10, we utilize ResNet20 and ResNet44 (He et al., 2016). 370

371 State-of-the-art for comparison. We compare FedOLF with the following representative methods 372 for resource-constrained FL: 1. Federated Dropout (Feddrop) (Caldas et al., 2018) randomly 373 prunes tensors in the global model and sends the remaining sub-model to clients for training. 2. 374 **FjORD** (Horváth et al., 2021) prunes the rightmost tensors of the global model. **3. HeteroFL** (Diao 375 et al., 2021) prunes the rightmost filters in convolution layers similar to FjORD, but keep the FC layers unchanged. 4. DepthFL (Kim et al., 2023) applies a top-first layer pruning method, and 376 adds extra classifiers to clients with fewer layers to distill knowledge. 5. CoCoFL (Pfeiffer et al., 377 2023b) let all clients store a full model locally and randomly freeze layers in training. 6. Successive 391

392

393

394 395 396

397

398

399 400

401

402

404

406

407

431

378	Dataset Model		EMNIST	CIFAR-10	CIFAR-100		CINIC-10	
379			CNN	AlexNet	ResNet20	ResNet44	ResNet20	ResNet44
381	Feddrop		32.11	14.33	17.02	6.2	9.87	10.31
382	FjORD		7.55	46.3	12.7	14.68	16.55	20.08
202	HeteroFL		17.4	54.79	12.32	12.96	10.69	10.03
303	DepthFL		60.25	16.74	24.87	37.82	9.97	34.28
384	CoCoFL		83.71	61.83	22.16	27.56	25.66	26.67
385	SLT		60.72	30.47	25.04	43.73	24.11	33.63
386		no TOA	84.02	68.27	37.85	48.15	32.27	35.57
387	FedOLF	TOA(0.75)	-	66.6	36.04	40.72	31.85	32.52
388		TOA(0.5)	-	63.12	24.93	29.68	31.92	30.89
389 390	FedAvg		84.42	69.22	46.01	49.46	36.32	37.59

Table 1: Comparison of the final test accuracy (in %) for T = 500 iterations in the non-iid case. Note that for EMNIST where the number of frozen layers is at most one, FedOLF+TOA is not evaluated as TOA only works with at least two frozen layers.

Layer Training (SLT) (Pfeiffer et al., 2023a) mandates all clients to sequentially train each layer from bottom to top, while freezing the parameters of the remaining layers. We also include the standard FedAvg benchmark (McMahan et al., 2017) for reference. All methods have run for three independent trials with their mean performance being recorded.

Parameter settings and system implementation. The experiment runs on a virtual network consisting of K = 100 clients operating on a desktop computer with one NVIDIA GeForce GTX 1650 GPU. The number of participants per round is $|C_t| = 10$ following the settings in (Horváth et al., 2021; Caldas et al., 2018). The maximum global iteration is set to T = 500 and the local training 403 epoch is E = 5 for all clients (Horváth et al., 2021; Li et al., 2021a). The learning rate is set to $\eta = 0.0001$ for EMNIST, $\eta = 0.001$ for CIFAR-10, and $\eta = 0.01$ for CIFAR-100 and CINIC-10 405 (Luo et al., 2021). The batch size is set to 16 for EMNIST and 128 for the remaining datasets (Li et al., 2021a; Horváth et al., 2021). The experiment is implemented with PyTorch 2.0.0 and Flower 1.4.0(Beutel et al., 2022). 408







Figure 6: The theoretical context-independent memory footprint (MB) among clients.

Data and system heterogeneity. We evaluate FedOLF in both iid and non-iid environments. For the 452 iid case, data are allocated to clients uniformly. For the non-iid case, we follow (Luo et al., 2021) and 453 allocate data to clients based on an extreme Dirichlet distribution with parameter 0.1. To emulate 454 system heterogeneity, we divide all clients into c uniform clusters that represent c different degrees of 455 device capability and resource constraints, as per (Horváth et al., 2021; Diao et al., 2021; Kim et al., 456 2023; Pfeiffer et al., 2023b). Specifically, following (Horváth et al., 2021), for CNN on EMNIST, c is 457 set to 2, wherein the numbers of pruned/frozen layers are 0 and 1 respectively for DepthFL, CoCoFL 458 and FedOLF; for Feddrop and FjORD, the sub-model ratios (i.e. the percentage of left neurons 459 per layer) are 0.5 and 1.0 for each cluster; for AlexNet on CIFAR-10 or ResNet20 on CIFAR-460 100/CINIC-10, c = 5 and the sub-model ratios are $\{0.2, 0.4, 0.6, 0.8, 1.0\}$ for Feddrop/FjORD; and the number of pruned/frozen layers or blocks are $\{4, 3, 2, 1, 0\}$ for DepthFL, CoCoFL, and FedOLF 461 (see Appendix D for details). For SLT that conducts universal successive training on all clients, the 462 scaling factor for the partial training procedure is set to 0.5 (Pfeiffer et al., 2023a). 463

464 465

449

450 451

5.2 EXPERIMENT RESULTS

466 Accuracy. Table 1 shows the accuracy comparison in the non-iid case (see Appendix C for the 467 iid case). As shown in Table 1, FedOLF achieves the highest final accuracy among all methods 468 on all datasets, which demonstrates the strength of FedOLF in preserving accuracy on resource-469 constrained devices. By looking through all methods, we find that dropout (Feddrop, FjORD) per-470 forms poorly with non-iid data, as training a sub-model cannot extract sufficient knowledge from 471 the local dataset to construct an accurate global model (Pfeiffer et al., 2023b). Besides, sub-models 472 with inconsistent architectures usually learn divergent parameter updates in training, and aggre-473 gating these updates altogether will inevitably compromise the global model's performance (Jiang 474 et al., 2022). Although existing layer freezing approaches (CoCoFL, SLT) improve accuracy by maintaining the full model architecture on all clients, it still lags behind FedOLF in accuracy. Be-475 cause in CoCoFL or SLT, the gradient loss caused by frozen layers does not decay as in FedOLF, 476 and impedes performance more straightforwardly. 477

Energy consumption and overall efficiency. Combining accuracy and energy consumption, we can derive the energy efficiency of each method. As shown in Figure 4, FedOLF significantly improves energy efficiency by obtaining the highest accuracy with the same amount of energy expenditure. The specific computation and communication costs, including FLOPs, data transmission volume and energy, can be found in Appendix C.

¹We merge the curves of Feddrop, FjORD, HeteroFL for brevity as their memory footprints are very close.

486 the real memory usage is usually context-dependent (physical device, programming language, etc), 487 we also calculate their theoretical memory usage following Equation (4). As shown in Figures 5 and 488 6, FedOLF effectively reduces the memory footprint both theoretically and practically. 489 490 Effect of s on accura Effect of s on the size of transmitted frozen parame Accuracy vs. communication efficient 491 492 493 494 495 496 497 498 (b) s vs. percentage of the frozen (c) TOA vs. QSGD. (a) s vs. accuracy. 499 parameters. 500 501 Figure 7: Effect of TOA and the scaling of factor s. 502 504 AlexNet memory footprint ResNet 44 memory footprint ResNet 20 memory footprint 505 s=0.5 s=0.5 s=1506 0.75 s=0.5 s=0.7 200 150 s=0.25 s=0.75 s=1.0 507 s=0.25 300 s=1.0 150 s=0.25 1000 200 100 509 100 510 50 Cluster 2 Cluster 3 Cluster 4 Cluster 5 Cluster 1 Cluster 2 Cluster 3 Cluster 4 Cluster 5 Cluster 1 Cluster 2 Cluster 3 Cluster 4 Cluster 5 511 Cluster 1 512 (c) ResNet44. (a) AlexNet. (b) ResNet20. 513 Figure 8: Effect of the scaling of factor s on the practical memory footprint (MB). 514 515 516

Hyperparameter tuning and ablation study. We tune the scaling factor of TOA s using a grid 517 search within $\{0.25, 0.5, 0.75, 1.0\}$, where s = 1 is equivalent to FedOLF without TOA. Results in 518 Table 1 and Figures 7a and 7b reveal that TOA effectively reduces the downstream communication 519 cost without degrading much accuracy (except for CIFAR-100 with ResNet-44). For example, a 520 scaling factor s = 0.25 can reduce the size of the transmitted frozen parameters by utmost 84% with a minor 5.56% accuracy loss compared with FedOLF sole (AlexNet). Besides, TOA further 521 reduces the practical memory footprint as Figure 8 shows. Additionally, we compare TOA with the 522 well-recognized quantized SGD (QSGD) method (Alistarh et al., 2017) for AlexNet on CIFAR-10. 523 As shown in Figure 7c, TOA achieves much higher accuracy than QSGD given the same degree of 524 communication efficiency. Specifically, TOA (s = 0.5) is compared with QSGD with 8 bits and 525 TOA (s = 0.75) is compared with QSGD with 16 bits so that their reductions of communication 526 cost are approximately equal. 527

528 529

530

6 CONCLUSION

531 This paper proposed Federated Learning with Ordered Layer Freezing (FedOLF), an efficient FL 532 framework where edge devices only train the top-level layers of the model to accommodate resource 533 constraints. The OLF strategy can minimize the backpropagation path length and the gradient er-534 ror, which significantly reduces the memory requirement and improves accuracy. We also enhance 535 FedOLF with the Tensor Operation Approximation (TOA) technique (Adelman et al., 2021), further 536 alleviating energy consumption and memory footprint with less accuracy sacrifice. In the future, 537 we plan to explore the similarities of local clients by using techniques like learning vector quantization (Qin & Suganthan, 2005), to make similar clients share the same layer freezing and TOA 538 settings. We also plan to enhance the engagement of FedOLF in IoT applications such as mobile edge networks (Jin et al., 2024) and video surveillance (Zhang et al., 2022a).

540 REFERENCES

546

547

548

560

561

562 563

564

565

566

570

571

572

575

- Durmus Alp Emre Acar, Yue Zhao, Ramon Matas, Matthew Mattina, Paul Whatmough, and
 Venkatesh Saligrama. Federated learning based on dynamic regularization. In *International Con- ference on Learning Representations*, 2021.
 - Menachem Adelman, Kfir Levy, Ido Hakimi, and Mark Silberstein. Faster neural network training with approximate tensor operations. *Advances in Neural Information Processing Systems*, 34: 27877–27889, 2021.
- 549 Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd:
 550 Communication-efficient sgd via gradient quantization and encoding. In I. Guyon, U. Von
 551 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances
 552 in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017.
- Daniel J. Beutel, Taner Topal, Akhil Mathur, Xinchi Qiu, Javier Fernandez-Marques, Yan Gao, Lorenzo Sani, Kwing Hei Li, Titouan Parcollet, Pedro Porto Buarque de Gusmão, and Nicholas D. Lane. Flower: A friendly federated learning research framework, 2022.
- Sebastian Caldas, Jakub Konečny, H. Brendan McMahan, and Ameet Talwalkar. Expanding the
 reach of federated learning by reducing client resource requirements. In *NeurIPS Workshop on Federated Learning for Data Privacy and Confidentiality*, 2018.
 - Xiaoyu Cao, Minghong Fang, Jia Liu, and Neil Zhenqiang Gong. Fltrust: Byzantine-robust federated learning via trust bootstrapping. *arXiv preprint arXiv:2012.13995*, 2020.
 - Yuanyuan Chen, Zichen Chen, Pengcheng Wu, and Han Yu. Fedobd: Opportunistic block dropout for efficiently training large-scale neural networks through federated learning. *arXiv preprint arXiv:2208.05174*, 2022a.
- Zhengsu Chen, Lingxi Xie, Jianwei Niu, Xuefeng Liu, Longhui Wei, and Qi Tian. Network adjust ment: Channel and block search guided by resource utilization ratio. *International Journal of Computer Vision*, 130:1–16, 03 2022b. doi: 10.1007/s11263-021-01566-5.
 - Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. Emnist: an extension of mnist to handwritten letters. *arXiv preprint arXiv:1702.05373*, 2017.
- Luke N. Darlow, Elliot J. Crowley, Antreas Antoniou, and Amos J. Storkey. Cinic-10 is not imagenet
 or cifar-10, 2018.
- Enmao Diao, Jie Ding, and Vahid Tarokh. Heterofl: Computation and communication efficient
 federated learning for heterogeneous clients. In *International Conference on Learning Representations*, 2021.
- Anis Elgabli, Chaouki Ben Issaid, Amrit Singh Bedi, Ketan Rajawat, Mehdi Bennis, and Vaneet Aggarwal. Fednew: A communication-efficient and privacy-preserving newton-type method for federated learning. In *International conference on machine learning*, pp. 5861–5877. PMLR, 2022.
- Farzin Haddadpour, Mohammad Mahdi Kamani, Aryan Mokhtari, and Mehrdad Mahdavi. Feder ated learning with compression: Unified analysis and sharp guarantees. In *International Confer- ence on Artificial Intelligence and Statistics*, pp. 2350–2358. PMLR, 2021.
- Meng Hao, Hongwei Li, Guowen Xu, Hanxiao Chen, and Tianwei Zhang. Efficient, private and robust federated learning. In *Proceedings of the 37th Annual Computer Security Applications Conference*, ACSAC '21, pp. 45–60, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450385794. doi: 10.1145/3485832.3488014.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

594 Samuel Horváth, Stefanos Laskaridis, Mario Almeida, Ilias Leontiadis, Stylianos Venieris, and 595 Nicholas Lane. Fjord: Fair and accurate federated learning under heterogeneous targets with or-596 dered dropout. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan 597 (eds.), Advances in Neural Information Processing Systems, volume 34, pp. 12876–12889. Curran 598 Associates, Inc., 2021. Ahmed Imteaj and M Hadi Amini. Fedparl: Client activity and resource-oriented lightweight feder-600 ated learning model for resource-constrained heterogeneous iot environment. Frontiers in Com-601 munications and Networks, 2:657653, 2021. 602 Ahmed Imteaj, Urmish Thakker, Shiqiang Wang, Jian Li, and M. Hadi Amini. A survey on federated 603 learning for resource-constrained iot devices. IEEE Internet of Things Journal, 9(1):1-24, 2022. 604 605 Sohei Itahara, Takayuki Nishio, Yusuke Koda, Masahiro Morikura, and Koji Yamamoto. 606 Distillation-based semi-supervised federated learning for communication-efficient collaborative 607 training with non-iid private data. IEEE Transactions on Mobile Computing, 22(1):191–205, 608 2023. 609 Yuang Jiang, Shiqiang Wang, Victor Valls, Bong Jun Ko, Wei-Han Lee, Kin K Leung, and Leandros 610 Tassiulas. Model pruning enables efficient federated learning on edge devices. *IEEE Transactions* 611 on Neural Networks and Learning Systems, 2022. 612 613 Zhida Jiang, Yang Xu, Hongli Xu, Zhiyuan Wang, Jianchun Liu, Qian Chen, and Chunming Qiao. Computation and communication efficient federated learning with adaptive model pruning. IEEE 614 Transactions on Mobile Computing, pp. 1–18, 2023. 615 616 Huiying Jin, Pengcheng Zhang, Hai Dong, Xinmiao Wei, Yuelong Zhu, and Tao Gu. Mobility-aware 617 and privacy-protecting gos optimization in mobile edge networks. IEEE Transactions on Mobile 618 Computing, 23(2):1169-1185, 2024. doi: 10.1109/TMC.2022.3230856. 619 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and 620 Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for federated learning. 621 In Proceedings of the 37th International Conference on Machine Learning, volume 119 of Pro-622 ceedings of Machine Learning Research, pp. 5132–5143. PMLR, 2020. 623 Minjae Kim, Sangyoon Yu, Suhyun Kim, and Soo-Mook Moon. DepthFL : Depthwise federated 624 learning for heterogeneous clients. In The Eleventh International Conference on Learning Repre-625 sentations, 2023. 626 627 Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural 628 network representations revisited. In International conference on machine learning, pp. 3519-629 3529. PMLR, 2019. 630 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep con-631 volutional neural networks. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger (eds.), 632 Advances in Neural Information Processing Systems, volume 25. Curran Associates, Inc., 2012. 633 634 Alex Krizhevsky et al. Learning multiple layers of features from tiny images. University of Toronto, 2009. 635 636 Ang Li, Jingwei Sun, Pengcheng Li, Yu Pu, Hai Li, and Yiran Chen. Hermes: An efficient fed-637 erated learning framework for heterogeneous mobile clients. In Proceedings of ACM Mobi-638 Com, pp. 420–437, New York, NY, USA, 2021a. Association for Computing Machinery. ISBN 639 9781450383424. 640 Ang Li, Jingwei Sun, Xiao Zeng, Mi Zhang, Hai Li, and Yiran Chen. Fedmask: Joint compu-641 tation and communication-efficient personalized federated learning via heterogeneous masking. 642 In Proceedings of the 19th ACM Conference on Embedded Networked Sensor Systems, SenSys 643 '21, pp. 42-55, New York, NY, USA, 2021b. Association for Computing Machinery. ISBN 644 9781450390972. doi: 10.1145/3485730.3485929. 645 Chenning Li, Xiao Zeng, Mi Zhang, and Zhichao Cao. Pyramidfl: A fine-grained client selection 646 framework for efficient federated learning. In Proceedings of ACM MobiCom, pp. 158–171, New 647 York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450391818.

651

680

- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith.
 Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, pp. 429–450, 2020.
- Mi Luo, Fei Chen, Dapeng Hu, Yifan Zhang, Jian Liang, and Jiashi Feng. No fear of heterogeneity: Classifier calibration for federated learning with non-IID data. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
 Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pp. 1273–1282, 2017.
- Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training of machine learning models. In *International Conference on Machine Learning*, pp. 6950–6960.
 PMLR, 2020.
- Dinh C. Nguyen, Quoc-Viet Pham, Pubudu N. Pathirana, Ming Ding, Aruna Seneviratne, Zihuai
 Lin, Octavia Dobre, and Won-Joo Hwang. Federated learning for smart healthcare: A survey.
 ACM Computing Survey, 55(3), feb 2022. ISSN 0360-0300.
- Ziru Niu, Hai Dong, A. K. Qin, and Tao Gu. Flrce: Resource-efficient federated learning with
 early-stopping strategy. *IEEE Transactions on Mobile Computing*, 2024.
- Kilian Pfeiffer, Ramin Khalili, and Joerg Henkel. Aggregating capacity in FL through successive
 layer training for computationally-constrained devices. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023a.
- Kilian Pfeiffer, Martin Rapp, Ramin Khalili, and Joerg Henkel. CocoFL: Communication- and computation-aware federated learning via partial NN freezing and quantization. *Transactions on Machine Learning Research*, 2023b. ISSN 2835-8856.
- 676 PyTorch, 2023. URL https://pytorch.org/docs/stable/generated/torch. cuda.max_memory_allocated.html#torch.cuda.max_memory_allocated.
- A. K. Qin and P. N. Suganthan. Initialization insensitive LVQ algorithm based on cost-function adaptation. *Pattern Recognition*, 38(5):773–776, 2005.
- Daniel Rothchild, Ashwinee Panda, Enayat Ullah, Nikita Ivkin, Ion Stoica, Vladimir Braverman,
 Joseph Gonzalez, and Raman Arora. Fetchsgd: Communication-efficient federated learning with
 sketching. In *International Conference on Machine Learning*, pp. 8253–8265. PMLR, 2020.
- Felix Sattler, Simon Wiedemann, Klaus-Robert Müller, and Wojciech Samek. Robust and communication-efficient federated learning from non-i.i.d. data. *IEEE Transactions on Neural Networks and Learning Systems*, 31(9):3400–3413, 2020.
- Yue Tan, Guodong Long, Lu Liu, Tianyi Zhou, Qinghua Lu, Jing Jiang, and Chengqi Zhang. Fed proto: Federated prototype learning across heterogeneous clients. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 8432–8440, 2022.
- Sijia Zhang, Maoguo Gong, Yu Xie, A. K. Qin, Hao Li, Yuan Gao, and Yew-Soon Ong. Influence aware attention networks for anomaly detection in surveillance videos. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(8):5427–5437, 2022a.
- Weishan Zhang, Tao Zhou, Qinghua Lu, Yong Yuan, Amr Tolba, and Wael Said. Fedsl: A communication efficient federated learning with split layer aggregation. *IEEE Internet of Things Journal*, pp. 1–1, 2024. doi: 10.1109/JIOT.2024.3350241.
- Zaixi Zhang, Xiaoyu Cao, Jinyuan Jia, and Neil Zhenqiang Gong. Fldetector: Defending federated learning against model poisoning attacks via detecting malicious clients. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '22,
 pp. 2545–2555, New York, NY, USA, 2022b. Association for Computing Machinery. ISBN 9781450393850.

Tianyue Zheng, Ang Li, Zhe Chen, Hongbo Wang, and Jun Luo. Autofed: Heterogeneity-aware federated multimodal learning for robust autonomous driving. In *Proceedings of the 29th Annual International Conference on Mobile Computing and Networking*, pp. 1–15, 2023.

A VANISHING REPRESENTATION ERROR AND BOUNDED TRAINING LOSS

We find that in FedOLF, the negative impact of low-level frozen layers tends to vanish. In formulation, let $\sigma_{l,l+1} \in \mathbb{R}^+$ denote the ratio between the representation error between two consecutive frozen layers W_l and W_{l+1} , that is:

$$\sigma_{l,l+1} := \frac{\|\boldsymbol{\sigma}_{l+1}\|}{\|\boldsymbol{\sigma}_{l}\|} = \frac{\|\boldsymbol{x}_{l+1}' - \boldsymbol{x}_{l+1}\|}{\|\boldsymbol{x}_{l}' - \boldsymbol{x}_{l}\|} = \frac{\|W_{l+1}(\boldsymbol{x}_{l}', \boldsymbol{\theta}_{l+1}) - W_{l+1}(\boldsymbol{x}_{l}', \boldsymbol{\theta}_{l+1})\|}{\|\boldsymbol{x}_{l}' - \boldsymbol{x}_{l}\|}.$$
 (13)

In addition, we rewrite Assumption 1 here for better illustration:

Assumption 1. The intrinsic function of each layer W_l is B_l -Lipschitz continuous with $B_l > 0$:

$$\forall \boldsymbol{x}_1, \boldsymbol{x}_2, \|W_l(\boldsymbol{x}_1, \boldsymbol{\theta}_l) - W_l(\boldsymbol{x}_2, \boldsymbol{\theta}_l)\| \le B_l \|\boldsymbol{x}_1 - \boldsymbol{x}_2\|.$$
(14)

By induction:

$$\|\boldsymbol{\sigma}_N\| \le \prod_{q=l}^{N-1} B_q \|\boldsymbol{\sigma}_l\|.$$
(15)

Empirically, we find that the accumulative product $\|\sigma_N\| \leq \prod_{q=l}^{N-1} B_q \|\sigma_l\|$ across layers is usually *shrinking*. For example, for AlexNet on CIFAR-10, we compute the error ratios among all convolution layers for a randomly selected client who only freezes the first layer, as shown in Figure 9:





Figure 9: The ratios of the representation error between two consecutive layers in AlexNet with only the first layer frozen. Layers 1-5 are convolution layers, layers 6 and 7 are fully-connected layers and layer 8 is the classifier.

In this scenario, each boundary B_l is highly likely to be smaller than one, which indicates that the representation error $\|\sigma_{l+d}\| = \prod_{q=l}^{l+d-1} B_q \|\sigma_l\|$ tends to vanish as *d* increases. Consequently, the top level learns relatively accurate parameter updates by forwarding representations with lower error rates. Empowered by this property, FedOLF is able to achieve higher accuracy compared to other layer freezing methods.

For models like ResNet, the error ratios across layers are not consistently less than 1 as Figure 10 illustrates. We attribute this phenomenon to the unique architecture of ResNet, i.e. it adds connections between residual blocks so that the representation errors vanish less slowly. However, the term $\prod_{q=1}^{l-1} B_q$ still exhibits an overall vanishing trend as l increases, because the remaining bounds B_q for q > 2 are likely to be less than one. To further support the validity of this assumption, please refer to (Mirzasoleiman et al., 2020), where an equivalent assumption has also been made.



Figure 10: The ratios of the representation error between two consecutive residual blocks in ResNet20 with only the first layer frozen.

B THEORETICAL PROOF

This section presents the detailed proof of Corollary 1 and Theorems 1 and 2 in Section 4. First, we rewrite Assumptions 2 and 3 here:

779 Assumption 2 (smoothness). The local objective function f_k is L-smooth:

$$\forall w_1, w_2, \ \|\nabla f_k(w_1) - \nabla f_k(w_2)\| \le L \|w_1 - w_2\|.$$
(16)

Assumption 3 (Bounded variance). The variance of local gradients to the global gradient is bounded:

$$\forall k, w, \mathbb{E}(\|\nabla f_k(w) - \nabla f(w)\|^2) \le \gamma^2.$$
(17)

786 B.1 PROOF OF COROLLARY 1

Based on Assumptions 1 and 2, we can derive that the gradient divergence caused by layer freezing is bounded, as defined in Corollary 1:

Corollary 1. For any client k, the divergence between the local gradient with and without layer freezing is bounded:

$$\forall k, w, \|\nabla f'_k(w) - \nabla f_k(w)\|^2 \le D^2.$$
 (18)

793 Proof.

 We can represent a model w in the format of the set of all layers' parameters, i.e.:

$$w := (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, ..., \boldsymbol{\theta}_N). \tag{19}$$

798 Accordingly, for the gradient $\nabla f_k(w)$ we have:

$$\nabla f_k(w) = \nabla \theta_1 + \nabla \theta_2 + \dots + \nabla \theta_N.$$
⁽²⁰⁾

where $\nabla \theta_l = \nabla f_k(\boldsymbol{x}_{l-1}, \theta_l)$ for any $l \ (1 \le l \le N)$.

Furthermore, we define $\nabla \theta'_l := \nabla f_k(x'_{l-1}, \theta_l)$. Then we can use $\|\nabla \theta'_l - \nabla \theta_l\|$ to represent the gradient error on the *l*-th layer caused by layer freezing.

804 Since f_k is *L*-smooth, we have:

$$\|\nabla \theta_{l}' - \nabla \theta_{l}\| = \|\nabla f_{k}(\boldsymbol{x}_{l-1}', \theta_{l}) - \nabla f_{k}(\boldsymbol{x}_{l-1}, \theta_{l})\| \le L \|\boldsymbol{x}_{l-1}' - \boldsymbol{x}_{l-1}\| = L \|\boldsymbol{\sigma}_{l-1}\|.$$
(21)

By induction, we have:

$$\|\nabla \boldsymbol{\theta}_{l+d}' - \nabla \boldsymbol{\theta}_{l+d}\| \le L \|\boldsymbol{\sigma}_{l+d-1}\| \le \prod_{q=l}^{l+d-1} B_q L \|\boldsymbol{\sigma}_l\|.$$
(22)

Based on Equations (20) and (22), for the gradient difference $\|\nabla f'_k - \nabla f_k\|$, we have:

$$\|\nabla f'_{k}(w) - \nabla f_{k}(w)\|$$

$$= \|\sum_{l=1}^{N} \nabla \theta'_{l} - \sum_{l=1}^{N} \nabla \theta_{l}\|$$

$$= \|\sum_{l=1}^{N} (\nabla \theta'_{l} - \nabla \theta_{l})\|$$

$$\leq \sum_{l=1}^{N} \|\nabla \theta'_{l} - \nabla \theta_{l}\|$$

$$\leq \sum_{l=1}^{N} \|\nabla \theta'_{l} - \nabla \theta_{l}\|$$

$$\leq \sum_{l=1}^{N} \prod_{q=l}^{l-1} B_{q}L\|\sigma_{1}\|.$$

$$\leq \sum_{l=1}^{N} \prod_{q=l}^{l-1} B_{q}L\|\sigma_{1}\|.$$

Note that the first term $\|\nabla \theta'_1 - \nabla \theta_1\|$ equals zero and gets eliminated in Equation (23), because $\nabla \theta'_1 = \nabla \theta_1 = \mathbf{0}$ when the number of frozen layers is at least one.

Given that $\prod_{q=l}^{l-1} B_q$ is gradually vanishing as shown in Appendix A, the summation $\sum_{l=2}^{N} B^{l-1}L \| \sigma_1 \|$ must be finite and can be upper-bounded, alogn with $\| \nabla f'_k(w) - \nabla f_k(w) \|$. Therefore, Corollary 1 is naturally proven by setting D as the boundary.

B.2 PROOF OF THEOREM 1 AND THEOREM 2

Based on Assumptions 1-3 and Corollary 1, we derive Theorem 1 and Theorem 2:

Theorem 1. When the learning rate η satisfies $\frac{1}{L} < \eta < \frac{3}{2L}$, we have:

$$f(w^{t+1}) - f(w^{t}) \leq \frac{\eta}{2} (2\eta L - 3) (\mathbb{E}[\|\nabla f(w^{t})\|])^{2} + \eta D(\eta L - 1) \mathbb{E}[\|\nabla f(w^{t})\|] + \frac{\eta}{2} (2\eta L \gamma^{2} - \gamma^{2} + \eta L D^{2} + 2\eta L D\gamma)$$
(24)

Theorem 2. When the learning rate η satisfies $\eta \leq \frac{1}{L}$, we have:

$$f(w^{t+1}) - f(w^t) \le \frac{\eta}{2} \times (-\mathbb{E}[\|\nabla f(w^t)\|^2] + D^2 + \gamma^2 + 2D\gamma).$$
(25)

Proof.

 Since every f_k is L-smooth based on Assumption 2, f is also L-smooth, so that we have:

$$f(w^{t+1}) - f(w^t) \le \langle w^{t+1} - w^t, \nabla f(w^t) \rangle + \frac{L}{2} \|w^{t+1} - w^t\|^2.$$
(26)

In the setting of layer freezing, we have $w^{t+1} = w^t - \eta \nabla f'(w^t)$ and $\nabla f'(w^t) = \mathbb{E}[\nabla f'_k(w^t)]$. Therefore:

$$f(w^{t+1}) - f(w^{t})$$

$$\leq -\eta \langle \mathbb{E}[\nabla f'_{k}(w^{t})], \nabla f(w^{t}) \rangle + \frac{L}{2} \| -\eta \mathbb{E}[\nabla f'_{k}(w^{t})] \|^{2}$$

$$= -\eta \mathbb{E} \left[\langle \nabla f'_{k}(w^{t}), \nabla f(w^{t}) \rangle \right] + \frac{L\eta^{2}}{2} \| \mathbb{E}[\nabla f'_{k}(w^{t})] \|^{2}$$

$$\leq -\eta \mathbb{E} \left[\langle \nabla f'_{k}(w^{t}), \nabla f(w^{t}) \rangle \right] + \frac{L\eta^{2}}{2} \mathbb{E}(\| \nabla f'_{k}(w^{t}) \|^{2}).$$

$$(27)$$

864

Since $\|\nabla f'_{k}(w^{t}) - \nabla f(w^{t})\|^{2} = \|\nabla f'_{k}(w^{t})\|^{2} - 2 \langle \nabla f'_{k}(w^{t}), \nabla f(w^{t}) \rangle + \|\nabla f(w^{t})\|^{2}$, Equation 865 (27) can be written as: 866 $f(w^{t+1}) - f(w^t)$ 867 $\leq \frac{\eta}{2} \mathbb{E}(\|\nabla f_k'(w^t) - \nabla f(w^t)\|^2 - \|\nabla f_k'(w^t)\|^2 - \|\nabla f(w^t)\|^2) + \frac{L\eta^2}{2} \mathbb{E}(\|\nabla f_k'(w^t)\|^2)$ 868 $= \frac{\eta}{2} \mathbb{E}(\|\nabla f_k'(w^t) - \nabla f(w^t)\|^2) + \frac{\eta}{2}(\eta L - 1)\mathbb{E}[\|\nabla f_k'(w^t)\|^2] - \frac{\eta}{2}\mathbb{E}[\|\nabla f(w^t)\|^2]$ 870 (28)871 872 $= \frac{\eta}{2} \mathbb{E} \| (\nabla f_k'(w^t) - \nabla f_k(w^t) + \nabla f_k(w^t) - \nabla f(w^t) \|^2)$ 873 $+ \frac{\eta}{2} (\eta L - 1) \mathbb{E}[\|\nabla f'_k(w^t)\|^2] - \frac{\eta}{2} \mathbb{E}[\|\nabla f(w^t)\|^2].$ 874 875 According to Cauchy-Schwarz inequality, $\|\nabla f'_k(w^t) - \nabla f_k(w^t) + \nabla f_k(w^t) - \nabla f(w^t)\|^2 \leq (\|\nabla f'_k(w^t) - \nabla f_k(w^t)\| + \|\nabla f_k(w^t) - \nabla f(w^t)\|)^2$. Therefore, from Equation (28) we get: 876 877 878 $f(w^{t+1}) - f(w^t)$ 879 $\leq \frac{\eta}{2} \mathbb{E}[(\|\nabla f_k'(w^t) - \nabla f_k(w^t)\| + \|\nabla f_k(w^t) - \nabla f(w^t)\|)^2]$ 880 $+ \frac{\eta}{2} (\eta L - 1) \mathbb{E}[\|\nabla f'_k(w^t)\|^2] - \frac{\eta}{2} \mathbb{E}[\|\nabla f(w^t)\|^2]$ 883 $= \frac{\eta}{2} \mathbb{E}(\|\nabla f_k'(w^t) - \nabla f_k(w^t)\|^2 + \|\nabla f_k(w^t) - \nabla f(w^t)\|^2)$ (29)884 $+ 2\mathbb{E}(\|\nabla f_k'(w^t) - \nabla f_k(w^t)\| \times \|\nabla f_k(w^t) - \nabla f(w^t)\|)$ 885 $+ \frac{\eta}{2} (\eta L - 1) \mathbb{E}[\|\nabla f'_k(w^t)\|^2] - \frac{\eta}{2} \mathbb{E}[\|\nabla f(w^t)\|^2]$ 887 888 $\leq \frac{\eta}{2}(D^2 + \gamma^2 + 2D\gamma) + \frac{\eta}{2}(\eta L - 1)\mathbb{E}[\|\nabla f'_k(w^t)\|^2] - \frac{\eta}{2}\mathbb{E}[\|\nabla f(w^t)\|^2].$ 889 890 The last inequality in Equation (29) results from Corollary 1 and Assumption 3. 891 When the learning rate $\eta > \frac{1}{T}$, we have $\eta L - 1 > 0$. In this case, we can upper bound $\frac{\eta}{2}(\eta L - 1)$ 892 1) $\mathbb{E}[\|\nabla f'_k(w^t)\|^2]$ by upper bounding $\mathbb{E}[\|\nabla f'_k(w^t)\|^2]$. 893 First, we bound $\|\nabla f'_k(w^t)\|$. Based on Corollary 1 and triangle inequality, we have: 894 895 $\|\nabla f'_k(w^t)\| - \|\nabla f_k(w^t)\| \le \|\nabla f'_k(w^t) - \nabla f_k(w^t)\| \le D.$ (30)896 That is: 897 898 $\|\nabla f_k'(w^t)\|^2 \le (\|\nabla f_k(w^t)\| + D)^2 = \|\nabla f_k(w^t)\|^2 + D^2 + 2D\|\nabla f_k(w^t)\|.$ (31)899 By taking the expectation on Equation (31), we get: 900 $\mathbb{E}[\|\nabla f_k'(w^t)\|^2] \le \mathbb{E}[\|\nabla f_k(w^t)\|^2] + D^2 + 2D \mathbb{E}[\|\nabla f_k(w^t)\|].$ (32)901 902 Because of the triangle inequality, we have: 903 $\mathbb{E}[\|\nabla f_k(w^t)\|] = \mathbb{E}[\|\nabla f_k(w^t) - \nabla f(w^t) + \nabla f(w^t)\|]$ 904 $\leq \mathbb{E}[\|\nabla f_k(w^t) - \nabla f(w^t)\|] + \mathbb{E}[\|\nabla f(w^t)\|]$ 905 (33)906 $\leq \mathbb{E}[\|\nabla f(w^t)\|] + \gamma.$ 907 The last inequality in Equation (33) holds because $\mathbb{E}[\|\nabla f_k(w^t) - \nabla f(w^t)\|] \le \gamma$ as $(\mathbb{E}[\|\nabla f_k(w^t) - \nabla f(w^t)\|] \le \gamma$ 908 $\nabla f(w^t) \|])^2 \leq \mathbb{E}[\|\nabla f_k(w^t) - \nabla f(w^t)\|^2] \leq \gamma^2$ by Assumption 3. Moreover, by expanding As-909 sumption 3, we have: 910 911 $\mathbb{E}[\|\nabla f_k(w^t)\|^2]$ 912 $= \mathbb{E}[\|\nabla f_k(w^t) - \nabla f(w^t) + \nabla f(w^t)\|^2]$ 913 $= \mathbb{E}[\|\nabla f(w^t)\|^2] + \mathbb{E}[\|\nabla f_k(w^t) - \nabla f(w^t)\|^2] + 2\mathbb{E}(\langle \nabla f_k(w^t) - \nabla f(w^t), \nabla f(w^t) \rangle)$ 914 (34)915 $<\mathbb{E}[\|\nabla f(w^{t})\|^{2}] + \gamma^{2} + 2\mathbb{E}(\langle \nabla f_{k}(w^{t}) - \nabla f(w^{t}), \nabla f(w^{t}) \rangle)$ 916 $<\mathbb{E}[\|\nabla f(w^{t})\|^{2}] + \gamma^{2} + \mathbb{E}(\|\nabla f_{k}(w^{t}) - \nabla f(w^{t})\|^{2} + \|\nabla f(w^{t})\|^{2})$ 917 $= 2\mathbb{E}[\|\nabla f(w^t)\|^2] + 2\gamma^2.$

By combining Equations (32), (33), (34) altogether, we get:

Accordingly, we can rewrite Equation (29) as:

f(a,t+1) f(a,t)

$$\begin{aligned} &\int (w^{t-1}) - f(w^{t}) \\ &\leq -\frac{\eta}{2} \mathbb{E}[\|\nabla f(w^{t})\|^{2}] + \frac{\eta}{2} (D^{2} + \gamma^{2} + 2D\gamma) + \frac{\eta}{2} (\eta L - 1) \mathbb{E}[\|\nabla f_{k}'(w^{t})\|^{2}] \\ &\leq -\frac{\eta}{2} \mathbb{E}[\|\nabla f(w^{t})\|^{2}] + \frac{\eta}{2} (D^{2} + \gamma^{2} + 2D\gamma) \\ &+ \frac{\eta}{2} (\eta L - 1) \times (2\mathbb{E}[\|\nabla f(w^{t})\|^{2}] + 2D \mathbb{E}[\|\nabla f(w^{t})\|] + 2\gamma^{2} + D^{2} + 2D\gamma) \\ &= \frac{\eta}{2} (2\eta L - 3) \mathbb{E}[\|\nabla f(w^{t})\|^{2}] + \eta D(\eta L - 1) \mathbb{E}[\|\nabla f(w^{t})\|] + \frac{\eta}{2} (2\eta L\gamma^{2} - \gamma^{2} + \eta LD^{2} + 2\eta LD\gamma). \end{aligned} \tag{36}$$

937 When $2\eta L - 3 < 0$, i.e. $\eta < \frac{3}{2L}$, we have $(2\eta L - 3)\mathbb{E}[\|\nabla f(w^t)\|^2] \le (2\eta L - 3)(\mathbb{E}[\|\nabla f(w^t)\|)^2$. 938 In this case, Equation (29) can be written as: $f(w^{t+1}) - f(w^t)$

$$\leq \frac{\eta}{2} (2\eta L - 3) (\mathbb{E}[\|\nabla f(w^{t})\|])^{2} + \eta D(\eta L - 1) \mathbb{E}[\|\nabla f(w^{t})\|] + \frac{\eta}{2} (2\eta L \gamma^{2} - \gamma^{2} + \eta L D^{2} + 2\eta L D\gamma)$$
(37)

Which successfully proves Theorem 1. Furthermore, if we take $\mathbb{E}[\|\nabla f(w^t)\|$ as a variable, $f(w^{t+1}) - f(w^t)$ is deemed to be upper bounded by a **polynomial function** of $\mathbb{E}[\|\nabla f(w^t)\|$. In this case, we can naturally find $\epsilon_1 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}$ by letting the polynomial function equal to zero, with $a = 2\eta L - 3$, $b = 2D(\eta L - 1)$ and $c = 2\eta L\gamma^2 - \gamma^2 + \eta LD^2 + 2\eta LD\gamma$.

After calculation, we can get $\epsilon_1 =$

$$\frac{D(\eta L - 1) + \sqrt{\eta D^2 L + 8\eta L \gamma^2 + 6\eta D L \gamma + D^2 - 3\gamma^2}}{3 - 2\eta L}$$
(38)

Similarly, when the learning rate $\eta \leq \frac{1}{L}$, we have $\eta L - 1 \leq 0$. In this case, $\frac{\eta}{2}(\eta L - 1)\mathbb{E}[\|\nabla f'_k(w^t)\|^2]$ is naturally upper bounded by zero, so that Equation (29) can be written as:

$$f(w^{t+1}) - f(w^t) \le \frac{\eta}{2} (D^2 + \gamma^2 + 2D\gamma) - \frac{\eta}{2} \mathbb{E}[\|\nabla f(w^t)\|^2].$$
(39)

Which successfully proves Theorem 2. By letting $\frac{\eta}{2}(D^2 + \gamma^2 + 2D\gamma) - \frac{\eta}{2}\mathbb{E}[\|\nabla f(w^t)\|^2]$ equal to zero we naturally get $\epsilon_2 = D + \gamma$.

C SUPPLEMENTRAY EXPERIMENT RESULTS

C.1 COMPUTATION AND COMMUNICATION OVERHEADS

The single-sided comparison of computation cost (in FLOPs) and communication cost (in size of data transmission) are shown in Figures 11 and 12.

967 C.2 ENERGY CONSUMPTION

The overall energy consumptions, including the computation energy for training, and the communication energy for parameter transmission, are shown in Figure 13. The energy consumption is measured using a plug-in power monitor ².

²https://www.amazon.com.au/Electricity-Monitor-PIOGHAX-Overload-Protection/dp/B09SFSB66M

972 973 Total number of FLOPs Total number of FLOPs 974 Total number of FLOPs 42 975 6.2 28 (40 (₀01× u) 976 6.0 (₀1 26 5.8 977 ij) SdOTJ9 34 978 979 5.2 32 980 5. Feedorf+TOND 5) 30 Cocofi FedOLF FedOLF+TOAIO FedOLF 55 Fedolif+TOA cocs Fedolf+TOAL Depth Heter oe edolf*T 981 FedOLS 982 (a) EMNIST. (b) CIFAR-10. (c) CIFAR-100 on ResNet20. 983 Total number of FLOP Total number of FLOPs Total number of FLOPs 984 34 20 985 32 12 18 986 (j 30 9 16 a 10 987 £ 14 0 1 26 12 988 24 989 22 990 edolf*T 991 dOLF*1 adolf . dol 992 993 (d) CIFAR-100 on ResNet44. (e) CINIC-10 on ResNet20. (f) CINIC-10 on ResNet44. 994 Figure 11: Comparison of the overall computation cost, which is measured in total Floating Point 995 Operations (FLOPs) of all clients. This is the average result for the three trials. 996 997 998 999 1000 of data trans Total size of data transmission Total size of data transmission 1001 15 20 1002 14 70 17.5 8 ¹³⁰ ⁽¹⁾
 1003 120 50 12.5 1004 40 110 10.0 1005 10 Data 7.5 Data 1006 5.0 2.5 10 1007 0.0 70 FedOLE TONOS FedOL COCOFL FedOLF *TORIO.51 *TOALOS Heteroft Depthft Cocoft 55 55 1008 FIORD Heteroft Depth e4015+TOA FIORE Heteroft Deptifit Cocoft ŝ FedOL edOLF*TOR cedd 1009 1010 (a) EMNIST. (b) CIFAR-10. (c) CIFAR-100 on ResNet20. 1011 size of data tra Total size of data transmission Total size of data transmission 160 170 1012 80 155 160 1013 70 150 150 60 1014 8 145 140 50 140 130 1015



1016

1017

1018

1019

1020 1021

1022

Figure 12: Comparison of the overall communication cost, which is measured in the total size of 1023 parameters transmitted across the network. This is the average result for the three trials. 1024 1025

edolf*



Figure 13: An overview of the overall energy consumption (kJ) of all clients, including the computation energy for local training (green) and the communication energy for global communication (yellow).

Dataset Model		EMNIST	CIFAR-10 CIFAR-100		R-100	CINIC-10	
		CNN	AlexNet	ResNet20	ResNet44	ResNet20	ResNet44
Feddrop		16.42	14.33	6.05	6.15	9.71	10.81
FjORD		12.68	27.8	11.14	9.09	22.22	11.26
HeteroFL		12.88	58.03	7.02	14.32	13.46	12.0
DepthFL		83.0	10.52	5.05	39.88	10.31	33.44
CoCoFL		81.98	53.92	26.95	31.1	31.81	31.68
SLT		81.04	49.73	45.80	39.60	21.63	36.20
	no TOA	84.98	66.98	48.49	44.12	40.66	37.33
FedOLF	TOA(0.75)	-	63.7	40.49	42.16	33.96	31.51
	TOA(0.5)	-	62.05	36.19	38.29	33.42	28.42
FedAvg		85.04	68.41	51.11	52.13	40.80	39.88

Table 2: Comparison of the final test accuracy (in %) for T = 500 iterations in the iid case.

1070 C.3 ACCURACY IN THE IID CASE 1071

The accuracy comparison in the iid case is listed in Table 2. As shown in Table 2, FedOLF still outperforms the baselines in the iid case, and maintains a competitive accuracy against the FedAvg benchmark.

1077 C.4 ACCURACY VS. ROUND

The curves of accuracy with respect to training rounds are shown in Figures 14 and 15.

1080 D DEMONSTRATION OF LAYER FREEZING IN FEDOLF

 Figures 16a, 16b, 16c illustrate how FedOLF freezes layers among the heterogeneous clients. Each bar represents a cluster of clients, dividing the model into two segments. On one side of the bar, denoted by "F," clients freeze the corresponding layers, while on the other side, denoted by "T,"

clients actively train the model. Specifically, clients within Cluster 2 for EMNIST and Cluster 5 for CIFAR-10/CIFAR-100/CINIC-10 are assumed to possess the capability to train the entire model. It is important to note that convolution layers are typically followed by activation functions (e.g., ReLU), batch normalization or max-pooling layers, though these are not depicted in the figures for brevity.



Figure 14: Accuracy vs. round in the non-iid case.





Figure 15: Accuracy vs. round in the iid case.





Figure 16: A specific demonstration of how FedOLF freezes layers among the heterogeneous clients.

¹¹⁸⁸ E DISCUSSION

1190 E.1 LIMITATIONS

1192 One of the major limitations of this paper is the lack of theoretical support of TOA's application on FedOLF. Even though TOA seems to work well based on the obvious reduction of energy con-1193 sumption/memory footprint as shown in the experiment results, the specific relationship between 1194 performance degradation and the scaling factor s remains unexplored. In this case, the only way 1195 to determine the optimal value of s is through an experiment (i.e. hyperparameter tuning), which 1196 is usually computationally expensive. To enhance the usefulness of FedOLF w. TOA in practical 1197 applications, a close-formed representation of the effect of s on accuracy is required, so that we can 1198 determine the optimal value of s given the particular energy budget and accuracy requirement. 1199

The other limitation of FedOLF is the additional communication overhead as shown in Figure 12. 1200 Even with TOA, the communication overhead of FedOLF is still higher than other methods in most 1201 cases. In the environment of our experiment, the connection between clients and the server is rela-1202 tively stable, so that the extra communication overhead of FedOLF does not generate too much en-1203 ergy consumption. However, in a real-world system with underprivileged network conditions, such 1204 as a mobile-edge network Jin et al. (2024), the negative impact of the increased communication over-1205 head becomes severe, resulting in much higher communication costs. To promote the application 1206 of FedOLF in bandwidth-constrained systems, addressing the concern of increased communication 1207 cost becomes a vital matter.

1208

1209 E.2 BROADER IMPACT 1210

FedOLF has a positive social impact on boosting fair FL training among heterogeneous clients. FedOLF proposes that powerful clients take more responsibility in training (i.e. train more layers), and share the low-level layers with weak clients for forwardpropagation in their local training tasks. This significantly improves FL's accuracy, efficiency and robustness in resource-constrained settings.

Furthermore, FedOLF alleviates privacy concerns compared with traditional FL frameworks such as
Fedavg McMahan et al. (2017). As resource-constrained clients only communicate the active layers
with the server. Compared with the full model, transmitting partial active layers reduces the risks of
several types of attacks such as byzantine attack and privacy inference Hao et al. (2021); Cao et al.
(2020); Zhang et al. (2022b).

As for the negative impact, the increased communication overhead might restrict the usefulness of FedOLF in mobile networking systems with insufficient bandwidth support. Except for TOA, possible solutions to address this problem include: 1. Periodical downward communication:
Clients download the frozen layers periodically rather than every round. 2. Clustering: Clients download from a proximal header rather than the remote server.

- 1226 1227 1228 1229 1230 1231 1232 1233 1234 1235 1236 1237
- 1238
- 1230
- 1240
- 1241