

EMO: Edge Model Overlays to Scale Model Size in Federated Learning

Di Wu, Weibo He, Wanglei Feng, Zhenyu Wen, Bin Qian, and Blesson Varghese

Abstract—Federated Learning (FL) trains machine learning models on edge devices with distributed data. However, the computational and memory limitations of these devices restrict the training of large models using FL. Split Federated Learning (SFL) addresses this challenge by distributing the model across the device and server, but it introduces a tightly coupled data flow, leading to computational bottlenecks and high communication costs. We propose EMO as a solution to enable the training of large models in FL while mitigating the challenges of SFL. EMO introduces Edge Model Overlay(s) between the device and server, enabling the creation of a larger ensemble model without modifying the FL workflow. The key innovation in EMO is Augmented Federated Learning (AFL), which builds an ensemble model by connecting the original (smaller) FL model with model(s) trained in the overlay(s) to facilitate horizontal or vertical scaling. This is accomplished through three key modules: a hierarchical activation replay cache to decouple AFL from FL, a convergence-aware communication controller to optimize communication overhead, and an ensemble inference module. Evaluations on a real-world prototype show that EMO improves accuracy by up to 17.77% compared to FL, and reduces communication costs by up to 7.17 \times and decreases training time by up to 6.9 \times compared to SFL.

I. INTRODUCTION

Machine Learning (ML) applications in language and vision tasks usually use centralized data for training models [10]. However, distributed and private data generated or stored on devices at the edge of the network, such as on mobile phones, is expected to play an increasingly important role in training [4]. Federated Learning (FL) is a proposed privacy-aware solution for training models using distributed data.

However, the *models produced by existing FL solutions are limited in several ways*. Firstly, the size of the model is bound by hard constraints, such as the memory of the device participating in FL [16]. Secondly, the relatively limited computing capability on devices results in impractical training times [19], [21]. For instance, devices used in FL, such as Nvidia Nano, have 8 to 50 times less memory and 10 to 1,000 times fewer FLOPS than resources available in centralized servers, such as Nvidia A100 GPU [14].

Corresponding author: Zhenyu Wen.
Di Wu and Blesson Varghese are with the School of Computer Science, University of St Andrews, UK.

Weibo He and Wanglei Feng are with the Institute of Cyberspace Security and College of Information Engineering, Zhejiang University of Technology, China.

Zhenyu Wen is with the Institute of Cyberspace Security and College of Information Engineering, Zhejiang University of Technology, China; University of Science and Technology of China, China.

Bin Qian is with the State Key Laboratory of Industrial Control Technology, Zhejiang University, China.

Therefore, *the models trained using edge devices in classic FL systems are small and have limited learning performance*.

Split Federated Learning (SFL) has been proposed to mitigate the above shortcoming by offloading certain layers of the model from the device to a server [18]. This reduces the memory use and the computational workload on the device. However, SFL introduces the following two challenges:

Challenge 1 - Low parallel efficiency due to the strong dependency between the computations carried out on the devices and server, referred to as computational locking [16], [18], resulting in a lower system throughput (see Section III-A).

Challenge 2 - High communication costs are introduced due to the frequent exchange of activations and gradients between the devices and the server, thereby increasing latency and bandwidth requirements (see Section III-B).

We propose a novel system, EMO that enables large models to be trained in FL without introducing the abovementioned challenges of SFL. At its core, EMO introduces edge model overlay(s) between the device and server to facilitate the creation of a larger ensemble model that can be scaled both horizontally (model boosting) and vertically (model bagging) [2]. The ensemble model is created by augmenting the original (small) model trained by the FL system with any or all of the larger models trained in the overlays. This is referred to as ‘*Augmented FL*’ (AFL). Unlike SFL, AFL does not modify or replace the FL workflow and the proposed overlays are decoupled from and executed in parallel to the original FL system. To achieve the above, EMO incorporates three modules: the *hierarchical activation replay cache*, the *convergence-aware communication controller* and the *ensemble inference module*. The *hierarchical activation replay cache* enables EMO to decouple AFL from FL. The *convergence-aware communication controller* reduces the AFL communication costs. Finally, the *ensemble inference module* augments the original FL models at the end of training.

This work makes the following contributions:

- (1) A new system, EMO to create larger ensemble models that have up to 17.77% higher accuracy than FL models without increasing computational overheads on the device.
- (2) A novel method, AFL to reduce computational dependencies in SFL accelerating training by up to 6.9 \times and reducing communication by 7.17 \times .

II. RELATED WORK

Training large models in FL remains challenging due to the computational and memory constraints of participating devices. There are two categories of methods presented in the

literature to reduce the computational workload and memory requirements of on-device models in FL. The first is partial training and the second is SFL.

In **partial training** methods in FL, devices extract a sub-model from the global model using various sampling strategies, such as random selection [5], rolling selection [1], or pruning-based selection [9], and train only the sub-model. Despite their computational efficiency, partial training methods introduce further challenges. Since each device trains only a subset of parameters from the global model, each parameter receives fewer updates per FL round, resulting in slower convergence [3]. Moreover, aggregating sub-models extracted from different locations further degrades the model’s accuracy compared to FL training. In some cases, these methods offer no advantage over training smaller models in FL, limiting their effectiveness in scaling FL to larger models [3].

The **SFL** method partitions the on-device model across the device and a server unlike classic FL, where each device trains the entire model locally. Specifically, the earlier layers of the global model are trained on devices, while the later layers are offloaded to a cloud server for training [18]. Despite its advantages, SFL suffers from low parallel efficiency due to computational locking—a strong interdependence between device and server computations. In addition, SFL incurs high communication overhead due to frequent exchanges of activations and gradients, leading to reduced system throughput, increased latency, and higher bandwidth usage.

Unlike the existing methods discussed above, the proposed system **EMO** improves the learning performance of FL by creating a larger ensemble model without additional computational overhead to devices. Furthermore, the training of augmented models used for the ensemble avoids computational locking and incurs lower communication costs than those seen in SFL.

III. COMPUTATIONAL LOCKING AND COMMUNICATION OVERHEAD IN SFL

In SFL, the later layers of the model trained on a device are offloaded to the server to alleviate the computation burden on devices [19]. However, compared to classic FL, SFL systems have strong computational dependencies and experience high communication overhead between the devices and the server.

A. Computational Locking

Model splitting in SFL. In SFL, the full-size model Θ is divided into two parts: a device-side part, $\Theta_{[:j]}$, and a server-side part, $\Theta_{[j:]}$, where j represents the split point. Specifically, layer j acts as the boundary, with $\Theta_{[:j]}$ containing all layers up to and including layer j , and $\Theta_{[j:]}$ comprising all layers following layer j . Therefore, the full model is represented as $\Theta = \{\Theta_{[:j]}, \Theta_{[j:]}\}$. $\Theta_{[:j]}$ is trained on the device, while $\Theta_{[j:]}$ is offloaded to and trained on the server. In contrast, classic FL trains the entire model Θ locally on a device.

Forward and backward locking. The mini-batch gradient descent algorithm is usually the method for updating parameters during model training [6]. Figure 1 illustrates the computational dependencies in SFL and EMO. Given a batch

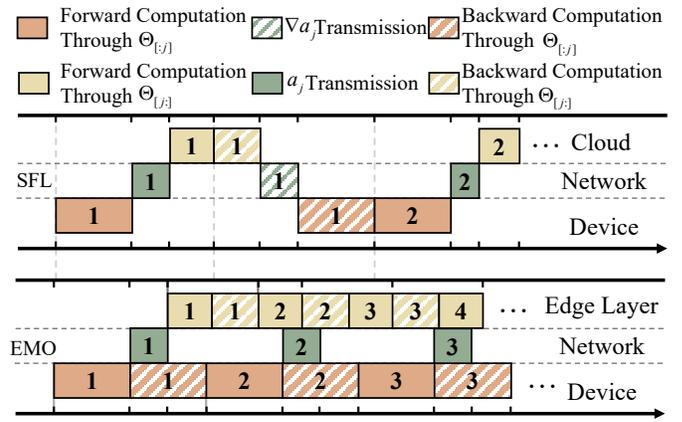


Fig. 1: Computational dependencies in SFL and EMO.

of input data, devices perform forward computation on $\Theta_{[:j]}$. The activations generated by $\Theta_{[:j]}$, denoted as a_j , are then transmitted to the cloud server for forward computation using $\Theta_{[j:]}$. The cloud server must wait for the devices to complete the forward computation using $\Theta_{[:j]}$ before it can execute forward computation using $\Theta_{[j:]}$, causing *forward locking*. Similarly, after completing the forward computation and moving to the backward computation, the server first calculates the gradients $\nabla \Theta_{[j:]}$. The gradient of a_j , denoted as ∇a_j , is then required to be sent back to devices for computing $\nabla \Theta_{[:j]}$. This introduces *backward locking*, as devices must wait to receive ∇a_j before calculating $\nabla \Theta_{[:j]}$. *Forward locking* and *backward locking* reduce parallel efficiency compared to classic FL, thereby decreasing throughput in SFL systems. This motivates the design of the *hierarchical activation replay cache* module in EMO to decouple computations between devices and edge layers. Consequently, EMO eliminates forward and backward locking.

B. Communication Overhead

Another challenge in SFL systems is the communication overhead, which increases latency and raises bandwidth usage.

Communication of activation and gradient in SFL. Compared to classic FL, SFL introduces additional communication costs for transferring the activations (a_j) and corresponding gradients (∇a_j). These costs are substantial and proportional to the combined size of the datasets on devices and the number of local training iterations [17]. The communication costs also depend on the size of the activations and gradients, which may be larger than the original data [20].

Communication latency and bandwidth consumption. Transferring the activation a_j and gradients ∇a_j increases both communication latency and bandwidth consumption. For instance, SFL can introduce up to $800\times$ more cumulative communication time and $610\times$ more bandwidth consumption compared to FL for transferring activations and gradients to the cloud (see Section V-C). This motivates the design of a *convergence-aware communication controller* module in EMO to reduce the communication overhead.

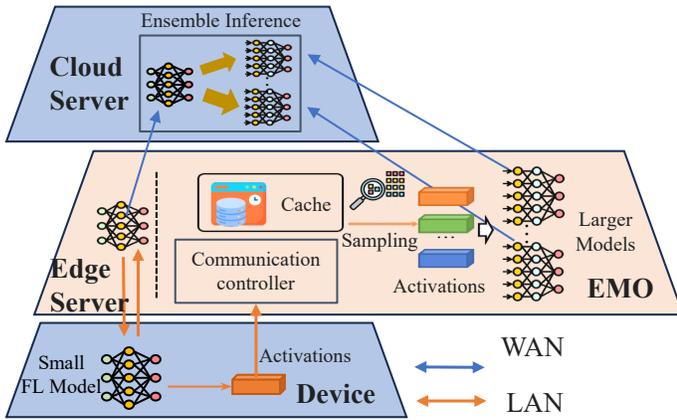


Fig. 2: The EMO architecture.

IV. EMO: EDGE MODEL OVERLAY(S)

A. Overview

The EMO architecture is shown in Figure 4. Augmented FL (AFL) in EMO is hosted alongside devices on edge servers within a local area network (LAN). AFL runs parallel with the FL and trains larger models - the original FL model can be connected to the larger models. The *activation replay cache* module (see Section IV-B) decouples AFL from the original FL system. EMO leverages the available activations stored in the *activation replay cache*, eliminating the need to wait for new device activations (forward locking).

During AFL, activations generated by the FL model must be transferred from devices to edge servers to update the *activation replay cache*. The transmission in EMO occurs within a LAN. Given a LAN’s higher bandwidth than a wide area network (WAN), this approach is more communication-efficient than classic SFL, which transmits activations to a cloud server via a WAN. Additionally, the *communication controller* module in EMO further optimizes activation transfer by adjusting communication intervals (see Section IV-C).

After completing AFL, the original FL model and the larger model(s) trained in EMO are sent to the cloud server for aggregation into an ensemble model. Both horizontal scaling (model boosting) and vertical scaling (model bagging) are employed on the cloud server to create the ensemble model (see Section IV-D). The number of models trained in AFL depends on the number of overlays on the edge servers.

B. Activation Replay Cache

The first module considered is the *activation replay cache*, which decouples AFL from FL. EMO eliminates forward locking seen in SFL by reusing activations received from previous communications. Classic SFL discards activations after a single use. EMO preserves the original FL training workflow, thus avoiding backward locking. Two methods underpin the activation replay cache: hierarchical activation caching and activation replay with cache sampling.

Hierarchical activation caching. The challenge in building an activation cache using previous communications is to store

activations on the edge server efficiently. The size of the activation cache can be substantial, as it depends on the number of activations from all devices, which is determined by both the data size and the number of local training iterations. EMO employs a hierarchical caching architecture to mitigate the potentially large storage requirements. As shown in Figure 3, EMO indexes activations by device ID and batch ID. For each index (a device-batch ID pair), the most recently updated activation is stored in the h5 file format¹. Additionally, EMO builds a second-level activation cache on-disk containing all activations while maintaining a first-level GPU memory cache generated by a sampling policy. This hierarchical caching method allows EMO to store all activations on large disks while minimizing the I/O overhead associated with transferring activations from disk to GPU memory.

Activation replay with cache sampling. The second key consideration in the *activation replay cache* module is to update the first-level GPU cache. In EMO, the update involves sampling activations from the disk cache to the GPU cache. Specifically, after each training iteration using activations stored in the GPU cache, the cache is refreshed by sampling new activations from the disk cache. EMO employs a mixed sampling policy: newly received activations not loaded into the GPU cache are directly selected for training, ensuring the prioritization of fresh data. A random sampling strategy is used if no new activations are available (i.e., all activations have previously been loaded into the GPU cache). This policy prioritizes new activations and switches to uniform sampling once all have been utilized. The random sampling policy can be replaced with importance sampling, where activations are selected based on their associated training loss [11].

C. Convergence-aware Communication Controller

To reduce the communication overhead, EMO incorporates a *convergence-aware communication controller* to monitor activation convergence and dynamically adjust the communication frequency based on the state of convergence of the activations.

Activation convergence. Recent studies have demonstrated that, during model training, the earlier layers converge faster than the later layers — a phenomenon referred to as the bottom-up learning dynamic [15]. This observation provides an opportunity for EMO to reduce activation transfer costs by limiting the transfer of activations from the original FL model as it gradually converges. During the early stages of training, when the distribution of activations undergoes significant changes, EMO collects activations from devices more frequently and updates the disk cache accordingly. As training progresses and the distribution of activations stabilizes, the intervals for activation collection and cache updates are gradually extended. During these intervals, EMO reuses the existing cache, reducing communication overhead. To assess activation convergence, EMO employs Singular Vector Canonical Correlation Analysis (SVCCA) [15]. The SVCCA score is a normalized metric ranging from 0 to 1 to quantify the

¹An h5 file is a file format used to store large amounts of numerical data.

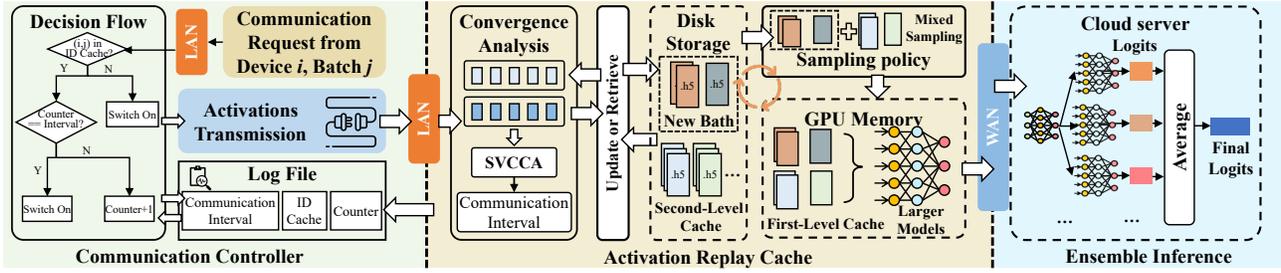


Fig. 3: The activation replay cache, convergence-aware communication controller and ensemble inference modules in EMO.

Algorithm 1: Convergence-aware activation transmission protocol

```

1 Set up: Receive the transfer request for activation
  batch  $(i, j)$  from device  $i$  for batch  $j$ . The current log
  file consists of an ID cache  $\mathbb{B}$ , the communication
  interval for this batch  $i_{(i,j)}$ , and the interval counter
   $c_{(i,j)}$ .
2 if  $(i, j) \notin \mathbb{B}$  then
3   Require the activation batch  $a_{(i,j)}$  to be sent;
4   Initialize  $s_{(i,j)} = 0, c_{(i,j)} = 0$ .
5 end
6 else
7   if  $c_{(i,j)} == i_{(i,j)}$  then
8     Require the activation batch  $a_{(i,j)}$  to be sent.
9     Extract activation  $a_{(i,j)}^c$  from the disk cache;
10    Compute the SVCCA [15] score of  $a_{(i,j)}$  by
11     $s_{(i,j)} \leftarrow SVCCA(a_{(i,j)}^c, a_{(i,j)})$ ;
12    Replace the activation in disk cache
13     $a_{(i,j)}^c \leftarrow a_{(i,j)}$ ;
14    Update the communication interval
15     $i_{(i,j)} \leftarrow \frac{1}{1-s_{(i,j)}}$ ;
16    Reset interval counter  $c_{(i,j)} = 0$ .
17  end
18  else
19    Skip sending this batch activation  $a_{(i,j)}$ .
20     $c_{(i,j)} \leftarrow c_{(i,j)} + 1$ ;
21  end
22 end

```

correlation between two activations; a higher score indicates higher convergence [15].

Convergence-aware activation transmission protocol. Algorithm 1 outlines the activation transmission protocol implemented in the *communication controller*, which dynamically adjusts activation transmission intervals based on the convergence of activations. The decision flow for transferring the activation batch (i, j) is as follows: the controller first checks whether the batch ID is already stored in the ID cache \mathbb{B} . If the batch ID is not cached, the device must transmit this batch. Next, the controller checks the interval between the last activation transmission and the current attempt. The device

sends activations $a_{(i,j)}$ only if the predefined transmission interval has been reached; otherwise, the activations are not sent to EMO, and the communication interval counter is incremented. The communication interval, ID cache, and counter are logged in a log file for querying. Upon receiving the activation, EMO checks whether the batch is already cached. If it is, EMO retrieves the previously stored activation $a_{(i,j)}^c$ and calculates the SVCCA score $s_{(i,j)}$ between the new and cached activations. The cached activation is then updated with the new one, and the communication interval is adjusted using $i_{(i,j)} = \frac{1}{1-s_{(i,j)}}$. This convergence-aware interval adaptation mechanism effectively minimizes redundant communication of similar activations while ensuring that non-trivial activation updates are transmitted.

D. Ensemble inference

After both FL and AFL are completed, the original FL model, denoted as M_{FL} , along with the models trained in the overlay M_{EMO}^i generated on edge server i by EMO, are sent to the cloud for aggregation into an ensemble model. AFL can be executed across N edge servers, resulting in the generation of multiple models within EMO.

Two categories of model aggregation are employed. The first is *horizontal aggregation*, where the original FL model M_{FL} is horizontally connected to each overlay model M_{EMO}^i , also referred to as model boosting [2]. This results in a deeper model (M_{FL}, M_{EMO}^i) , increasing the number of layers in the original FL model M_{FL} . The second is *vertical aggregation*, where the final outputs (logits) from the N pairs of original FL models and their corresponding overlay models (M_{FL}, M_{EMO}^i) are averaged before being used for class prediction. This approach is referred to as vertical aggregation or model bagging [2]. The final ensemble model is deployed in the cloud for inference.

V. EVALUATION

We evaluated EMO against two baselines to assess both the learning and system performance using metrics such as accuracy, communication cost, and training time.

A. Experimental Setup

Testbed. We developed a three-tier system comprising end devices, two edge servers, and a cloud server. The end devices are organized into two clusters: Raspberry Pi 4 Model

TABLE I: Models used for evaluating different methods. The number following each model represents the total number of layers. **C** denotes convolutional layers and **F** represents fully connected layers (both are preceded by the number of layers).

| Methods | CIFAR-10 | | CIFAR-100 | |
|------------|---------------------|--------------------|----------------------------|--------------------------|
| | Device | Server | Device | Server |
| FL-Small | VGG-6 (4C + 2F) | N/A | MobileNet-6 (5C + 1F) | N/A |
| FL-Large | VGG-11 (8C + 3F) | N/A | MobileNet-12 (11C + 1F) | N/A |
| SFL | VGG-4 (4C) | VGG-7 (4C + 3F) | MobileNet-5 (5C) | MobileNet-7 (6C + 1F) |
| EMO | VGG-6 (4C + 2F) | VGG-7 (4C + 3F) | MobileNet-6 (5C + 1F) | MobileNet-7 (6C + 1F) |

TABLE II: Highest test accuracy of EMO and baselines.

| Dataset | Methods | | | |
|-----------|----------|----------|--------|---------------|
| | FL-Small | FL-Large | SFL | EMO |
| CIFAR-10 | 76.57% | 76.22% | 77.13% | 79.51% |
| CIFAR-100 | 24.92% | 33.63% | 41.31% | 42.69% |

B single-board computers and NVIDIA Jetson NX devices equipped with GPUs. The edge servers are powered by NVIDIA RTX 3070 Ti GPUs, while the cloud server operates on an NVIDIA P100 GPU hosted on Alibaba Cloud. All end devices and edge servers are interconnected via a LAN with symmetrical upload/download bandwidths of 800 Mbps. Additionally, both the devices and edge servers are connected to the cloud server over a WAN, with an upload/download bandwidth of 100 Mbps.

Baselines. We evaluated EMO against two baselines: **FL**: The classic FL approach is applied in two configurations: **FL-Small**, where a small model is trained on devices, and **FL-Large**, where a large model is fully trained on devices. In both cases, the entire model resides on devices during training. **SFL**: In the SFL baseline, the top seven layers of the large model are offloaded and trained on a server. We implement two SFL variants: **SFL-Cloud**, where the later layers of the model are offloaded to a cloud server, and **SFL-Edge**, where they are offloaded to an edge server. **EMO**: For EMO, we implement two edge model overlays on separate edge servers, covering 40% and 60% of the devices, respectively.

Training Setup. The training tasks are conducted on the CIFAR-10 [12] and CIFAR-100 [12] datasets, using VGG and MobileNet [7] models, respectively. The datasets are partitioned across devices in a non-independent and non-identically distributed (non-IID) manner using the Dirichlet distribution method [8]. In each training round, 20 devices are randomly selected from a total of 100 devices. Training is carried out over 200 rounds with a learning rate of 0.01. The batch size is set to 16, and the aggregation method used for both FL and SFL is FedAvg [13]. Table I summarizes the model configurations employed for different methods and datasets.

B. Learning Performance

Figure 4 shows the test accuracy curves for the baselines and EMO. On CIFAR-10, FL-Large, SFL, and EMO achieve higher

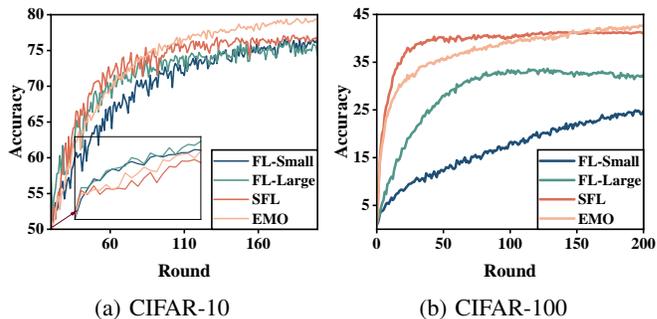


Fig. 4: Test accuracy curves of EMO and baselines for CIFAR-10 and CIFAR-100 datasets.

TABLE III: Average communication cost per round.

| Methods | Communication cost | |
|------------|--------------------|----------------|
| | CIFAR-10 | CIFAR-100 |
| FL-Small | 0.60 GB | 0.004 GB |
| FL-Large | 2.91 GB | 0.12 GB |
| SFL | 0.90 GB | 2.44 GB |
| EMO | 0.83 GB | 0.34 GB |

accuracy in the early stages of training. However, on CIFAR-100, SFL and EMO significantly outperform the FL methods during initial training. In the later stages, EMO consistently surpasses all baselines on both CIFAR-10 and CIFAR-100. Table II reports the highest test accuracy achieved by EMO and the baseline methods. On CIFAR-10, FL-Large achieves accuracy comparable to FL-Small. However, on CIFAR-100, FL-Large outperforms FL-Small by 8.71%, demonstrating the benefits of increased model size in FL training. Across both datasets, EMO outperforms all baselines by a considerable margin. Specifically, it achieves up to 3.29% and 17.77% higher accuracy than the FL methods on CIFAR-10 and CIFAR-100, respectively. Compared to SFL, EMO improves accuracy by 2.38% on CIFAR-10 and 1.38% on CIFAR-100.

C. System Performance

For system performance, we report the average communication cost and training time per round, where a round is defined as a complete cycle of local training on the devices followed by aggregation on the server.

Communication cost. The communication cost of EMO and the four baselines for a single training round are presented in Table III. On CIFAR-10, the size of the VGG model is the main contributor to communication overhead, requiring 0.6 GB for FL-Small and 2.91 GB for FL-Large. SFL introduces an additional 0.3 GB of communication overhead compared to FL-Small, totaling 0.9 GB, while EMO reduces this to 0.83 GB. Using CIFAR-100 and given that the MobileNet model is lightweight, SFL results in the highest communication cost at 2.44 GB due to the communication of activations and gradients. In contrast, EMO achieves a substantial reduction in communication overhead, lowering it to 0.34 GB, which is a 7.17 \times decrease compared to SFL.

Training time. The average training time of a round and communication time for EMO and the baselines are shown in

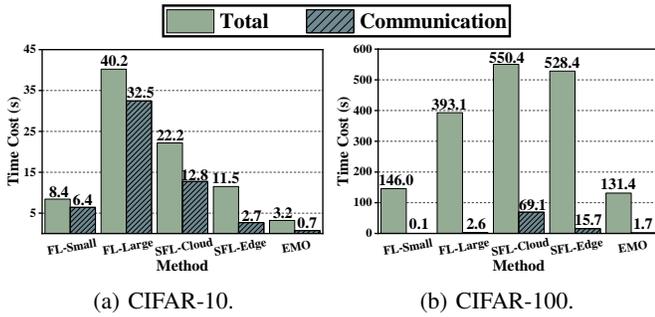


Fig. 5: Average training time per round for different methods, including its communication and total time.

Figure 5. Using both CIFAR-10 and CIFAR-100, EMO achieves the lowest training time. For CIFAR-10, communication time is the dominant factor for FL-Small, FL-Large, and SFL-Cloud. However, SFL-Edge and EMO benefit from higher bandwidth between devices and edge servers, resulting in reduced communication time. On CIFAR-100, communication time is not the bottleneck for FL-Small and FL-Large due to the smaller size of the MobileNet model. Nevertheless, both SFL-Cloud and SFL-Edge introduce additional communication time as a result of transferring activations and gradients to remote servers. Overall, compared to FL and SFL, EMO reduces round training time by up to $12.5\times$ and $6.9\times$, respectively.

VI. CONCLUSION AND FUTURE WORK

We propose EMO, Edge Model Overlay(s), a system that augments FL by scaling the size of trained models while addressing the challenges associated with SFL. Unlike SFL methods, EMO is decoupled from the original FL system and is executed in parallel without modifying the FL workflow. The overlays train larger models on edge servers that can be connected to the smaller FL models trained on devices. Our experiments demonstrate that EMO improves FL training accuracy with augmented models by up to 17.77% compared to FL training with a small model, while also reducing training time. Additionally, EMO decreases communication overhead by up to $7.17\times$ and cuts training time by $6.9\times$ compared to SFL. In future work, we will investigate techniques to enhance privacy of the activations cached in EMO.

ACKNOWLEDGMENT

This work was supported by UK Research and Innovation grant EP/Y028813/1; National Nature Science Foundation of China under Grant 62472387, China; Postdoctoral Science Foundation under Grant 2023M743403 and Zhejiang Provincial Natural Science Foundation of Major Program (Youth Original Project) under Grant LDQ24F020001.

REFERENCES

[1] Samiul Alam, Luyang Liu, Ming Yan, and Mi Zhang. Fedrolex: Model-heterogeneous Federated Learning with Rolling Sub-model Extraction. *Advances in Neural Information Processing Systems*, 35:29677–29690, 2022.

[2] Devansh Arpit, Huan Wang, Yingbo Zhou, and Caiming Xiong. Ensemble of Averages: Improving Model Selection and Boosting Performance in Domain Generalization. *Advances in Neural Information Processing Systems*, 35:8265–8277, 2022.

[3] Gary Cheng, Zachary Charles, Zachary Garrett, and Keith Rush. Does Federated Dropout Actually Work? In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3387–3395, 2022.

[4] Shuiguang Deng, Hailiang Zhao, Weijia Fang, Jianwei Yin, Schahram Dustdar, and Albert Y Zomaya. Edge Intelligence: The Confluence of Edge Computing and Artificial Intelligence. *IEEE Internet of Things Journal*, 7(8):7457–7469, 2020.

[5] Enmao Diao, Jie Ding, and Vahid Tarokh. Heteroft: Computation and Communication Efficient Federated Learning for Heterogeneous Clients. *arXiv preprint arXiv:2010.01264*, 2020.

[6] Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Neural Networks for Machine Learning Lecture 6a Overview of Mini-batch Gradient Descent. *University of Toronto*, 14(8):2, 2012.

[7] Andrew G Howard. Mobilenets: Efficient Convolutional Neural Networks for Mobile Vision Applications. *arXiv preprint arXiv:1704.04861*, 2017.

[8] Tzu-Ming Harry Hsu, Hang Qi, et al. Measuring the Effects of Non-identical Data Distribution for Federated Visual Classification. *arXiv:1909.06335*, 2019.

[9] Zhida Jiang, Yang Xu, Hongli Xu, Zhiyuan Wang, Chunming Qiao, and Yangming Zhao. Fedmp: Federated Learning through Adaptive Model Pruning in Heterogeneous Edge Computing. In *IEEE International Conference on Data Engineering*, pages 767–779. IEEE, 2022.

[10] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling Laws for Neural Language Models. *arXiv preprint arXiv:2001.08361*, 2020.

[11] Angelos Katharopoulos and François Fleuret. Not All Samples are Created Equal: Deep Learning with Importance Sampling. In *International Conference on Machine Learning*, pages 2525–2534. PMLR, 2018.

[12] Alex Krizhevsky, Geoffrey Hinton, et al. Learning Multiple Layers of Features from Tiny Images. 2009.

[13] 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.

[14] Kilian Pfeiffer, Martin Rapp, Ramin Khalili, and Jörg Henkel. Federated Learning for Computationally Constrained Heterogeneous Devices: A Survey. *ACM Computing Surveys*, 55(14s):1–27, 2023.

[15] Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular Vector Canonical Correlation Analysis for Deep Learning Dynamics and Interpretability. *Advances in Neural Information Processing Systems*, 30, 2017.

[16] Dhyanjay Saikumar and Blesson Varghese. NeuroFlux: Memory-Efficient CNN Training Using Adaptive Local Learning. In *Proceedings of the Nineteenth European Conference on Computer Systems*, pages 999–1015, 2024.

[17] Abhishek Singh, Praneeth Vepakomma, Otkrist Gupta, and Ramesh Raskar. Detailed Comparison of Communication Efficiency of Split Learning and Federated Learning. *arXiv preprint arXiv:1909.09145*, 2019.

[18] Chandra Thapa, Pathum Chamikara Mahawaga Arachchige, Seyit Camtepe, and Lichao Sun. Splitfed: When Federated Learning Meets Split Learning. In *AAAI Conference on Artificial Intelligence*, volume 36, pages 8485–8493, 2022.

[19] Di Wu, Rehmat Ullah, Paul Harvey, Peter Kilpatrick, Ivor Spence, and Blesson Varghese. FedAdapt: Adaptive Offloading for IoT Devices in Federated Learning. *IEEE Internet of Things Journal*, 9(21):20889–20901, 2022.

[20] Di Wu, Rehmat Ullah, Philip Rodgers, Peter Kilpatrick, Ivor Spence, and Blesson Varghese. EcoFed: Efficient Communication for DNN Partitioning-Based Federated Learning. *IEEE Transactions on Parallel and Distributed Systems*, 2024.

[21] Tuo Zhang, Lei Gao, Chaoyang He, Mi Zhang, Bhaskar Krishnamachari, and A Salman Avestimehr. Federated Learning for The Internet of Things: Applications, Challenges, and Opportunities. *IEEE Internet of Things Magazine*, 5(1):24–29, 2022.