FedRolex: Model-Heterogeneous Federated Learning with Rolling Sub-model Extraction

Anonymous Author(s) Affiliation Address email

Abstract

1	Federated learning (FL) is a collaborative machine learning paradigm to train
2	models from decentralized private data. Most FL research focuses on model-
3	homogeneous setting where models deployed across all the participating clients
4	and the server are required to be identical. However, in real-world scenarios,
5	such a requirement acts as a constraint that restricts the outreach to clients with
6	heterogeneous device resources and unfairly excludes users with low-end devices
7	who would otherwise benefit from FL. In this work, we propose a simple yet
8	effective model-heterogeneous FL method named FedRolex to tackle this con-
9	straint. Unlike the model-homogeneous scenario, the fundamental challenge of
10	model heterogeneity in FL is that different parameters of the global model are
11	trained on heterogeneous data distributions. Our method addresses this challenge
12	by rolling the sub-model in each federated iteration so that the parameters of the
13	global model are evenly trained on the global data distribution across all devices,
14	making it more akin to model-homogeneous training. Our experiments show that
15	FedRolex outperforms other state-of-the-art model-heterogeneous FL methods,
16	especially under high data-heterogeneity scenarios. We have conducted ablation
17	studies to show that submodel rolling is an effective technique to reduce the gap
18	between model-heterogeneous and standard model-homogeneous settings. Lastly,
19	we consider the distribution of client capabilities that is similar to real-world in-
20	come distribution instead of the uniform distribution used in existing works. Our
21	results show a consistent improvement in the accuracies on low-end devices, which
22	enhances the inclusiveness of federated learning.

23 1 Introduction

Cross-device federated learning (FL) was originally proposed by McMahan et al. (2017) as a privacypreserving machine learning paradigm to train a machine learning model on a federation of resourceconstrained edge devices without accessing their data. This idea has far-reaching applications, which are made urgent by modern demands for privacy. Indeed, the goal of strong data privacy is a central motivation for FL. By storing data locally, instead of replicating them on a remote server, the attack surface of the system is decreased (Jere et al., 2020), and by using focused ephemeral updates and early aggregation, the communication cost is also minimized (Nishio and Yonetani, 2019).

Majority of the existing cross-device FL studies focus on *model-homogeneous* setting, where models deployed across all the participating clients and the server are required to be identical. In practice, device heterogeneity is a realistic constraint to consider for deploying in real-world environments. Different edge devices have different on-device resources and, hence, are only capable of training models with capacities that match their on-device resources. Having the same model on all the devices forces a limitation that restricts the outreach to clients with heterogeneous device resources and unfairly excludes users with low-end devices who would otherwise benefit from FL.

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

Table 1: Comparison of FedRolex with model-homogeneous and model-heterogeneous FL methods.

	Model Heterogeneity	Need of Public Data	Scaling to Large Federations	Server Model Size
FedAvg McMahan et al. (2017)	No	No	Scalable	= Client Model
FedProx Li et al. (2020)	No	No	Scalable	= Client Model
SCAFFOLD Karimireddy et al. (2020)	No	No	Scalable	= Client Model
FedBE Chen and Chao (2020)	No	Unlabeled	Computationally Expensive	= Client Model
FedDF Lin et al. (2020)	Yes	Unlabeled	Computationally Expensive	= Largest Client Model
FedGKT He et al. (2020)	Yes	No	Computationally Expensive	≥ Largest Client Model
FedGEMS Cheng et al. (2021)	Yes	Labeled	Computationally Expensive	≥ Largest Client Model
Fed-ET Cho et al. (2022)	Yes	Unlabeled	Computationally Expensive	≥ Largest Client Model
Federated Dropout Caldas et al. (2018)	Yes	No	Large Federation Required	≥ Largest Client Model
HeteroFL Diao et al. (2020)	Yes	No	Scalable	= Largest Client Model
FjORD Horvath et al. (2021)	Yes	No	Scalable	= Largest Client Model
FedRolex	Yes	No	Scalable	\geq Largest Client Model

To tackle this constraint, recent works proposed model-heterogeneous FL, where models with different 38 capacities are trained on edge devices with heterogeneous resources during the federated training 39 process. One fundamental challenge in model-heterogeneous FL is the aggregation of heterogeneous 40 client models. Lopes et al. (2017); He et al. (2020); Cho et al. (2022) proposed to use knowledge 41 distillation (KD) (Hinton et al., 2015), where the client models serve as teachers and server ensembles 42 the knowledge distilled from the individual client models. However, KD-based methods generally 43 require access to public data to achieve competitive model accuracy, where the accessibility of high-44 quality public data is not always practical. Moreover, KD becomes computationally expensive when 45 the cohort size in each FL round scales up. To remove the dependency on public data, another set of 46 studies focuses on designing partial training (PT) based techniques such as federated dropout (Caldas 47 et al., 2018), ordered dropout (Horvath et al., 2021), and static sub-model extraction (Diao et al.) 48 2020), where each client trains a smaller sub-model of the larger global model on the server. However, 49 the fundamental issue of existing PT-based methods is that the sub-models are extracted in ways such 50 that the parameters of the global model obtained by aggregating the extracted sub-models are not 51 evenly trained, making the global model vulnerable to a phenomenon called *client drift* (Wang et al.) 52 2021). This phenomenon becomes more prominent when the cohort size in each communication 53 round is small or data heterogeneity across clients becomes more severe, resulting in lower global 54 model accuracy. 55

In this work, we propose a simple yet effective PT-based model-heterogeneous FL method named 56 57 FedRolex that tackles the fundamental issue of existing PT-based methods. Compared to existing PT-based methods, FedRolex proposes a rolling sub-model extraction scheme, where the sub-model 58 is extracted from the global model using a rolling window that advances in each communication round. 59 As such, sub-models are extracted from different parts of the global model in different rounds. Since 60 the extraction window is rolling in each round, parameters of the global model are evenly trained to 61 minimize the client drift. Throughout communication rounds, all the parameters of the global model 62 are updated over the entire data distribution. As we show in §4.2, such rolling sub-model extraction 63 strategy significantly reduces the gap between model-heterogeneous and model-homogeneous setting. 64 We evaluate the performance of FedRolex on two datasets - CIFAR-10 and CIFAR-100 - and 65

compare it against state-of-the-art KD-based model-heterogeneous FL methods FedDF (Lin et al., 66 2020), DS-FL (Itahara et al., 2020) and Fed-ET (Cho et al., 2021) as well as PT-based model-67 heterogeneous FL methods HeteroFL (Diao et al., 2020) and Federated Dropout (Caldas et al., 2018) 68 Our results show that FedRolex consistently outperforms state-of-the-art PT-based methods in both 69 high and low data heterogeneity scenarios. FedRolex also outperforms state-of-the-art KD-based 70 methods in low data heterogeneity scenario as well as on the more challenging CIFAR-100 dataset 71 72 in high data heterogeneity scenario without public data. Furthermore, we show that our proposed 73 rolling sub-model extraction scheme is flexible to train a global model that can be much larger than the largest client model. Lastly, we conducted an experiment that uses real-world household income 74 distribution to emulate real-world device distribution. Our results show that FedRolex effectively 75 boosts the model accuracy of low-end devices, which enhances the inclusiveness of federated learning. 76 The code of FedRolex will be publicly available on Github. 77

¹We did not compare with FjORD because its code is not open-source and we could not reproduce their results following the paper.

78 2 Related Work

Our work focuses on model-heterogeneous FL where existing methods can be generally grouped into
 two categories: knowledge distillation (KD) based and partial training (PT) based methods.

Knowledge Distillation (KD) based Model-Heterogeneous FL. In knowledge distillation (Hinton 81 et al., 2015; Mirzadeh et al., 2020; Liu et al., 2019), the core idea is to compress a large pre-trained 82 model by teaching a smaller network, step by step, exactly what to do using the larger pre-trained 83 network. However, this has been adapted in FL systems to train a single model from a federation of 84 client models trained on private distributed data. In FedDF, (Lin et al., 2020) uses KD from a set of 85 classifiers pooled from a federation of client devices with their own private data. The logit outputs 86 of each of the classifiers against a public dataset input are then used to train a student model at the 87 server and thereby distill knowledge. In DS-FL (Itahara et al.) 2020), an unlabeled public dataset 88 is similarly used. The logit output for this dataset is aggregated and averaged at the server and is 89 broadcast back to the clients. Local models now train on the additional pseudo-labeled data, and 90 performance is enhanced owing to the data augmentation effect. Fed-ET (Cho et al., 2022) considers 91 a weighted consensus distillation approach with diversity regularization that enables training of a 92 large server model with smaller models at the clients. These models have several drawbacks. KD is 93 94 generally a computationally expensive task and so the costs blow up as the system scales up for larger models and more clients. These are also typically restricted to classification tasks as they depend on 95 logits for implementing KD. The methods are also likely to face deployment issues, especially for 96 tasks with a large number of possible output classes. 97 Partial Training (PT) based Model-Heterogeneous FL. There have been several studies employing 98

different methods to train models of different sizes. These generally follow PT. HeteroFL (Diao 99 et al., 2020) uses this concept to aggregate models of different sizes. Fjord (Horvath et al., 2021) uses 100 ordered dropout, which similarly samples submodels from a global model, trains them, and then does 101 a weighted aggregation. PruneFL (Jiang et al., 2022) uses model pruning to extract submodels. These 102 algorithms are easily scalable and have minimum overheads on the server. However, these methods 103 have a significant drawback in the sense that different parts of the model are restricted to seeing 104 updates from a fixed pool of client data. This limits the performance as the global knowledge is 105 concentrated on a smaller portion of the model. Furthermore, the server model capability is restricted 106 to the same capability as the largest client model. Servers are usually much more computationally 107 capable and so should be able to use a larger model. In comparison, Federated Dropout (Caldas 108 et al., 2018; Guliani et al., 2022; Ro et al., 2022) can train on larger global models and use its entire 109 110 capability, especially when the dataset and the number of clients are large. Here, a portion of the parameters in each layer of the global model is dropped, leaving a smaller submodel. Submodels are 111 trained on the client devices according to their capacities and the updates are aggregated via weighted 112 averaging similar to traditional dropout. However, as reported by Bouacida et al. (2020), the approach 113 works well only with less heterogeneous data and large client pools. This is because of high variance 114 in heterogeneous data when only a small subset of the client pool is selected. 115

116 **3 Our Method**

117 3.1 Problem Formulation of Model-Heterogeneous FL

Consider N different client devices, \mathcal{N} , with different compute capacities $\{C_1, C_2, ..., C_N\}$ and private data $\{D_1, D_2, ..., D_N\}$, each with different distributions. The objective is to train a large global model characterized by the set of trainable model parameters θ . Therefore, model-heterogeneous FL can be formulated as the following distributed optimization problem:

$$\min_{\theta} \left\{ F\left(\theta\right) \triangleq \sum_{n=1}^{N} p_n F_n(\theta) \right\}$$
(1)

where $F(\theta) = \mathbf{E}_{\xi \sim D} l(\theta; \xi)$, $F_n(\cdot)$ is the local objective function on the client dataset $D_n = \{d_{n,1}, d_{n,2}, d_{n,3}...d_{n,m_n}\}$ with respect to the user defined loss function $l(\cdot; \cdot)$ and p_n is the corresponding weight of the n^{th} client such that $p_n \ge 0$ and $\sum_{n=1}^{N} p_n = 1$. The n^{th} device contains m_n data points and has model parameters θ_n (this can be changed from one iteration to another). The size of θ_n depends on the client's model capacity C_n and are extracted from the global model i.e.,



Figure 1: Illustrated overview of rolling window for round numbers j = 0, 1, and 2 for client capacities set $\{1, \frac{1}{2}, \frac{1}{4}\}$. The orange squares indicate kernels extracted and trained on client device whereas the gray ones indicate deselected kernels. For $\beta_n = 1$, all kernels are used for all convolution layers in all rounds. For $\beta_n = \frac{1}{2}$, the first 4 kernels are extracted for the 1st convolution layer and the first 8 are extracted from the next at round j = 0. In round j = 1, kernels 1 through 5 are selected from the first convolution layer and 1 through 9 for the second. For j = 2, kernels 2 through 6 are selected from the first convolution layer and 2 through 10 for the second. This is similarly played out for $\beta_n = \frac{1}{4}$ except here only 2 kernels are selected from the first convolution layer and 4 from the second.

127 $\theta_n \subset \theta$. The local objective function $F_n(\cdot)$ is thereby rewritten as:

$$F_n(\theta_n) \triangleq \frac{1}{m_n} \sum_{k=1}^{m_n} l(\theta_n; d_{n,k}).$$
⁽²⁾

For simplicity, we abuse the notation of l by ignoring the other parameters in θ but not in θ_n . The core of the optimization problem is to select subset θ_n from the global model parameters θ based on model capacity C_n .

131 3.2 FedRolex: Model-Heterogeneous FL with Rolling Sub-model Extraction

We propose a more optimal way to extract and train the sub-models that are more consistent with 132 practical applications than previous methods. In our work, we extract subsets from a rolling window 133 that advances and loops over all the kernels of each convolution layer in the model. Consider 134 $\theta_n^{(j)}$ selected from θ for client n at the j^{th} communication round. Assume that the proportion of kernels extracted from global parameter set θ at client n be β_n , which belongs to the set of unique model capacities β . Let the total number of kernels in the i^{th} convolution layer be K_i . Then 135 136 137 the parameters $\theta_{n,[i,\cdot]}^{(j)}$ for the i^{th} convolution layer in extracted model parameters $\theta_n^{(j)}$ consist of 138 $\{\theta_{n,[i,j \mod K_i]}^{(j)}, \theta_{n,[i,(j+1) \mod K_i]}^{(j)}, \dots, \theta_{n,[i,(j+\beta_n K_i) \mod K_i]}^{(j)}\}$ where mod is the modulus operator. For j = 0, the window starts from the very first kernel index and kernel parameters for the sub-model 139 140 are extracted starting from index j = 0. After each round, the value j is incremented by 1, and so the 141 window advances forward over each round and loops back from the beginning after one complete 142 cycle through all the kernels of a particular convolution layer. Over a sufficiently large number of 143 communication rounds, all parts of the model will see the entire data distribution across the N clients. 144 The process over three rounds is illustrated in Figure 1. 145

Let, $\mathcal{M} \subset \mathcal{N}$ be the set of selected clients from the client pool from which the server pulls model parameters at round *j*. Let $\theta_{[i,k]}$ be the k^{th} trainable weight of the i^{th} layer of global model and $\theta_{m,[i,k]}; m \in \mathcal{M}$ be the k^{th} trainable weight of the i^{th} layer of the client *m*. The model parameters are aggregated as follows:

$$\theta_{[i,k]} = \frac{1}{\sum_{m \in \mathcal{M}} p_m} \sum_{m \in \mathcal{M}} p_m \theta_{m,[i,k]}$$
(3)

Here, p_m is an empirical weight given to a client. This can be based on a variety of factors like the client model capability, the number of data points a client has, etc. However, there is no theoretical foundation in the literature to quantify the relative importance of a client update nor are there empirical studies that accurately study this. We consider this tangent to our goals in this study and so, throughout the paper, unless otherwise stated, the weight of all clients are assumed to be the same, i.e, $p_m = 1$. Finally, the optimal sub-model is extracted from the global model for each client according to client capability.

Why Rolling Sub-model Extracting is better? PT-based algorithms HeteroFL and FjORD extract 157 the model parameters θ_n from specific locations in the global model. Here the number of kernels 158 extracted from each convolution layer for creating a client model is based on the client capacity, and 159 the kernels are selected starting from the first index. The extracted sub-models are then trained on 160 selected client devices from the client pool, and the updates for each parameter are aggregated and 161 averaged based on the number of clients that participated in training for that parameter. The kernels 162 in the beginning indices of each convolution layer will therefore see updates that were averaged 163 across more clients than end indices. This means the expected value of the updates will tend more 164 toward the optimal update. In comparison, the kernels at the end indices will see updates averaged 165 across much fewer clients, and therefore the updates to these kernels will have more uncertainty. 166 Essentially, the kernels will be unevenly trained. Dropout-based algorithms like Federated Dropout 167 Caldas et al. (2018) propose random dropout to create smaller sub-models. Though the expected 168 value of the frequency for updating each index is the same for all the indices in Federated Dropout, 169 their exact frequencies are not the same in each experiment because of the randomness. In fact, given 170 I indices, if one index is chosen at each round, the probability to update all indices once in I rounds is 171 $\prod_{i=1}^{l} \frac{i}{I}$. Furthermore, the expected number of rounds to go through all indices at least once is close 172 to $I \ln(I)$, see Appendix A.1 It shows that random dropout can not balance the update frequencies 173 of the indices, and it takes a large number of rounds to update all the indices. 174

Algorithm 1: FedRolex

Initialization ; $\theta^{(0)}$, \mathcal{N} 1 Input : $D_n \ \beta_n \ \forall n \in \mathcal{N}$, Output : θ^{J} 2 Server Executes 3 for $j \leftarrow 0$ to J do Broadcast $\theta_{n,[i,j \mod K_i \dots (j+\beta_n K_i) \mod K_i]}^{(j)} \forall i \text{ and } n \in \mathcal{N}$ 4 Sample subset \mathcal{M} from \mathcal{N} 5 for each client $m \in \mathcal{M}$ do 6 clientStep($\theta_n^{(j)}, D_n$) 7 end 8 Aggregate $\theta_{[i,k]}^{(j+1)} = \frac{1}{\sum_{m \in \mathcal{M}} 1} \sum_{m \in \mathcal{M}} \theta_{m,[i,k]}$ Ģ 10 end 11 Subroutine clientStep($\theta_n^{(j)}, D_n$) $\begin{array}{l} m_n \longleftarrow len(D_n) \\ \text{for } k \leftarrow 0 \text{ to } m_n \text{ do} \\ \mid \theta_n \longleftarrow \theta_n - \eta \nabla l(\theta_n; d_{n,k}) \end{array}$ 12 13 14 end 15 return θ_n 16

175 4 Experiments and Analysis

176 4.1 Experimental Setup

Datasets and Platform. We validate the performance of FedRolex on two public open-source datasets – CIFAR-10 and CIFAR-100 (Krizhevsky, 2009). We apply bounding box crop (Zoph et al., 2020) to augment the images. To simulate the non-iid nature of the data, we use balanced non-iid data distribution following HeteroFL. We skew the data distribution by restricting the clients to having

a fixed number of unique labels L. We define L = 2 as high data heterogeneity (i.e., high non-iid) 181 and L = 5 as low data heterogeneity (i.e., low non-iid) for CIFAR-10. Similarly, we use L = 20182 as high non-iid and L = 50 as low non-iid for CIFAR-100. These roughly correspond to Dirichlet 183 distribution $Dir_K(\alpha)$ with α set to 0.1 and 0.5, respectively. The data is partitioned into 100 non-iid 184 partitions and assigned to the client pool. In each communication round, 10% of the clients are 185 randomly selected from the client pool of 100. We run our experiments on 8 Nvidia A6000 GPUs 186 187 using the PyTorch library for machine learning and the Ray library for simulating the distributed system. 188

Models. We use pre-activated ResNet18 (PreResNet18) models (He et al., 2016). We replace the 189 batch Normalization in PreResNet18 with static batch normalization (Diao et al., 2020; Andreux 190 et al., 2020) and add a scaler module after each convolution layer (Diao et al., 2020). We list the exact 191 architecture of the model in Appendix A.2. To simplify our experiments and ablation studies, the 192 model sizes are restricted to a set of 5, i.e., β_n would take any value in $\{1, \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}\}$. Note that 193 we only vary the number of kernels in convolution layers in our method, HeteroFL, and Federated 194 Dropout to keep the comparison fair. We provide pseudocode for both HeteroFL and our Federated 195 Dropout variant in Appendix A.3 Unless otherwise stated, the client capacity distribution is assumed 196 to be uniform over the entire client pool for all methods. 197

Baselines. We compare FedRolex against state-of-the-art PT-based methods HeteroFL and Federated
Dropout and state-of-the-art KD-based methods FedDF, DS-FL, and Fed-ET obtained from Lin et al.
(2020). For a fair comparison, all the baselines from the PT-based methods are trained on the same
learning rate, the same number of communication rounds, and the same multi-step learning rate decay
schedule. The exact schedule for each dataset and experiment is given in Appendix A.2

Evaluation Metrics. We use local and global model accuracy as our evaluation metrics. Specifically, global model accuracy is defined as the mean accuracy on the test set totally held out during training; local model accuracy is defined as the mean accuracy of the server model on each of the client datasets. Since there is some randomness due to the non-iid partitioning, we run our experiments on 5 different seeds and list the mean and standard deviation of our metrics. Note that the results of KD-based methods were obtained from Cho et al. (2022) using 3 different seeds.

209 4.2 Performance Comparison with State-of-the-Art Model-Heterogeneous FL Methods

Table 2 shows the global model accuracy comparison between FedRolex and the baselines under 210 both high and low data heterogeneity settings. We also include the global model accuracy of two 211 model-homogeneous cases where all clients have the highest capacity ($\beta = 1.0$) and lowest capacity 212 $(\beta = 0.0625)$ representing the upper and lower-bound performance, respectively. We have three 213 key observations. (1) In comparison with state-of-the-art PT-based methods, FedRolex consistently 214 outperforms HeteroFL and Federated Dropout in both high and low data heterogeneity scenarios. In 215 particular, under high data heterogeneity, Federated Dropout that uses random dropout has drastically 216 worse performance than FedRolex and HeteroFL which both extract sub-models in a deterministic 217 manner. (2) In comparison with state-of-the-art KD-based methods, FedRolex only performs worse 218 than Fed-ET and FedDF on CIFAR-10 under high data heterogeneity, but outperforms all the state-of-219 the-arts in more challenging CIFAR-100 dataset that has many more classes than CIFAR-10 under 220 both high and low data heterogeneity scenarios. It is important to emphasize that KD-based methods 221 leverage public data to boost their model accuracy while FedRolex does not. (3) In comparison with 222 model-homogeneous cases, compared to other PT-based methods, FedRolex considerably reduces 223 the gap between model-heterogeneous and model-homogeneous setting, particularly on CIFAR-10 224 under low data heterogeneity. Note that Fed-ET achieves a higher global model accuracy than the 225 226 model-homogeneous upper bound on CIFAR-10 under high data heterogeneity, which showcases the 227 advantage of having access to public data.

228 4.3 Impact of Model Capacity Distribution

In our previous experiment, in line with existing model-heterogeneous FL studies, we have performed our experiment where the distributions of model capacities are uniform for a fair comparison. In practical scenarios, however, the distributions of model capacities will likely be skewed. Therefore, we conduct an experiment to understand how global model accuracy changes when the proportion of higher capacity clients varies. To do so, we set the proportion of large to small models to ρ and the

	Method	High Data H	Ieterogeneity	Low Data He	Low Data Heterogeneity	
	method	CIFAR-10	CIFAR-100	CIFAR-10	CIFAR-100	
KD based	FedDF	73.81 (± 0.42)	31.87 (± 0.46)	76.55 (± 0.32)	37.87 (± 0.31)	
	DS-FL	65.27 (± 0.53)	29.12 (± 0.51)	68.44 (± 0.47)	33.56 (± 0.55)	
	Fed-ET	78.66 (± 0.31)	35.78 (± 0.45)	81.13 (± 0.28)	41.58 (± 0.36)	
PT based	HeteroFL	63.90 (± 2.74)	52.38 (± 0.80)	73.19 (± 1.71)	57.44 (± 0.42)	
	Federated Dropout	46.64 (± 3.05)	45.07 (± 0.07)	76.20 (± 2.53)	46.40 (± 0.21)	
	FedRolex	69.44 (± 1.50)	56.57 (± 0.15)	84.45 (± 0.36)	58.73 (± 0.33)	
	Homogeneous (Smallest)	38.82 (± 0.88)	12.69 (± 0.50)	46.86 (± 0.54)	19.70 (± 0.34)	
	Homogeneous (Largest)	75.74 (± 0.42)	60.89 (± 0.60)	84.48 (± 0.58)	62.51 (± 0.20)	

Table 2: The global model accuracy comparison on CIFAR-10 and CIFAR-100 under both high and low data heterogeneity.

set of unique client capacities is $\beta = \{1, 1/16\}$ where $\rho = 1$ represents the edge case in which the models of all clients have the architecture of the full ResNet18 model, and $\rho = 0$ is the other edge case in which all clients have the smallest model.

Figure 2 shows how global model accuracy changes when ρ varies from 0 to 1 for CIFAR-10 and 237 CIFAR-100 under both high and low data heterogeneity. We have three key observations. (1) For 238 CIFAR-10, there is a large gap in global model accuracy between high and low data heterogeneity 239 for a wide range of ρ (from 0.1 to 1). This is because CIFAR-10 is a relatively simple task and 240 hence the global model accuracy is bottlenecked by the level of data heterogeneity instead of model 241 capacity. This result indicates that having more high-capacity models in the cohort has limited 242 contribution to global model accuracy. (2) For the more challenging CIFAR-100, the gap in global 243 model accuracy is much lower between high and low data heterogeneity. In contrast to CIFAR-10, the 244 global model accuracy is bottlenecked by the highest capacity of the models rather than the level of 245 data heterogeneity. (3) For both CIFAR-10 and CIFAR-100, we observe that having a small fraction 246 of high-capacity models is enough to significantly boost the global model accuracy. 247



Figure 2: Impact of model capacity distribution on global model accuracy under low and high data heterogeneity for (i) CIFAR-10 and (ii) CIFAR-100.

248 4.4 Training Large Server Models

One limitation of HeteroFL and FjORD is that the server model size is restricted to the highest capacity client model. In contrast, FedRolex is flexible to train a global model that can be much larger than the largest client model. This is more akin to real-world settings where the server has a much higher capacity.

We consider the case where the global model size is γ times the size of the highest capacity client model. Figure 3 shows the global model accuracy of $\gamma = \{2, 4, 8, 16\}$ for CIFAR-10 and CIFAR-100 under high and low data heterogeneity. As shown, for both CIFAR-10 and CIFAR-100, the global model accuracy drastically drops when γ increases from 2 to 4 but stay relatively stable when γ increases from 4 to 8 and from 8 to 16. This result indicates that there is a minimum client model capacity below which the server model cannot be trained effectively.



Figure 3: The global model accuracy when the server model is γ times the size of the largest client model under low and high data heterogeneity for (i) CIFAR-10 and (ii) CIFAR-100.

259 4.5 Enhance Inclusiveness of FL in Real-world Distribution

A primary goal of our method is to enhance the inclusiveness of FL. If the highest model capacity is 260 used for deployment, a lot of data from the low-end devices will be unused. The exclusion of low-end 261 devices limits the usefulness of the global model for these low-end devices. For example, low-end 262 devices may have more grainy, low-resolution photos which will not be used in training, and so, the 263 server model cannot make accurate predictions on such images and hence limiting its usefulness 264 for these devices. If we consider homogeneous FL where we want to include all the client data in 265 training, the smallest model would need to be used. The accuracy of the model is then limited by the 266 lowest capacity client. 267

In this experiment, we show that FedRolex is a valid solution that enhances the inclusiveness of federated learning. To illustrate our point, we conducted an experiment that uses real-world household income distribution to emulate real-world device distribution. Specifically, we retrieve household income distribution information from Bureau (2021). We map $\beta_n = 1/16$ with the income group with earning less than \$75000 and assign proportions of remaining groups in \$25000 increments with increasing values of β_n in the statistic. Detailed mapping of this distribution to the corresponding income distribution is given in Appendix A.2

	Method	High Het	erogeneity	Low Heterogeneity	
Dataset		Local Accuracy	Global Accuracy	Local Accuracy	Global Accuracy
CIFAR-10	Homogeneous (smallest)	85.90 (± 0.46)	38.82 (± 0.88)	66.02 (± 0.52)	46.86 (± 0.54)
	Homogeneous (largest)	95.54 (± 0.26)	75.74 (± 0.41)	93.54 (± 0.44)	84.48 (± 0.58)
	FedRolex	94.05 (± 1.01)	63.17 (± 1.45)	91.03 (± 0.36)	80.14 ± 0.52)
CIFAR-100	Homogeneous(smallest)	34.51 (± 0.56)	12.69 (± 0.50)	33.22 (± 0.10)	19.70 (± 0.34)
	Homogeneous(largest)	81.99 (± 0.78)	60.89 (± 0.60)	76.43 (± 0.54)	62.51 (± 0.20)
	FedRolex	73.33 (± 0.96)	45.78 (± 1.71)	66.31 (± 0.34)	48.44 (± 0.51)

Table 3: Performance of FedRolex under Real-World Distribution.

Table 3 shows the global and local model accuracy of FedRolex on this emulated setting compared 275 to two model-homogeneous cases where all clients have the smallest and largest model capacities, 276 representing lower and upper bound on accuracy respectively. We make two key observations (1) 277 FedRolex consistently and significantly outperforms the lower bound model-homogeneous case 278 in terms of local model accuracy for CIFAR-10 and CIFAR-100 under both high and low data 279 heterogeneity. This result indicates that FedRolex is able to effectively boost the local model 280 accuracy of low-end devices, which would otherwise not benefit from FL. (2) FedRolex achieves 281 a local model accuracy that is close to the upper bound model-homogeneous case, especially for 282 CIFAR-10 under both high and low data heterogeneity. This result indicates that FedRolex is able 283 to significantly reduce the gap with the best-performing scenario without being constrained to only 284 using high-end devices. 285



Figure 4: The local model accuracy distribution of FedRolex (in orange color) and the smallest homogeneous case (in blue color) for CIFAR-10 and CIFAR-100 under low and high data heterogeneity.

Finally, Figure A provides a detailed illustration of how the local model accuracy shifts for CIFAR-10 and CIFAR-100 under low and high data heterogeneity when FedRolex is used compared to the smallest homogeneous case with the same client outreach. As shown, the local model accuracy of individual clients is significantly boosted by FedRolex.

290 5 Conclusion

In this work, we present a partial training (PT)-based model-heterogeneous method named FedRolex 291 that tackles the fundamental issue of existing PT-based methods. At the core of FedRolex is a rolling 292 sub-model extraction scheme that enables parameters of the global model to be evenly trained to 293 minimize the client drift. Our experimental results show that FedRolex consistently outperforms 294 state-of-the-art PT-based methods in both high and low data heterogeneity scenarios. We demonstrate 295 its effectiveness in real-world scenarios by showing its performance on a more practical model 296 capacity distribution and show FedRolex contributes to making FL more inclusive. FedRolex 297 however is not tested on other tasks like language modeling and is left for future work. Future 298 directions, therefore, include validating its performance on different kinds of tasks. 299

300 **References**

- 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*, pages 1273–1282. PMLR, 2017.
- Malhar S Jere, Tyler Farnan, and Farinaz Koushanfar. A taxonomy of attacks on federated learning.
 IEEE Security & Privacy, 19(2):20–28, 2020.

Takayuki Nishio and Ryo Yonetani. Client selection for federated learning with heterogeneous
 resources in mobile edge. In *ICC 2019-2019 IEEE international conference on communications* (*ICC*), pages 1–7. IEEE, 2019.

- 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*,
 2:429–450, 2020.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and
 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In
 International Conference on Machine Learning, pages 5132–5143. PMLR, 2020.
- Hong-You Chen and Wei-Lun Chao. Fedbe: Making bayesian model ensemble applicable to federated
 learning. *arXiv preprint arXiv:2009.01974*, 2020.
- Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model
 fusion in federated learning. *Advances in Neural Information Processing Systems*, 33:2351–2363,
 2020.
- Chaoyang He, Murali Annavaram, and Salman Avestimehr. Group knowledge transfer: Federated
 learning of large cnns at the edge. *Advances in Neural Information Processing Systems*, 33:
 14068–14080, 2020.

Sijie Cheng, Jingwen Wu, Yanghua Xiao, and Yang Liu. Fedgems: Federated learning of larger server models via selective knowledge fusion. *arXiv preprint arXiv:2110.11027*, 2021.

Yae Jee Cho, Andre Manoel, Gauri Joshi, Robert Sim, and Dimitrios Dimitriadis. Heterogeneous ensemble knowledge transfer for training large models in federated learning. *arXiv preprint arXiv:2204.12703*, 2022.

- Sebastian Caldas, Jakub Konečny, H Brendan McMahan, and Ameet Talwalkar. Expanding the reach
 of federated learning by reducing client resource requirements. *arXiv preprint arXiv:1812.07210*,
 2018.
- Enmao Diao, Jie Ding, and Vahid Tarokh. Heterofl: Computation and communication efficient federated learning for heterogeneous clients. *arXiv preprint arXiv:2010.01264*, 2020.
- Samuel Horvath, Stefanos Laskaridis, Mario Almeida, Ilias Leontiadis, Stylianos Venieris, and
 Nicholas Lane. Fjord: Fair and accurate federated learning under heterogeneous targets with
 ordered dropout. *Advances in Neural Information Processing Systems*, 34, 2021.
- Raphael Gontijo Lopes, Stefano Fenu, and Thad Starner. Data-free knowledge distillation for deep
 neural networks. *arXiv preprint arXiv:1710.07535*, 2017.
- Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. Distilling the knowledge in a neural network. *arXiv* preprint arXiv:1503.02531, 2(7), 2015.
- Jianyu Wang, Zachary Charles, Zheng Xu, Gauri Joshi, H Brendan McMahan, Maruan Al-Shedivat,
 Galen Andrew, Salman Avestimehr, Katharine Daly, Deepesh Data, et al. A field guide to federated
 optimization. *arXiv preprint arXiv:2107.06917*, 2021.
- Sohei Itahara, Takayuki Nishio, Yusuke Koda, Masahiro Morikura, and Koji Yamamoto. Distillation based semi-supervised federated learning for communication-efficient collaborative training with
 non-iid private data. *arXiv preprint arXiv:2008.06180*, 2020.
- Yae Jee Cho, Jianyu Wang, Tarun Chiruvolu, and Gauri Joshi. Personalized federated learning for heterogeneous clients with clustered knowledge transfer. *arXiv preprint arXiv:2109.08119*, 2021.
- Seyed Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, Nir Levine, Akihiro Matsukawa, and Hassan
 Ghasemzadeh. Improved knowledge distillation via teacher assistant. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5191–5198, 2020.
- Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. Improving multi-task deep neu ral networks via knowledge distillation for natural language understanding. *arXiv preprint arXiv:1904.09482*, 2019.
- Yuang Jiang, Shiqiang Wang, Victor Valls, Bong Jun Ko, Wei-Han Lee, Kin K Leung, and Leandros
 Tassiulas. Model pruning enables efficient federated learning on edge devices. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- Dhruv Guliani, Lillian Zhou, Changwan Ryu, Tien-Ju Yang, Harry Zhang, Yonghui Xiao, Françoise
 Beaufays, and Giovanni Motta. Enabling on-device training of speech recognition models with
 federated dropout. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8757–8761. IEEE, 2022.
- Jae Hun Ro, Theresa Breiner, Lara McConnaughey, Mingqing Chen, Ananda Theertha Suresh, Shankar Kumar, and Rajiv Mathews. Scaling language model size in cross-device federated learning. *arXiv preprint arXiv:2204.09715*, 2022.
- Nader Bouacida, Jiahui Hou, Hui Zang, and Xin Liu. Adaptive federated dropout: Improving
 communication efficiency and generalization for federated learning. *CoRR*, abs/2011.04050, 2020.
 URL https://arxiv.org/abs/2011.04050.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University
 of Toronto, 2009.

- Barret Zoph, Ekin D Cubuk, Golnaz Ghiasi, Tsung-Yi Lin, Jonathon Shlens, and Quoc V Le. Learning
 data augmentation strategies for object detection. In *European conference on computer vision*,
- ³⁷¹ pages 566–583. Springer, 2020.

372 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*, pages 770–778, 2016.

Mathieu Andreux, Jean Ogier du Terrail, Constance Beguier, and Eric W Tramel. Siloed federated learning for multi-centric histopathology datasets. In *Domain Adaptation and Representation*

Transfer, and Distributed and Collaborative Learning, pages 129–139. Springer, 2020.

U.S. Census Bureau. Percentage distribution of household income in the u.s. in 2020. In Statista,
 September 2021. Retrieved May 18, 2022, from https://www.statista.com/statistics/

September 2021. Retrieved May 18, 2022, from https://www.statista.com
 203183/percentage-distribution-of-household-income-in-the-us

381 Checklist

382	1. For all authors
383 384 385 386 387	 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes] Our work' claims are verified through experiment in Section 4.2 and statistical study in Appendix A.1 (b) Did you describe the limitations of your work? [Yes] The limitations are mentioned in Section 5 (c) Did you discuss any potential negative societal impacts of your work? [No]
389 390	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
391	2. If you are including theoretical results
392 393 394 395	 (a) Did you state the full set of assumptions of all theoretical results? [N/A] We do not show any rigorous theoretical results nor claim as such (b) Did you include complete proofs of all theoretical results? [N/A] We do not show any rigorous theoretical results nor claim as such
396	3. If you ran experiments
397 398 399	 (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] Code is made available to reviewers. Data is open source. Instructions to run are available in code
400 401 402	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] All major training details are mentioned in Section 4 and fine details provided in Appendix A.2
403 404 405	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] All experiments were performed with 5 different seeds and the mean and standard deviations were reported when applicable.
406 407	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
408	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
409 410	(a) If your work uses existing assets, did you cite the creators? [Yes](b) Did you mention the license of the assets? [N/A]
411 412	(c) Did you include any new assets either in the supplemental material or as a URL? $[N/A]$
413 414	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
415 416	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
417	5. If you used crowdsourcing or conducted research with human subjects
418 419	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
420 421	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
422 423	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]