

RingSFL: An Adaptive Split Federated Learning Towards Taming Client Heterogeneity

Jinglong Shen, *Student Member, IEEE*, Nan Cheng, *Member, IEEE*,

Xiucheng Wang, *Student Member, IEEE*, Feng Lyu, *Member, IEEE*, Wenchao Xu, *Member, IEEE*,

Zhi Liu, *Member, IEEE*, Khalid Aldubaikhy, *Member, IEEE*, and Xuemin (Sherman) Shen, *Fellow, IEEE*

Abstract—Federated learning (FL) has gained increasing attention due to its ability to collaboratively train while protecting client data privacy. However, vanilla FL cannot adapt to client heterogeneity, leading to a degradation in training efficiency due to stragglers, and is still vulnerable to privacy leakage. To address these issues, this paper proposes RingSFL, a novel distributed learning scheme that integrates FL with a model split mechanism to adapt to client heterogeneity while maintaining data privacy. In RingSFL, all clients form a ring topology. For each client, instead of training the model locally, the model is split and trained among all clients along the ring through a pre-defined direction. By properly splitting the model and allocating it to heterogeneous clients, the straggler effect is mitigated, and the training efficiency of the system is significantly enhanced. Additionally, since the model is split, it is less likely for an eavesdropper to obtain the complete model and recover the raw data, thus improving data privacy. The experimental results on both simulation and prototype systems show that RingSFL can achieve better convergence performance than benchmark methods on independently identically distributed (IID) and non-IID datasets, while effectively preventing eavesdroppers from recovering training data.

Index Terms—Federated Learning, Split Learning, Heterogeneity, Ring Topology, Straggler Effect.



1 INTRODUCTION

RECENTLY, machine learning (ML) techniques have been widely applied to various domains such as computer vision [1], natural language processing [2], and speech recognition [3] due to their remarkable representation and learning capabilities [4]. Typically, ML necessitates a large amount of data and computational resources to train a model with satisfactory generalization performance. Consequently, centralized learning has been extensively adopted, where a central server owns all the data and trains a model with abundant computing resources. Nevertheless, in many applications, training data is generated by users, and uploading such raw data to the cloud server may compromise user privacy. Moreover, as the computational and storage capabilities of the user devices increase exponentially, it becomes feasible to leverage the local resources for training tasks. In 2016, Google proposed federated learning (FL) [5], which has been gaining increasing attention. With FL, the ML models are trained over user devices while keeping data localized. Rather than raw data, local updated parameters

are uploaded to the server for aggregation, which prevents data leakage.

Despite the potential of FL in edge networks, there remain numerous challenges [6], [7]. One of the most significant is the high heterogeneity among the clients involved in training [8]. This heterogeneity can manifest in terms of computational capability and battery level, both of which can have a significant impact on the efficiency of the FL system [9]. For instance, the different training times among clients due to uneven computational capabilities can lead to the *straggler effect*, where stragglers can bottleneck the efficiency of the FL system. Additionally, clients may have different battery levels, and the training process can drain the batteries of clients with lower levels, causing them to quit the FL system and making their training data inaccessible.

Another primary concern of FL is user privacy since sensitive information can still be revealed from model parameters/gradients by a third-party entity or the server. In [10], an optimization-based approach is proposed to recover user data from a single sample of gradients. This method involves the generation of random dummy inputs and labels locally, followed by the minimization of the distance between the dummy and actual gradients to recover the user data. Subsequently, [11] extends this approach by using a similar optimization method to achieve recovery of user data from the batch's average gradients. These works highlight the need for further improvements in the data privacy of FL.

In this paper, we propose a novel FL scheme, **RingSFL**, which integrates **FL** with a model **Split** mechanism to address the above issues of FL. Specifically, clients form a ring topology where adjacent clients can communicate

- J. Shen, N. Cheng, and X. Wang are with the School of Telecommunications Engineering, Xidian University, Xi'an, China. E-mail: jishen@stu.xidian.edu.cn, dr.nan.cheng@ieee.org, xcwang_1@stu.xidian.edu.cn
- F. Lyu is with the School of Computer Science and Engineering, Central South University, Changsha, China. E-mail: fenglyu@csu.edu.cn
- W. Xu is with the Department of Computing, Hong Kong Polytechnic University, Hongkong, China. E-mail: wenchao.xu@polyu.edu.hk
- Z. Liu is with the Department of Computer and Network Engineering, The University of Electro-Communications, Chofugaoka, Japan. E-mail: liu@ieee.org
- K. Aldubaikhy is with the Department of Electrical Engineering, Qassim University, Buraydah, Saudi Arabia. E-mail: khalid@qec.edu.sa
- X. Shen is with the Department of Electrical and Computer Engineering, University of Waterloo, Waterloo, Canada. E-mail: sshen@uwaterloo.ca

(Corresponding author: Nan Cheng.)

with each other through wireless links. The server assigns a *propagation length* (the number of neural layers to be processed during each forward and backward propagation) to each client according to the system heterogeneity, and the forward and backward propagation of each client starts from itself and traverses the ring topology under the constraint of propagation length to complete the training. After a round of training, each client has an updated model with different layers trained by data from different clients, i.e., a *blended model*, and the models are uploaded to the server for aggregation. RingSFL not only preserves the capability of FL to utilize the distributed computational power but also enhances data privacy since an eavesdropper can hardly recover the data from the blended model. Moreover, RingSFL can better adapt to the system heterogeneity by allocating computational loads according to the characteristics (such as computational capability or battery power) of different clients, which can significantly mitigate the straggler effect and improve the training efficiency of the system. In summary, the main contributions of this paper are as follows.

- *Novel adaptive and privacy-preserving FL scheme:* A novel FL scheme, termed RingSFL, is proposed. This scheme adaptively distributes the overall training load through a model split mechanism, thereby migrating the computational load from weak computational power clients to strong computational power clients, thus alleviating the straggler effect of vanilla FL, improving the computational efficiency of the system, and reducing the training time. Additionally, RingSFL not only does not transmit user data, but also does not transmit any complete model updates. All model updates uploaded to the server are blended through the model split mechanism, thus achieving enhanced privacy preservation.
- *Better Model Performance:* It has been observed that in RingSFL, the splitting and distribution of the model lead to the emergence of *overlapping layers* that are trained by multiple clients in parallel. These overlapping layers have a higher frequency of gradient aggregation, resulting in more reliable gradients. Consequently, a higher model accuracy can be attained by utilizing a dynamic learning rate to augment the step size of the gradients of the overlapping layers.
- *Simulation and Prototype System:* The efficacy and performance of the proposed RingSFL scheme were evaluated using real datasets in both a simulation environment and a prototype system containing two PC nodes and three Raspberry Pi nodes. The experimental results validate the effectiveness of the proposed RingSFL scheme.

The remainder of the paper is organized as follows. Section 2 provides the background and related work of the paper. Section 3 describes the proposed RingSFL, including the design of the training process, model split scheme, model aggregation scheme, improvement based on overlapping layers, and the discussion on the privacy enhancement. In section 4, extensive experimental results are given to evaluate RingSFL, followed by limitations and future works

discussed in section 5. Finally, concluding remarks are given in section 6.

2 BACKGROUND AND RELATED WORK

2.1 Federated Learning

FL [5] is a distributed machine learning approach that seeks to learn a centralized model of high quality by training on data distributed across a large number of devices, thereby addressing the issue of data silos. The server initiates the process by initializing a global model and sending it to the participating clients. Each client then trains the model using its own local dataset and uploads the updated model parameters to the server. The server aggregates the received model parameters to generate a new global model, which is then sent back to the clients. This process is repeated until the global model converges.

2.1.1 Heterogeneity in FL

The challenge of heterogeneity in FL has been extensively studied. In [12], [13], [14], client heterogeneity is mitigated by designing novel user scheduling mechanisms. Reinforcement learning, long-term perspective, and joint optimization approaches are adopted to select high-quality clients. Furthermore, client heterogeneity is alleviated by an optimization approach to assigning different resources to different clients that considers computation, communication, and energy consumption [15], [16]. Besides client heterogeneity, data heterogeneity is another important issue which refers to Non-IID data distribution among different clients. FedNova is proposed in [17] to address data heterogeneity, where the objective inconsistency due to data heterogeneity is mitigated by normalizing the gradients before aggregation. The problem of data heterogeneity is also mitigated in [18] by judiciously selecting clients to form an approximate IID dataset. In this paper, our proposed RingSFL scheme can alleviate the client heterogeneity issue through model split mechanism. Meanwhile, RingSFL is shown to achieve higher performance than vanilla FL on the Non-IID dataset, which proves the ability of RingSFL to mitigate data heterogeneity.

2.1.2 Privacy Protection in FL

To further improve the privacy of FL, many research works have made essential efforts. In [19], [20], [21], [22], differential privacy mechanisms are proposed to protect user privacy by adding noise to the model. In particular, channel noise is used in uncoded transmission wireless systems to protect user privacy, which reduces the impact of differential privacy on model accuracy [22]. Furthermore, differential privacy is combined with secure multiparty computation to reduce the growth of noise injection as the number of parties increases [23]. Secure aggregation is investigated in [24], [25] using an encryption mechanism that prevents any party, including the server, from accessing each client's model updates. In [26], clients protect their privacy by training and uploading only a portion of the parameters. In [27], a collaborative training model is introduced in which all participants work together to train a federated generative adversarial networks model that can generate artificial data

to replace participants' actual data, thus protecting participants' actual data privacy. In [28], [29], blockchain-based FL schemes are designed to enhance the FL system's security effectively. These approaches develop additional mechanisms to provide privacy at the expense of system efficiency or model performance. Alternatively, the proposed RingSFL ensures enhanced privacy relying only on the learning mechanism itself.

2.1.3 Benefits and Limitations

FL has been demonstrated to achieve accuracy comparable to centralized learning on IID datasets. Furthermore, FL has the advantage of being able to train in parallel, resulting in a relatively low training time consumption. Nevertheless, all training tasks in FL are conducted by clients, which can be challenging to meet in edge network environments due to the need for devices with large storage capacity and strong computational power.

2.2 Split Learning

Split Learning (SL) is proposed in [30] as a means of enhancing the efficiency and security of distributed learning systems, wherein the model is split into two parts and trained by a client and a server, respectively. During the training process, the client first feeds the local training data into the local model and initiates the forward propagation until the output of the split layer is obtained. Subsequently, the client sends the output of the split layer to the server, along with the corresponding labels. The server then inputs the client's output into the local model and continues the forward propagation to generate the predicted model output. Finally, the server utilizes the model output and the labels to calculate the loss value and initiate the back propagation.

In the vanilla SL architecture described above [30], the server does not have direct access to the raw data of clients, nor are the complete model parameters sent to the server. The only information exchanged between clients and the server is the output of the split layer from clients to the server, and the gradient of the split layer from the server to clients. However, this vanilla SL scheme necessitates the transmission of labels, thus risking data leakage. To address this issue, a U-shaped structure is proposed in [31], wherein the deepest and shallowest layers are kept in the client and the middle layers are kept in the server, allowing the labels to remain in the client for training and thus improving system security. Further, random perturbation techniques are proposed in [32] to prevent label information leakage. Additionally, [33] proposes broadcasting an average gradient at the split layer during back propagation, thus enabling scalable parallel SL.

2.2.1 Benefits and Limitations

The server-side model in vanilla SL is more robust and can achieve higher accuracy due to its ability to learn from multiple local datasets. Furthermore, vanilla SL is more suitable for deployment on edge devices than vanilla FL, as the models located on the client side are typically smaller. Nevertheless, vanilla SL often encounters convergence issues on non-IID datasets, and the training latency increases

linearly with the number of clients, making it difficult to train in parallel.

2.3 Integration of Federated Learning and Split Learning

Recently, a novel approach to FL has been proposed, which seeks to integrate FL with SL in order to leverage their respective strengths and create a new architecture for distributed learning that is suitable for edge network environments. SplitFed, as described in [34], is a pioneering and successful example of this approach. The framework necessitates the inclusion of two servers, a main server for splitting training and a fed server for aggregating models. SplitFed splits the neural network between the clients and the main server. The clients then use SL with the main server to update their local model parameters. Subsequently, the clients transmit the updated local model parameters to the fed server for aggregation, resulting in a new global model and the commencement of the next round of training.

In [35], [36], the authors introduced local loss signals into SplitFed, allowing the client's subnetwork to train directly using local losses without receiving loss gradients from the server, reducing the communication cost of the system. In [37], the authors integrate FL with SL in an edge unmanned aerial vehicle (UAV) network. During training, some of the drones in the system perform SL with the base station, while the remaining drones perform FL, and finally all the model parameters are aggregated to get an updated global model. By controlling the training method (FL or SL) of the UAVs, a stronger adaptation to the edge network environment is achieved. FedSL is proposed in [38] for distributed training of recurrent neural networks (RNNs). The authors split an RNN into multiple sub-networks and distribute them to different clients for training. During local training, the sub-networks on different clients communicate with each other to capture potential dependencies between data on different clients. Finally, all clients send their sub-networks to the server for aggregation. In [39], each client has its corresponding edge server, and the clients perform split learning with their edge servers. After local training, the model parameters on the edge server are sent to the parameter server for aggregation. In the field of pre-training large models, integrating FL with SL can also play an active role. FedBert, as proposed in [40], successfully deploys the pre-training task of large models at the edge of the network by splitting the large model into many smaller models. Different from previous discussed works that integrate SL with horizontal FL, PyVertical [41] effectively combines SL with vertical FL, enabling the training of neural networks on vertically split data features between multiple clients while keeping the original data on the owner's device.

2.3.1 Benefits and Limitations

By combining FL with SL, the aforementioned works can effectively retain the benefits of both. Compared to FL, the training load of clients can be effectively reduced; in comparison to SL, it has better parallel training capability and better convergence performance. However, in these existing schemes, all clients are assigned the same training load, and the training load allocation cannot be adjusted according to

TABLE 1
Notations Used in This Paper.

Notation	Explanation
\mathcal{U}	Client set containing clients participating in training.
u_i	Client indexed by i .
N	Number of clients participating in training.
\mathcal{D}_i	Local dataset of u_i .
D_i	Dataset size of \mathcal{D}_i .
C_i	Computational power of u_i .
$\mathcal{W}_{i,(j,k)}^t$	Layer j to k of u_i 's local model at communication round t .
$\tilde{\mathcal{W}}_{i,(j,k)}^t$	Back propagation process of $\mathcal{W}_{i,(j,k)}^t$.
\mathcal{W}^t	Global model at communication round t .
\mathbf{x}_i	Mini-batch sampled from \mathcal{D}_i .
y_i	Corresponding label sampled from \mathcal{D}_i .
L_i	Propagation length.
a_i	Aggregation weight.
$\tilde{\mathbf{x}}_i$	Feature map of \mathbf{x}_i .
$loss_fn$	Loss function.
\mathbf{g}_i^t	Gradient of \mathcal{W}_i^t .
p_i	Model split ratio.
M	Computation volume to complete mini-batch training.
$\mathcal{U}_{i,(j)}$	The set of clients who propagate through the j -th layer of u_i 's local model.
e'_i	The probability that the communication link between u_i and u_{i-1} is eavesdropped.
e_i	The probability that the communication link between u_i and server is eavesdropped.

the client's capability, resulting in an inefficient utilization of computational resources. Furthermore, multiple servers are often required to complete the training, which is difficult to achieve in certain scenarios. In this paper, we effectively realize dynamic allocation of training load by splitting the model between clients. Moreover, our scheme requires only one server to complete the training task.

3 PROPOSED RINGSFL SCHEME

3.1 Overview

The structure of the proposed RingSFL system is depicted in Fig. 1. This system consists of a server for model aggregation and a client set $\mathcal{U} = \{u_0, u_1, \dots, u_{N-1}\}$ of N clients for cooperative training. Each client u_i has a training dataset \mathcal{D}_i of size D_i , and the computational power of u_i is denoted by C_i . The devices form a ring topology, where adjacent devices can communicate with each other through direct communication technologies such as device-to-device (D2D) communication [42]. The devices can also communicate with the server for model downloading and uploading as in FL. It is worth noting that the construction of the ring topology will have a significant effect on the performance. Nevertheless, in this paper, we focus on the

fundamental properties of RingSFL and then arbitrarily set the ring topology, leaving the investigation on ring topology construction as future work. The notations used in this paper are summarized in Table 1.

During each round of training, the forward and backward propagation of a client is conducted cooperatively by all clients along the ring topology. To this end, the server assigns a *propagation length* to each client based on its capabilities, which is used to determine the number of neural layers each client needs to process during each forward and backward propagation. Subsequently, the forward propagation traverses the ring topology, and each client is responsible for propagating the number of neural layers specified by the propagation length over its local model. Similarly, the backward propagation is performed by each client to compute the gradients of the corresponding layers. For example, if there are three clients, denoted by u_0 , u_1 , and u_2 , in RingSFL. Then, the model of u_0 is trained traversing $u_0 \rightarrow u_1 \rightarrow u_2$ for forward propagation and $u_2 \rightarrow u_1 \rightarrow u_0$ for backward propagation. Similarly, the model of u_1 is trained traversing $u_1 \rightarrow u_2 \rightarrow u_0$ for forward propagation and $u_0 \rightarrow u_2 \rightarrow u_1$ for backward propagation, and the case is similar for u_2 . It is worth noting that the training process of different clients are independent and are conducted simultaneously. Finally, the trained local models are transferred to the server for aggregation. This process is repeated until the model converges, as detailed in Algorithm 1.

3.2 RingSFL Training Process

Without loss of generality, we present the detailed training process of RingSFL for an arbitrary number of clients in this subsection. The local model owned by client u_i in communication round t are denoted by \mathcal{W}_i^t with W layers. For a given mini-batch (\mathbf{x}_i, y_i) sampled from \mathcal{D}_i , the output of the model is denoted by $\mathcal{W}_i^t(\mathbf{x}_i) = \mathcal{W}_{i,(0,W-1)}^t(\mathbf{x}_i) = \mathcal{W}_{i,(W-1)}^t(\mathcal{W}_{i,(W-2)}^t(\dots \mathcal{W}_{i,(0)}^t(\mathbf{x}_i)))$, where $\mathcal{W}_{i,(j,k)}^t$ denotes layers j to k of model \mathcal{W}_i^t , $\mathcal{W}_{i,(j)}^t$ denotes layer j of model \mathcal{W}_i^t , \mathbf{x}_i denotes the input samples, and y_i denotes the corresponding labels. The backward propagation can be represented in the similar way: $\tilde{\mathcal{W}}_i^t(grad) = \tilde{\mathcal{W}}_{i,(0,W-1)}^t(grad) = \tilde{\mathcal{W}}_{i,(0)}^t(\tilde{\mathcal{W}}_{i,(1)}^t(\dots \tilde{\mathcal{W}}_{i,(W-1)}^t(grad)))$, where $\tilde{\mathcal{W}}_i^t$ denotes the back propagation process of \mathcal{W}_i^t .

3.2.1 Initialization

In the beginning, each client uploads its state information (C_i, D_i) to the server, where C_i and D_i are the computational power and dataset size of u_i , respectively. The server determines the propagation length $L_i = \frac{C_i}{\sum_{j=0}^{N-1} C_j} W$ and aggregation weight $a_i = \frac{D_i}{\sum_{j=0}^{N-1} D_j}$ for each client based on their state information. Finally, the server dispatches (L_i, a_i) to the corresponding client u_i along with the initialized global model \mathcal{W}^0 . The derivation of the propagation length is described in detail in subsection 3.3. And the aggregation weights are set in the same way as FedAvg [5], as this is not the focus of this paper.

3.2.2 Forward Propagation

In each communication round, each client completes forward and backward propagation in parallel. The forward

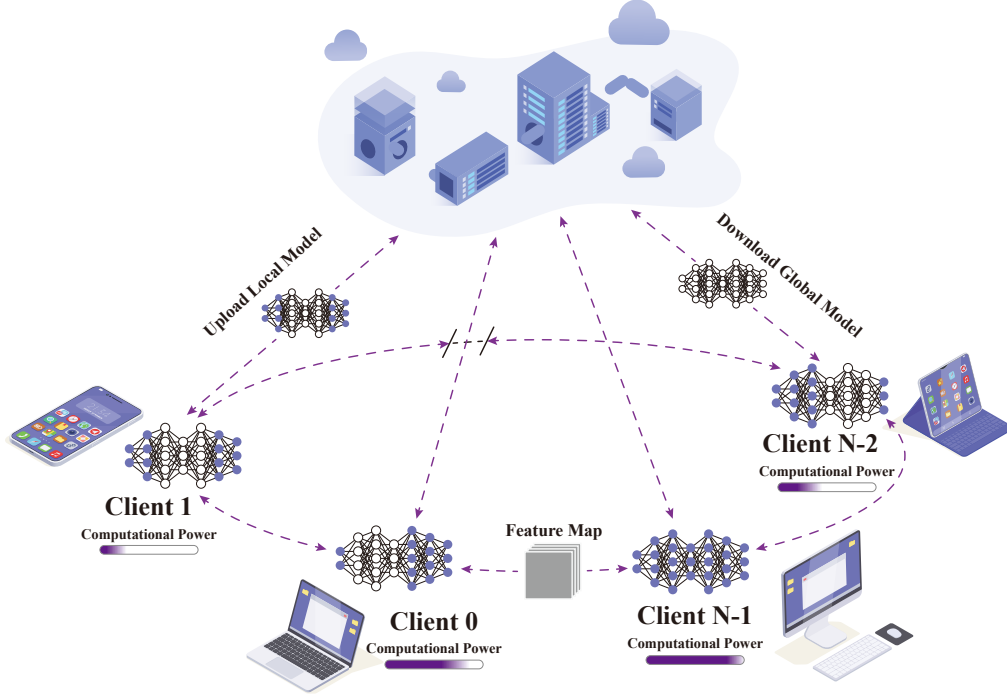


Fig. 1. RingSFL System Structure with N clients.

propagation process of client u_i can be divided into three phases: *Starting Phase*, *Relay Phase*, and *Stop Phase*.

- **Starting Phase:** Client u_i samples a mini-batch of input data (x_i, y_i) from its local dataset \mathcal{D}_i , feeds x_i into local model to get the middle feature map $\bar{x}_i = \mathcal{W}_{i,(0,L_i-1)}^t(x_i)$. Then, u_i sends (\bar{x}_i, l_{stop}) to the next client u_{i+1} in the ring topology, where $l_{stop} = L_i - 1$ denotes the index of the output layer.
- **Relay Phase:** For the next clients in the ring $\{u_{i+1}, \dots, u_{N-1}, u_0, \dots, u_{i-1}\}$, their job is to relay the feature maps from the previous client. Specifically, for each client $u_j \in \{u_{i+1}, \dots, u_{N-1}, u_0, \dots, u_{i-1}\}$, they receive (\bar{x}_i, l_{stop}) from the previous client u_{j-1} , feed it into the local model $\mathcal{W}_{j,(l_{stop}+1, l_{stop}+L_j)}^t$ to obtain a new feature map $\bar{x}_i = \mathcal{W}_{j,(l_{stop}+1, l_{stop}+L_j)}^t(\bar{x}_i)$ and send $(\bar{x}_i, l_{stop} + L_j)$ to the next client u_{j+1} .
- **Stop Phase:** Finally, client u_i will receive the corresponding model output \bar{x}_i from u_{i-1} , and the $loss = loss_fn(\bar{x}_i, y_i)$ is calculated based on it.

3.2.3 Backward Propagation

The backward propagation process is similar to the forward propagation and can be divided into three phases: *Starting Phase*, *Relay Phase*, and *Stop Phase*.

- **Starting Phase:** Client u_i takes the $loss$ calculated locally as the start point of backward propagation, sends $(grad, l_{stop}, a_i)$ to the previous client u_{i-1} along the ring topology, where $grad = loss$, $l_{stop} = W - 1$, and $a_i = \frac{D_i}{\sum_{j=0}^{N-1} D_j}$.
- **Relay Phase:** For the remaining clients in the ring $\{u_{i-1}, \dots, u_0, u_{N-1}, \dots, u_{i+1}\}$, their job is to relay the gradients from the previous client. Specifically,

for each client $u_j \in \{u_{i-1}, \dots, u_0, u_{N-1}, \dots, u_{i+1}\}$, they receive $(grad, l_{stop}, a_i)$ from u_{j+1} , back propagate L_j layers in $\mathcal{W}_{j,(l_{stop}-L_j+1, l_{stop})}^t$ to get the corresponding gradient $g_{j,(l_{stop}-L_j+1, l_{stop})}^t$, which will be weighted by a_i and cached. Then, u_j sends $(g_{j,(l_{stop}-L_j+1, l_{stop})}^t, l_{stop} - L_j, a_i)$ to u_{j-1} , where $g_{j,(l_{stop}-L_j+1, l_{stop})}^t = \mathcal{W}_{j,(l_{stop}-L_j+1, l_{stop})}^t(grad)$.

- **Stop Phase:** Finally, the backward propagation will stop at client u_i , and the gradients are weighted by its aggregation weight a_i and cached in all clients separately. Based on the gradients cached locally, each client updates their local model.

3.2.4 Model Aggregation

In each communication round, the trained local model parameters \mathcal{W}_i^{t+1} are uploaded to the server for aggregation. Since the gradients are already weighted during the training process, model aggregation can be achieved by direct averaging: $\mathcal{W}^{t+1} = \frac{1}{N} \sum_{i=0}^{N-1} \mathcal{W}_i^{t+1}$, which will be detailed in subsection 3.4. RingSFL performs consecutive training until the model converges. It can be seen from the forward/backward propagation that the local model of a client is trained by all clients, which leads to a blended model. This poses increasing difficulty in training data reconstruction from the eavesdropped parameters.

3.3 Model Split Scheme

In order to minimize the training time in RingSFL, model splitting is implemented by assigning distinct propagation lengths L_i to each client. This subsection details the procedure for determining L_i .

To minimize the training time, the training load of each client should match its computational power. We denote

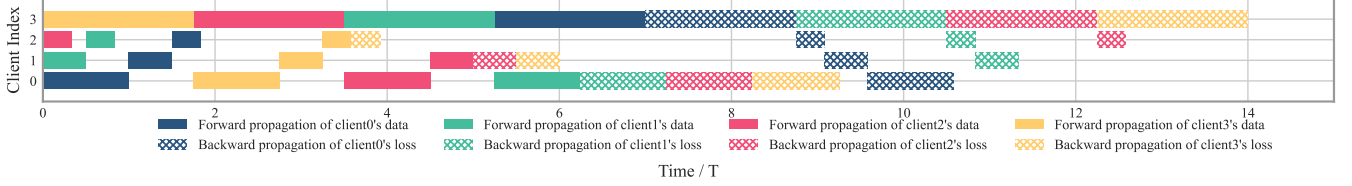


Fig. 2. Propagation Flow of RingSFL with $\frac{C_0}{C} = 0.1, \frac{C_1}{C} = 0.2, \frac{C_2}{C} = 0.3, \frac{C_3}{C} = 0.4, p_0 = 0.1, p_1 = 0.1, p_2 = 0.1, p_3 = 0.7$.

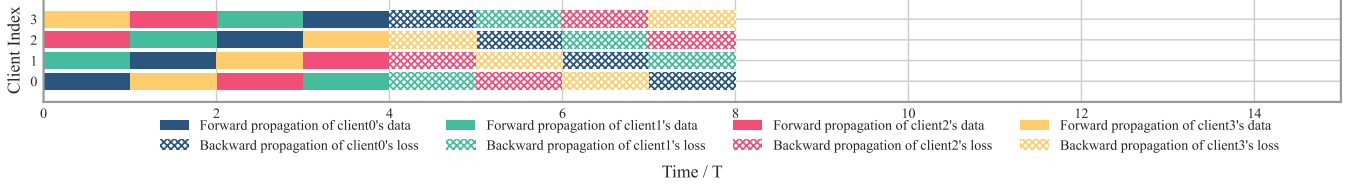


Fig. 3. Propagation Flow of RingSFL with $\frac{C_0}{C} = p_0 = 0.1, \frac{C_1}{C} = p_1 = 0.2, \frac{C_2}{C} = p_2 = 0.3, \frac{C_3}{C} = p_3 = 0.4$.

the computation required to complete training for a mini-batch (including forward propagation, backward propagation and parameter updates) by M GFLOPs, then the total computation required to complete a mini-batch of training for all users in the system is MN . And the computation undertaken by u_i can be denoted by $p_i MN$, where p_i denotes the ratio of the training load assigned to u_i to the total training load, $\sum_{i=0}^{N-1} p_i = 1$. Then, the computation time consumed by u_i to complete the training of a mini-batch is $\frac{p_i MN}{C_i}$.

Since there are N clients in the system, the computation time consumed by the straggler to complete the training of a mini-batch is $\max \left\{ \frac{p_0 MN}{C_0}, \frac{p_1 MN}{C_1}, \dots, \frac{p_{N-1} MN}{C_{N-1}} \right\}$. To minimize the time consumption of the straggler, we formulate the following optimization problem.

$$\min_{p_0, \dots, p_{N-1}} \max \left\{ \frac{p_0 MN}{C_0}, \frac{p_1 MN}{C_1}, \dots, \frac{p_{N-1} MN}{C_{N-1}} \right\} \quad (1)$$

$$\text{s.t.} \quad \sum_{i=0}^{N-1} p_i = 1, \quad (1a)$$

$$0 \leq p_i \leq 1, \quad \forall i = 0, \dots, N-1. \quad (1b)$$

Problem (1) is a MinMax problem. By introducing a new variable m , it can be rewritten as

$$\min_{p_0, \dots, p_{N-1}, m} m \quad (2)$$

$$\text{s.t.} \quad \frac{p_i MN}{C_i} - m \leq 0, \quad \forall i = 0, \dots, N-1, \quad (2a)$$

$$\sum_{i=0}^{N-1} p_i = 1, \quad (2b)$$

$$0 \leq p_i \leq 1, \quad \forall i = 0, \dots, N-1. \quad (2c)$$

Now, the original problem has been transformed into an ordinary linear programming problem, and it is easy to find

the optimal solution as

$$\begin{cases} p_i^* = \frac{C_i}{\sum_{j=0}^{N-1} C_j}, & \forall i = 0, \dots, N-1, \\ m^* = \frac{MN}{\sum_{j=0}^{N-1} C_j}. \end{cases} \quad (3)$$

The solution reveals that the optimal p_i^* should be equal to the ratio of the computational power $\frac{C_i}{\sum_{j=0}^{N-1} C_j}$. Since the training load of each client can be controlled by the number of locally trained neural layers, we set the propagation length to $L_i = p_i^* W = \frac{C_i}{\sum_{j=0}^{N-1} C_j} W$. However, in practice, the W layers of the model are of different sizes, and $p_i^* W$ is generally not an integer. Therefore, the server should adjust the propagation length L_i accordingly. This will affect the total training time, and we leave the details for future work.

To clearly show the relationship between the total training time and the training load of each client, Fig. 2 illustrates the computation flow of each client in a 4-client RingSFL system with $\frac{C_0}{C} = 0.1, \frac{C_1}{C} = 0.2, \frac{C_2}{C} = 0.3, \frac{C_3}{C} = 0.4$ and $p_0 = 0.1, p_1 = 0.1, p_2 = 0.1, p_3 = 0.7$, where $C = \sum_{i=0}^{N-1} C_i$ denotes the total computational power of all clients in the system. The unit of time axis in the figure is set to $T = \frac{M}{2C}$. The figure shows that although client 3 has more computational power than the other clients, it becomes a straggler in training due to the excessive training load. As a result, the time cost of the system to complete a mini-batch training is bottlenecked by this client, who spends the longest training time. Therefore, the computational load should be optimally allocated to suppress the straggler effect so that the training time is minimized.

Fig. 3 shows the computation flow of each client in a 4-client RingSFL system with $p_0 = \frac{C_0}{C} = 0.1, p_1 = \frac{C_1}{C} = 0.2, p_2 = \frac{C_2}{C} = 0.3, p_3 = \frac{C_3}{C} = 0.4$. From the figure, it can be seen that the straggler effect is significantly mitigated by optimally allocating p_i . The computation time consumed by each client is equal, and the time consumed by the system to complete a mini-batch training is reduced to 8 units of time.

3.4 Model Aggregation Scheme

The server receives the local models uploaded by clients and aggregates them, necessitating the design of a revised model aggregation scheme in RingSFL due to the blended model, which renders the traditional FedAvg algorithm [5] inapplicable.

Different from FedAvg, the weighting in RingSFL is achieved by each client during the training process. The aggregation weight a_i is transferred among users along with backward propagation, and the computed gradients are weighted by it. Since the layers of each client's local model may lie within the propagation range of multiple clients, the gradients computed using data from different clients may accumulate on the same layer. Denote the set of clients who propagate through the j -th layer of u_i 's local model by $\mathcal{U}_{i,(j)}$, where $\bigcup_{i=0}^{N-1} \mathcal{U}_{i,(j)} = \{0, \dots, N-1\}$. Then, in the t -th communication round, u_i 's training result is

$$\mathcal{W}_i^{t+1} = \begin{bmatrix} \mathcal{W}_{i,(0)}^t - \sum_{j \in \mathcal{U}_{i,(0)}} \eta a_j \mathbf{g}_{j,(0)}^t \\ \vdots \\ \mathcal{W}_{i,(k)}^t - \sum_{j \in \mathcal{U}_{i,(k)}} \eta a_j \mathbf{g}_{j,(k)}^t \\ \vdots \\ \mathcal{W}_{i,(W-1)}^t - \sum_{j \in \mathcal{U}_{i,(W-1)}} \eta a_j \mathbf{g}_{j,(W-1)}^t \end{bmatrix} \quad (4)$$

where $\mathbf{g}_{j,(k)}^t$ is the gradient of the k -th layer calculated based on the data from u_j 's dataset in communication round t , a_j is aggregation weight of u_j , and η is learning rate. Server can aggregate by simply averaging after receiving \mathcal{W}_i^{t+1} , $\forall i = 0, \dots, N-1$. The new global model obtained by server aggregation is

$$\begin{aligned} \mathcal{W}^{t+1} &= \frac{1}{N} \sum_{i=0}^{N-1} \mathcal{W}_i^{t+1} \\ &= \begin{bmatrix} \mathcal{W}_{i,(0)}^t - \frac{\eta}{N} \sum_{i=0}^{N-1} \sum_{j \in \mathcal{U}_{i,(0)}} a_j \mathbf{g}_{j,(0)}^t \\ \vdots \\ \mathcal{W}_{i,(k)}^t - \frac{\eta}{N} \sum_{i=0}^{N-1} \sum_{j \in \mathcal{U}_{i,(k)}} a_j \mathbf{g}_{j,(k)}^t \\ \vdots \\ \mathcal{W}_{i,(W-1)}^t - \frac{\eta}{N} \sum_{i=0}^{N-1} \sum_{j \in \mathcal{U}_{i,(W-1)}} a_j \mathbf{g}_{j,(W-1)}^t \end{bmatrix} \end{aligned} \quad (5)$$

Since $\bigcup_{i=0}^{N-1} \mathcal{U}_{i,(k)} = \{0, \dots, N-1\}$, we can derive that $\sum_{i=0}^{N-1} \sum_{j \in \mathcal{U}_{i,(k)}} a_j \mathbf{g}_{j,(k)}^t = \sum_{i=0}^{N-1} a_i \mathbf{g}_{i,(k)}^t$. Then, (5) can be rewritten as

$$\begin{aligned} \mathcal{W}^{t+1} &= \begin{bmatrix} \mathcal{W}_{i,(0)}^t - \frac{\eta}{N} \sum_{i=0}^{N-1} a_i \mathbf{g}_{i,(0)}^t \\ \vdots \\ \mathcal{W}_{i,(k)}^t - \frac{\eta}{N} \sum_{i=0}^{N-1} a_i \mathbf{g}_{i,(k)}^t \\ \vdots \\ \mathcal{W}_{i,(W-1)}^t - \frac{\eta}{N} \sum_{i=0}^{N-1} a_i \mathbf{g}_{i,(W-1)}^t \end{bmatrix} \\ &= \mathcal{W}^t - \frac{\eta}{N} \sum_{i=0}^{N-1} a_i \begin{bmatrix} \mathbf{g}_{i,(0)}^t \\ \vdots \\ \mathbf{g}_{i,(k)}^t \\ \vdots \\ \mathbf{g}_{i,(W-1)}^t \end{bmatrix} \\ &= \mathcal{W}^t - \frac{\eta}{N} \sum_{i=0}^{N-1} a_i \mathbf{g}_i^t \end{aligned} \quad (6)$$

The formulation of the aggregation result in (6) is similar to FedAvg [5], except that the learning rate is reduced by N times. To compensate for the discounted learning rate, we can manually multiply the learning rate by the number of clients participating in the training.

3.5 Overlapping Layers Can Improve Model Performance

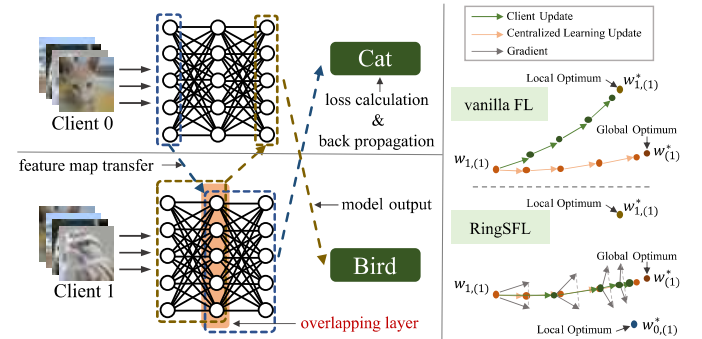


Fig. 4. RingSFL training over 2 clients; $L_0:L_1=1:2$.

As previously discussed in the subsections, due to the varying propagation lengths of each client, there are *overlapping layers* that are trained by multiple clients simultaneously. Fig. 4 illustrates a RingSFL system with two clients that jointly train a multilayer perceptron (MLP) with three layers. The local models of each client are denoted by \mathcal{W}_0^t and \mathcal{W}_1^t , with propagation lengths set to $L_0 = 1$ and $L_1 = 2$, respectively. Since layer $\mathcal{W}_{1,(1)}^t$ is within the propagation range of both client 0 and client 1, it is an overlapping layer. In each local iteration, this layer accumulates gradients from both two clients (the gradients of this layer are aggregated at each local iteration). This implies that the gradients of overlapping layers are more reliable and more capable of guiding the model parameters towards the global optimum. In contrast, in vanilla FL, the local model is updated using the gradients computed by itself only, which will lead the local model to converge towards the local optimum.

To better leverage the reliability of overlapping layer gradients, we propose RingSFLv2¹. We introduce a dynamic learning rate for overlapping layers, which is adjusted according to the number of gradients accumulated on it. Specifically, the larger the number of accumulated gradients, the larger the learning rate, thus enhancing the influence of accumulated gradients on the parameters of each overlapping layer. For example, let the j -th layer of client u_i be an overlapping layer. During training, the gradients generated by clients in the client set $\mathcal{U}_{i,(j)}$ are accumulated on this layer, and we denote the number of clients in $\mathcal{U}_{i,(j)}$ by $\|\mathcal{U}_{i,(j)}\|$. Then, the parameter update of the overlapping layer can be formulated as

$$\mathcal{W}_{i,(j)}^t = \mathcal{W}_{i,(j)}^t - \eta \|\mathcal{U}_{i,(j)}\| \sum_{k \in \mathcal{U}_{i,(j)}} a_k \mathbf{g}_{k,(j)}^t. \quad (7)$$

By altering the manner in which the overlapping layers are updated, we observe a marked improvement in the model accuracy, as discussed in greater detail in the experiment section.

3.6 Privacy Enhancement

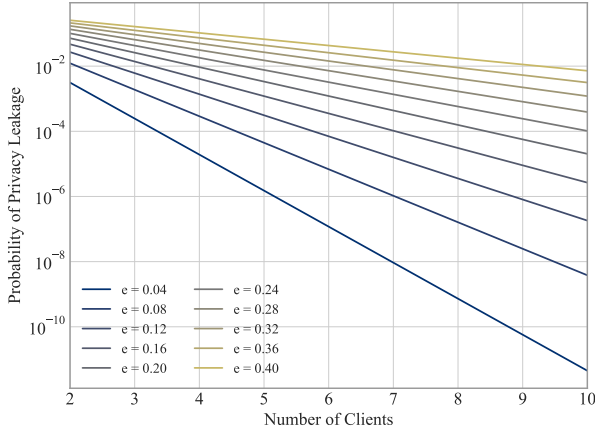


Fig. 5. Impact of Number of Clients and Eavesdropping on Privacy Leakage.

While a theoretical proof of privacy is beyond the scope of this paper, this subsection will provide a brief overview of the privacy-enhancing features of RingSFL.

3.6.1 For Eavesdropper

According to existing research [10], [11], the eavesdropper must obtain the complete model parameters or gradients to recover client data. Recovering client data from any partial or fractured model is complex, thus the enhanced privacy feature of RingSFL is proposed. In RingSFL, the communications among clients only contain the feature map of each client, while the models transferred from the client to the server are blended models. As the eavesdropper is not aware of the propagation lengths, it is difficult to obtain the complete model parameters by eavesdropping. However, when only one layer is trained per client, it is possible to

recover client data by eavesdropping. By eavesdropping on all the communication links in the system, the eavesdropper can obtain the gradients of each layer of the model and piece together a complete model. We then analyze the probability of privacy leakage in this case. Let e'_i denote the probability that the communication link between u_i and u_{i-1} is eavesdropped, and e_i denote the probability that the communication link between u_i and server is eavesdropped. The privacy leakage probability of client u_i can then be expressed as

$$P_i = e_i \prod_{j=0, \dots, N-1; j \neq i} (1 - (1 - e_j)(1 - e'_j)). \quad (8)$$

It is evident from (8) that the privacy leakage probability is influenced by the eavesdropping probabilities e_i, e'_i and the number of clients N . Fig. 5 illustrates the impact of different eavesdropping probabilities and the number of clients on the privacy leakage probability, where $e = e_i = e'_i$. It is observed that the probability of privacy leakage decreases exponentially with the increase in the number of clients. Even with a high probability of eavesdropping on each link (e.g., $e = 0.4$), the leakage probability still rapidly decreases to a near-zero value. This implies that when the number of clients in the RingSFL system is sufficiently large, a high level of security can be maintained even in extreme cases (each client trains only one layer at a time and has a high eavesdropping probability).

3.6.2 For Malicious Server

Given that the server is aware of the propagation lengths of each client, the process of recovering client data is slightly less complex for the server than for an eavesdropper. However, due to the fact that clients with a larger propagation length typically have overlapping layers during training, and the recovery of individual client gradients from these overlapping layers is difficult, RingSFL still provides enhanced privacy. Nevertheless, in certain unique propagation length settings, there are no overlapping layers, and it is possible for the server to recover data from all uploaded models. To address this issue, clients can negotiate with the server to ensure that the system has overlapping layers, thereby guaranteeing security.

4 EXPERIMENTAL RESULT

In this section, we assess the efficacy of the proposed scheme and analyze the performance improvements. We commence by introducing the experimental configuration and then present the experimental results.

4.1 Experimental Setup

4.1.1 Platforms

Our proposed RingSFL scheme was evaluated in both a simulated environment and a real prototype system.

- **Simulation Environment:** The software used in the simulation was written in Python 3.9.12, and PyTorch 1.11.0 was employed for model construction and training to ensure compatibility with the software environment of the prototype system. In

1. RingSFL without dynamic learning rate is referred to as RingSFLv1

Algorithm 1: RingSFL Training Progress

```

1 Server initializes model parameters  $\mathcal{W}^0$  with  $W$ 
  layers;
2 Clients upload state information ( $D_i$ : dataset size,
   $C_i$ : computational power) to the server;
3 Server calculates propagation length and
  aggregation weight for clients:
   $L_i = p_i W = \frac{C_i}{\sum_{j=0}^{N-1} C_j} W$ ;  $a_i = \frac{D_i}{\sum_{j=1}^N D_j}$ ;
4 Server sends  $(\mathcal{W}^0, L_i, a_i)$  to each client  $u_i$ ;
5 for each round  $t$  from 0 to  $T - 1$  do
6   for each client  $u_i$  parallelly do
7     Set local model  $\mathcal{W}_i^t = \mathcal{W}^t$ ;
8     for each epoch from 0 to  $E - 1$  do
9       for each batch  $(x_i, y_i)$  in  $\mathcal{D}_i$  do
10        BatchUpdate( $\mathcal{W}_i^t, (x_i, y_i)$ );
11      end
12    end
13    Upload training results  $\mathcal{W}_i^{t+1}$ ;
14  end
15  Server aggregates parameters:
   $\mathcal{W}^{t+1} = \frac{1}{N} \sum_{i=0}^{N-1} \mathcal{W}_i^{t+1}$ ;
16  Server sends  $\mathcal{W}^{t+1}$  to each client;
17 end

```

Algorithm 2: BatchUpdate($\mathcal{W}_i^t, (x_i, y_i)$)

```

1  $\bar{x}_i = x_i, l_{stop} = -1$ ;
2 for  $u_j$  in  $[u_i, \dots, u_{N-1}, u_0, \dots, u_{i-1}]$  do
3    $\bar{x}_i = \mathcal{W}_{j, (l_{stop}+1, l_{stop}+L_j)}^t(\bar{x}_i)$ ;
4    $l_{stop} = l_{stop} + L_j$ ;
5   Send  $(\bar{x}_i, l_{stop})$  to  $u_{j+1}$ ;
6 end
7  $loss = loss\_fn(\bar{x}_i, y_i)$ ;
8  $grad = loss, l_{stop} = W - 1$ ;
9 for  $u_j$  in  $[u_{i-1}, \dots, u_0, u_{N-1}, \dots, u_i]$  do
10   $grad = \mathcal{W}_{j, (l_{stop}-L_j+1, l_{stop})}^t(grad)$ ;
11  Cache gradients  $a_i \cdot g_{j, (l_{stop}-L_j+1, l_{stop})}^t$ 
    calculated in previous step;
12   $l_{stop} = l_{stop} - L_j$ ;
13  Send  $(grad, l_{stop}, a_i)$  to  $u_{j-1}$ ;
14 end
15 Each client  $u_i$  updates their local model  $\mathcal{W}_i^t$  based
    on the gradients cached;

```

the simulation, a ring topology was formed by 5 clients, who then conducted RingSFL with a server. The model performance of RingSFL was evaluated using different models and both IID and Non-IID datasets.

- **Prototype System:** As depicted in Fig. 6, our prototype system consists of three Raspberry Pi nodes (ARM Cortex-A72 @ 1.5GHz 6.4W) as weak computational power clients and two PC nodes (11th Gen Intel(R) Core(TM) i7-11700 @ 2.50GHz 65W) as strong computational power clients. A laptop is utilized as a parameter server for model aggregation. Each device operates at the 5440 MHz band with a

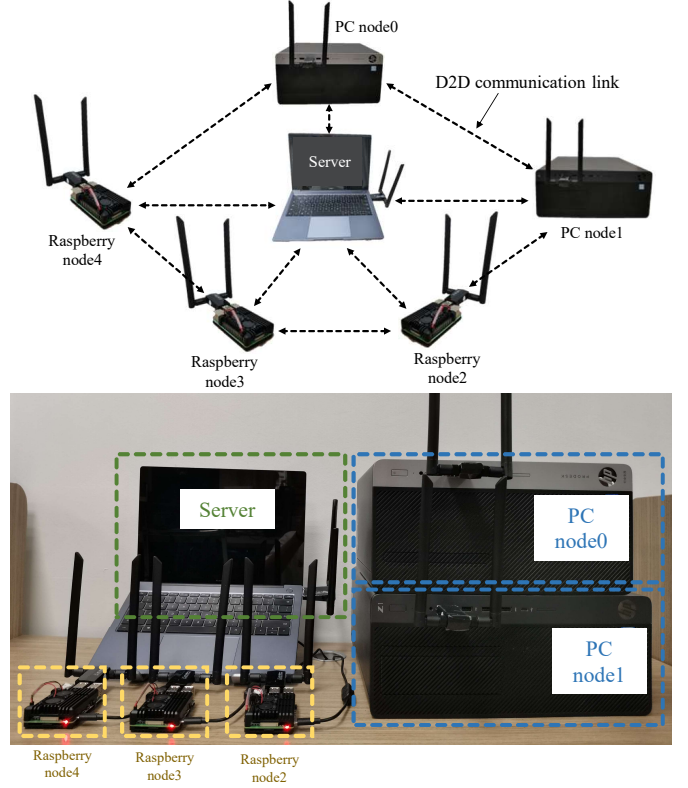


Fig. 6. Prototype System of RingSFL with 2 PC nodes and 3 raspberry Pi nodes.

channel of 40 MHz bandwidth, allowing D2D TCP rates of 135 Mbps with a standard deviation of 5.83. Through the prototype system, we mainly assess the capability of RingSFL to adapt to the client heterogeneity in real-world settings in terms of training time and energy consumption. As it is difficult to measure the energy consumption in the prototype system, we instead capture each client's actual CPU occupation time and convert it into consumed energy by multiplying it by the single-core full power. Additionally, we evaluate the effect of D2D communications on the proposed RingSFL scheme.

4.1.2 Datasets and Models

- **Datasets:** The experiments utilize the widely-used MNIST [43] and CIFAR10 [44] datasets. Specifically, MNIST consists of 70,000 square ($28 \times 28 = 784$ pixel) grayscale handwritten digital images, divided into 10 categories (60,000 for training and 10,000 for testing); Cifar10 consists of 60,000 square ($32 \times 32 = 1024$ pixel) 3-channel color images, divided into 10 categories (50,000 for training and 10,000 for testing). Each client's local dataset has the same number of samples. Furthermore, to evaluate the effect of the dataset distribution on RingSFL, experiments are conducted using both IID and non-IID data. Specifically, the Non-IID dataset is constructed by assigning each client a sub-dataset containing only two categories.
- **Models:** To validate the performance of our scheme, we employ mainstream networks for evaluation, in-

cluding ResNet18 [45], VGG16 [46], AlexNet [47], and LeNet-5 [43]. To ensure the model converges to the desired accuracy, we initialize the model parameters prior to training. Specifically, the Conv2d layer is initialized using the Kaiming normal method. The weights of the BatchNorm2d layer are initialized to 1, and the bias is initialized to 0. The weights of the Linear layer are initialized using a Gaussian distribution with a mean of 0 and standard deviation of 0.01, and the bias is initialized to 0.

4.1.3 Implementation

In our experiments, we use vanilla FL [5], vanilla SL [30], and SplitFed [34] as benchmarks. For each experiment, five clients (2 PC nodes u_0, u_1 and 3 raspberry Pi nodes u_2, u_3, u_4) are included for 100 training rounds. In each round, each client iterates two epochs on their local dataset. Since vanilla SL has no parameter aggregation process, the whole training process cannot be directly divided into rounds. Therefore, in our experiments, when vanilla SL updates $local_epoch_num^2 \times batch_num^3$ steps, we treat it as a round for a fair comparison with other algorithms. SGD optimizer is used for stochastic gradient descent to optimize the model. For the convenience of model split, we pre-divide models into blocks. Specifically, ResNet18 is divided into 10 blocks, VGG16 is divided into 16 blocks, AlexNet is divided into 17 blocks, and LeNet-5 is divided into 12 blocks. The complete training parameters are presented in Table 2.

4.2 Results and Discussion

We first evaluate the convergence performance of RingSFL and the impact of overlapping layers on model accuracy in the simulation environment. Subsequently, the impact of the D2D communication rate on RingSFL and the efficiency of RingSFL in terms of training time and energy consumption are evaluated in the prototype system. Finally, we validate the privacy enhancement achieved by RingSFL. Due to space limitations, we mainly discuss the performance of RingSFL with ResNet18 as an example. Please refer to the tables for the experimental results of other models.

4.2.1 Convergence Performance

This experiment serves to validate the convergence performance of RingSFL, with SplitFed, vanilla FL, and vanilla SL used as benchmarks for comparison. The propagation lengths of RingSFL were set to $L_0:L_1:L_2:L_3:L_4=6:1:1:1:1$.

As illustrated in Fig. 7, RingSFLv2 achieved the highest model accuracy on the CIFAR10 (IID) dataset after the same communication rounds, compared to the considered benchmarks. RingSFLv1 performed slightly worse than RingSFLv2, yet still converged to a higher model accuracy than vanilla FL and SplitFed. This is attributed to the fact that in RingSFL, gradients accumulated on overlapping layers are aggregated in each batch of training, thus leading to a higher capability to optimize the parameters of overlapping layers towards the global optimum, and consequently, better model performance. Furthermore, Fig. 8 demonstrates the

effect of different data distributions on the convergence performance of RingSFL when training on the CIFAR10 (Non-IID) dataset. It can be observed that RingSFLv2 still maintained the best model accuracy on the Non-IID dataset compared to the considered comparison algorithms. RingSFLv1 performed similarly to vanilla FL, while vanilla SL and SplitFed had difficulty in converging on the Non-IID dataset. This indicates that RingSFL is able to adapt well to the non-IID data distribution.

The results demonstrate that RingSFL has the most optimal convergence performance in comparison to the benchmarks evaluated. Notably, the distinction in performance between RingSFLv2 and RingSFLv1 further accentuates the significance of the overlapping layers in RingSFL, which will be discussed in greater detail in the subsequent subsection. The training results of the other models are presented in Table 3.

4.2.2 Effect of Overlapping Layers

In order to evaluate the impact of overlapping layers on model accuracy, we trained models with different propagation length configurations in RingSFLv2 and gradually made the propagation length distribution uneven to increase the number of overlapping layers. Specifically, we set the propagation length to $L_0:L_1:L_2:L_3:L_4=\{2:2:2:2:2, 3:2:2:2:1, 3:3:2:1:1, 4:3:1:1:1, 5:2:1:1:1, 6:1:1:1:1\}$ in each experiment to observe the effect of increasing the number of overlapping layers on model accuracy.

As demonstrated in Fig. 9, the effect of overlapping layers on the model accuracy is evident. The more unevenly the propagation lengths are allocated, the more overlapping layers there are, and the higher the model accuracy is. This phenomenon can be interpreted as *the gradient reliability of overlapping layer is higher*. As discussed in [17], due to the distributed characteristic, the optimal solution of the local objective function of each client is not consistent with the optimal solution of the global objective function. Thus, the gradients computed by each client in the process of local training will not optimize the local model towards the global optimum. In the RingSFL training process, the overlapping layers accumulate the gradients from different clients, which makes the gradients on the overlapping layer more reliable and has more tendency to optimize towards the global optimum. The training results of other models are presented in Table 4. This conclusion was further validated on the Non-IID dataset, as illustrated in Fig. 10. It can be seen that this conclusion still holds for training on the Non-IID dataset.

4.2.3 Effect of D2D Communication

In order to investigate the impact of inter-user D2D communications on the RingSFL system, we varied the rates of the D2D communication links to 135 Mbps, 40 Mbps, 15 Mbps, 10 Mbps, and 5 Mbps, respectively, while keeping the communication rates between the client and the server (both RingSFL and vanilla FL) fixed at 135 Mbps. Subsequently, we evaluated the convergence of the system under these different conditions. Additionally, we conducted experiments on the CIFAR10 dataset with propagation lengths of $L_0:L_1:L_2:L_3:L_4=4:3:1:1:1$ using the ResNet18 model structure to further explore the impact of inter-user

2. Number of iterations over the local dataset in each communication round.

3. Number of batches the local dataset can be divided into.

TABLE 2
Implementation Parameters.

Model	Dataset	Optimizer	Learning Rate (IID / Non-IID)	Blocks	Client Num	Local Epoch	Communication Round
ResNet18	Cifar10	SGD	0.1 / 0.04	10	5	2	100
VGG16	Cifar10	SGD	0.01 / 0.0015	16	5	2	100
AlexNet	MNIST	SGD	0.003 / 0.003	17	5	2	100
LeNet-5	MNIST	SGD	0.02 / 0.02	12	5	2	100

TABLE 3
Top-1 Accuracy (%) of Each Model under Different Algorithms.

	ResNet18 (IID / Non-IID)	VGG16 (IID / Non-IID)	AlexNet (IID / Non-IID)	LeNet-5 (IID / Non-IID)
RingSFLv1	82.35 \pm 0.36 / 48.30 \pm 0.57	79.30 \pm 0.20 / 40.35 \pm 0.99	98.83 \pm 0.11 / 89.58 \pm 0.55	98.82 \pm 0.19 / 94.34 \pm 0.56
RingSFLv2	84.57 \pm 0.17 / 56.80 \pm 0.78	84.33 \pm 0.10 / 41.26 \pm 1.29	99.13 \pm 0.07 / 94.31 \pm 0.88	99.10 \pm 0.04 / 95.75 \pm 0.73
SplitFed	75.92 \pm 0.51 / 30.16 \pm 4.49	72.86 \pm 0.62 / 28.17 \pm 2.15	98.76 \pm 0.09 / 84.00 \pm 4.39	98.74 \pm 0.24 / 93.64 \pm 0.70
vanilla FL	78.93 \pm 0.27 / 48.02 \pm 1.28	77.02 \pm 0.34 / 39.52 \pm 0.81	98.81 \pm 0.07 / 91.60 \pm 1.14	98.84 \pm 0.08 / 94.77 \pm 0.29
vanilla SL	83.41 \pm 0.44 / 26.96 \pm 3.58	78.50 \pm 0.69 / 35.33 \pm 1.29	98.69 \pm 0.10 / 98.84 \pm 0.08	98.80 \pm 0.14 / 98.86 \pm 0.09

TABLE 4
Top-1 Accuracy (%) of Each Model under Different Propagation Lengths.

Propagation Lengths	ResNet18 (IID / Non-IID)	Propagation Lengths	VGG16 (IID / Non-IID)	Propagation Lengths	AlexNet (IID / Non-IID)	Propagation Lengths	LeNet-5 (IID / Non-IID)
6:1:1:1:1	84.66 \pm 0.33 56.45 \pm 1.10	12:1:1:1:1	84.29 \pm 0.14 41.48 \pm 1.08	13:1:1:1:1	99.00 \pm 0.16 94.49 \pm 0.67	8:1:1:1:1	99.10 \pm 0.07 95.85 \pm 0.32
5:2:1:1:1	83.90 \pm 0.29 55.45 \pm 0.47	11:2:1:1:1	83.98 \pm 0.24 42.56 \pm 0.69	11:3:1:1:1	99.05 \pm 0.12 94.28 \pm 0.59	7:2:1:1:1	99.04 \pm 0.06 95.79 \pm 0.30
4:3:1:1:1	83.00 \pm 0.16 53.63 \pm 0.62	10:3:1:1:1	83.78 \pm 0.53 41.68 \pm 0.69	9:5:1:1:1	99.11 \pm 0.10 93.79 \pm 0.30	6:3:1:1:1	99.02 \pm 0.05 95.66 \pm 0.19
3:3:2:1:1	82.24 \pm 0.20 51.34 \pm 0.74	8:3:3:1:1	82.81 \pm 0.25 39.25 \pm 0.62	7:5:3:1:1	99.00 \pm 0.14 93.00 \pm 0.11	5:3:2:1:1	99.00 \pm 0.06 95.65 \pm 0.18
3:2:2:2:1	80.90 \pm 0.19 50.27 \pm 0.53	6:3:3:3:1	80.53 \pm 0.32 37.57 \pm 0.94	5:5:3:3:1	98.91 \pm 0.07 92.14 \pm 0.61	4:3:2:2:1	98.96 \pm 0.04 95.51 \pm 0.16
2:2:2:2:2	79.00 \pm 0.50 47.89 \pm 0.64	4:3:3:3:3	77.58 \pm 0.25 39.76 \pm 0.96	4:4:3:3:3	98.84 \pm 0.12 92.53 \pm 1.03	3:3:2:2:2	98.97 \pm 0.03 95.45 \pm 0.16

D2D communications. However, it is important to note that the results of this subsection are not absolute, as different model structures, datasets, and propagation lengths may lead to different outcomes.

The convergence curves of RingSFL at different D2D communication rates are depicted in Fig. 11. As the D2D communication rate decreases, the time required for RingSFL to reach the same accuracy increases accordingly. Nevertheless, it can be observed that as long as the D2D communication rate is not too bad, RingSFL can reach convergence in a shorter time compared with vanilla FL. Specifically, RingSFL takes less time to converge than vanilla FL when the D2D communication rate is greater than or equal to 10 Mbps. Conversely, when the D2D communication rate is less than or equal to 5 Mbps, RingSFL takes more time to converge than vanilla FL. This is further illustrated in Fig. 12, which shows the time consumption of RingSFL to complete a training round at different D2D communication

rates.

4.2.4 Convergence Time Reduction and Energy Efficiency

To further validate RingSFL's capability to accommodate system heterogeneity, we evaluate its convergence time reduction and energy efficiency on the prototype system. We utilize ResNet18 for training and set the propagation lengths of RingSFL to $L_0:L_1:L_2:L_3:L_4=4:3:1:1:1$.

As illustrated in Fig. 13, compared to vanilla FL, RingSFL reduces the training time by 61.6%, significantly enhancing the system's training efficiency and computational resource utilization. This is attributed to the adaptive model split scheme, which assigns more propagation lengths to clients with strong computational power and correspondingly more computational load, and vice versa. This heterogeneity-adaptive computational load balancing scheme effectively alleviates the straggler effect, thus reducing the training time. In the figure, the training time of

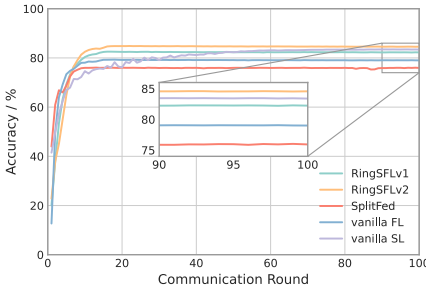


Fig. 7. Testing convergence of ResNet18 on CIFAR10 (IID) under different algorithms.

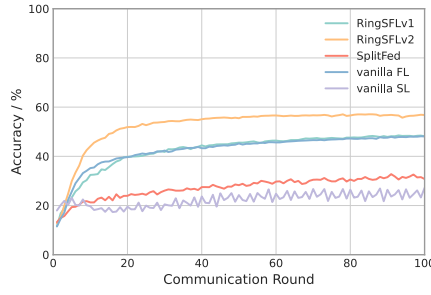


Fig. 8. Testing convergence of ResNet18 on CIFAR10 (Non-IID) under different algorithms.

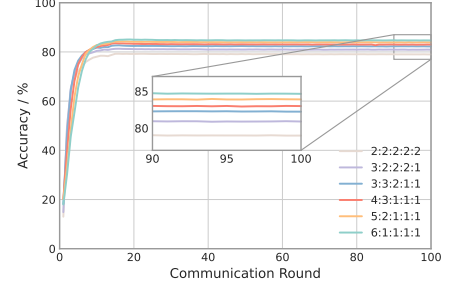


Fig. 9. Testing convergence of ResNet18 on CIFAR10 (IID) under different propagation lengths.

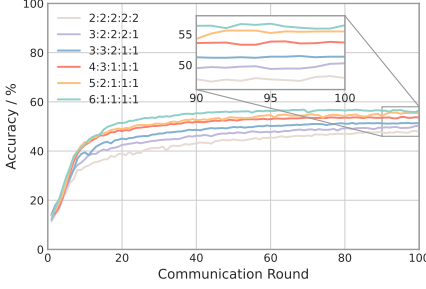


Fig. 10. Testing convergence of ResNet18 on CIFAR10 (Non-IID) under different propagation lengths.

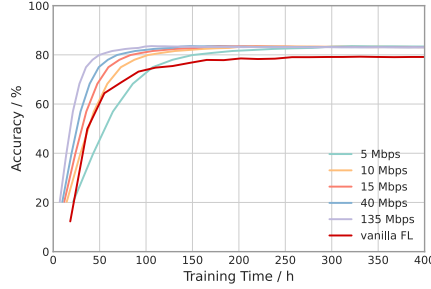


Fig. 11. Testing convergence of ResNet18 on CIFAR10 (IID) under different D2D communication rates.

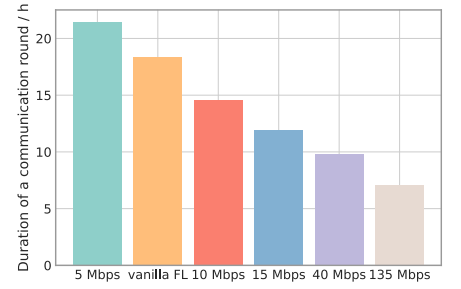


Fig. 12. Time cost of ResNet18 in a communication round under different D2D communication rates.

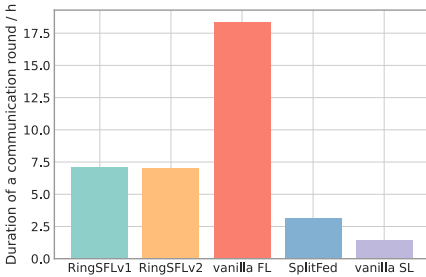


Fig. 13. Time cost of ResNet18 in a communication round under different algorithms.

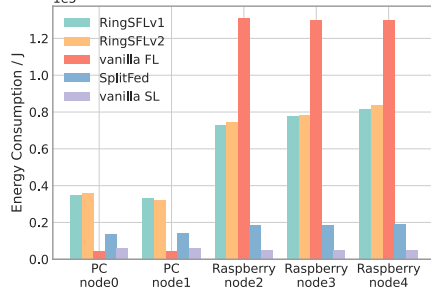


Fig. 14. Energy consumption of different devices in a communication round.

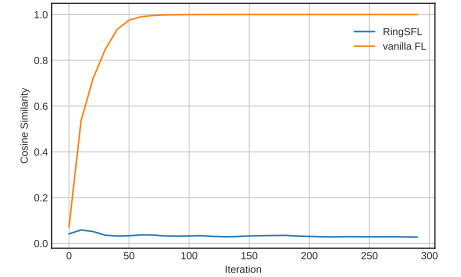


Fig. 15. Cosine similarity between reconstructed data and original data.

SplitFed and vanilla SL is lower than RingSFL due to the model split mechanism, which offloads a large portion of the training tasks to the server with high computation power. Specifically, in our experiments, the client computational volume of SplitFed in each round is 1/5 of RingSFL, and the client computational volume of vanilla SL is 1/25 of RingSFL. However, they cannot achieve comparable performance in terms of model accuracy and user privacy, as discussed above.

Fig. 14 illustrates the capability of RingSFL to regulate the power consumption of the FL system. It is evident that, in comparison to vanilla FL, RingSFL increases the energy consumption of PC nodes and decreases the energy consumption of Raspberry Pi nodes, which is attributed to the ability of RingSFL to balance the computational loads among clients. This observation demonstrates that RingSFL can effectively conserve the energy of the battery-powered devices, thus enhancing the life time of the RingSFL system.

4.2.5 Privacy Preservation

To assess RingSFL's capacity to withstand data reconstruction attacks, we conducted attack experiments on both vanilla FL and RingSFL using the methods outlined in [10]. We attempted to recover the training data by eavesdropping on the model parameters transmitted between the client and the server. To measure the efficacy of the data recovery, we calculated the cosine similarity of the recovered data to the original data, which is formulated as

$$s = \frac{\mathbf{x}^T \hat{\mathbf{x}}}{\|\mathbf{x}\|_2 \cdot \|\hat{\mathbf{x}}\|_2} \quad (9)$$

where \mathbf{x} denote the original data, $\hat{\mathbf{x}}$ denote the recovered data, and $\|\cdot\|_2$ denote the L2 norm of the vector. The cosine similarity between \mathbf{x} and $\hat{\mathbf{x}}$ can be used to measure the similarity between the two, with a value of 1 indicating perfect similarity. In this experiment, the neural network employed was LeNet and the dataset used was CIFAR10. The system was composed of five clients, and the propagation length was set to $L_0:L_1:L_2:L_3:L_4=3:3:2:2:2$.

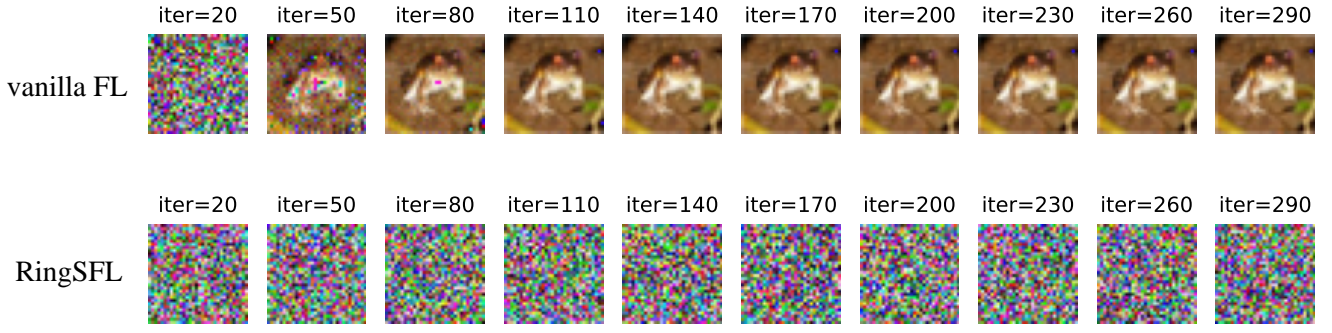


Fig. 16. Reconstructed data after attacking vanilla FL and RingSFL.

The results of the data reconstruction attack are depicted in Fig. 16, where the first row corresponds to the attack on vanilla FL and the second row to the attack on RingSFL. It is evident that, after 80 iterations, the training data in vanilla FL has been successfully recovered, while the attack on RingSFL fails to recover the training data. Additionally, Fig. 15 shows that the cosine similarity between the reconstructed data and the original data has converged to 1 in the attack on vanilla FL after 80 iterations, while the cosine similarity remains around 0 in the attack on RingSFL. This is attributed to the transfer of the complete model parameters by vanilla FL, which allows the attacker to extract the complete information about the training data, thus enabling the recovery of the original training data and resulting in user privacy leakage. In contrast, the model transmitted by RingSFL is a blended model, making it difficult for an attacker to extract information about the training data from such fragmented and mixed model parameters, thus protecting user privacy effectively.

5 LIMITATIONS AND FUTURE WORKS

5.1 Incentive Mechanism Design

Fairness is a fundamental principle of multi-client collaborative training, and clients are unlikely to participate in training when they perceive it to be unfair [48]. In the RingSFL system, however, a higher computational load is assigned to clients with higher computational power, resulting in increased consumption of computation resources and battery energy. To motivate clients to participate in the training and contribute more resources, effective incentive mechanisms are necessary, such as monetary rewards. Nevertheless, verifying the actual contribution from the training results is a challenge, as each client only trains a part of the model.

5.2 Communication Cost Reduction

Compared to traditional FL, RingSFL introduces additional communication between clients, which may render it inapplicable in scenarios where communication resources are limited. To address this issue, further improvement of RingSFL is necessary to reduce the total amount of communication. One potential solution is the introduction of the Ring-Allreduce mechanism in RingSFL, which aggregates

models by passing model parameters only among clients, thus avoiding the transmission of model parameters between clients and the server, reducing the number of parameters communicated, and enhancing the system's security.

5.3 Communication Topology Construction

In RingSFL, the ring topology is a critical factor in the training process and model performance. Consequently, it is essential to devise an effective algorithm for constructing the ring topology. As many factors, such as geographical location of users, channel conditions between users, the computational power of users, quality of user data sets, etc., can influence the communication topology, traditional optimization methods may not be applicable. Therefore, machine learning-based methods, such as reinforcement learning, can be explored as a potential research direction for constructing communication topologies.

6 CONCLUSION

In this paper, we have proposed a novel FL scheme, dubbed RingSFL, which integrates FL with an efficient model split mechanism to adapt to system heterogeneity while preserving data privacy. Experiments conducted on both the simulation environment and prototype system have demonstrated that RingSFL can achieve faster convergence and higher accuracy than the benchmarks, and have better privacy performance. Furthermore, RingSFL can be applied to scenarios with significant system heterogeneity to enhance the overall system efficiency. For future works, we will design an effective incentive mechanism and learning-based optimal topology construction scheme for RingSFL.

REFERENCES

- [1] S. Khan, M. Naseer, M. Hayat, S. W. Zamir, F. S. Khan, and M. Shah, "Transformers in vision: A survey," *ACM computing surveys (CSUR)*, vol. 54, no. 10s, pp. 1–41, 2022.
- [2] D. W. Otter, J. R. Medina, and J. K. Kalita, "A Survey of the Usages of Deep Learning for Natural Language Processing," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 2, pp. 604–624, 2021.
- [3] J. Li *et al.*, "Recent advances in end-to-end automatic speech recognition," *APSIPA Transactions on Signal and Information Processing*, vol. 11, no. 1, 2022.

- [4] Y. Wang, M. Chen, T. Luo, W. Saad, D. Niyato, H. V. Poor, and S. Cui, "Performance Optimization for Semantic Communications: An Attention-Based Reinforcement Learning Approach," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 9, pp. 2598–2613, 2022.
- [5] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, "Communication-Efficient Learning of Deep Networks from Decentralized Data," in *Proc. International Conference on Artificial Intelligence and Statistics*, 2017, pp. 1273–1282.
- [6] W. Y. B. Lim, N. C. Luong, D. T. Hoang, Y. Jiao, Y.-C. Liang, Q. Yang, D. Niyato, and C. Miao, "Federated Learning in Mobile Edge Networks: A Comprehensive Survey," *IEEE Communications Surveys Tutorials*, vol. 22, no. 3, pp. 2031–2063, 2020.
- [7] Z. Yang, M. Chen, K.-K. Wong, H. V. Poor, and S. Cui, "Federated learning for 6G: Applications, challenges, and opportunities," *Engineering*, vol. 8, pp. 33–41, 2022.
- [8] C. Yang, M. Xu, Q. Wang, Z. Chen, K. Huang, Y. Ma, K. Bian, G. Huang, Y. Liu, X. Jin *et al.*, "FLASH: Heterogeneity-Aware Federated Learning at Scale," *IEEE Transactions on Mobile Computing*, 2022.
- [9] N. Cheng, F. Lyu, W. Quan, C. Zhou, H. He, W. Shi, and X. Shen, "Space/Aerial-Assisted Computing Offloading for IoT Applications: A Learning-Based Approach," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 5, pp. 1117–1129, 2019.
- [10] L. Zhu, Z. Liu, and S. Han, "Deep Leakage from Gradients," in *Proc. Advances in Neural Information Processing Systems*, 2019, pp. 1–11.
- [11] H. Yin, A. Mallya, A. Vahdat, J. M. Alvarez, J. Kautz, and P. Molchanov, "See Through Gradients: Image Batch Recovery via GradInversion," in *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 16 337–16 346.
- [12] Y. Deng, F. Lyu, J. Ren, H. Wu, Y. Zhou, Y. Zhang, and X. Shen, "AUCTION: Automated and Quality-Aware Client Selection Framework for Efficient Federated Learning," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 8, pp. 1996–2009, 2022.
- [13] J. Xu and H. Wang, "Client Selection and Bandwidth Allocation in Wireless Federated Learning Networks: A Long-Term Perspective," *IEEE Transactions on Wireless Communications*, vol. 20, no. 2, pp. 1188–1200, 2021.
- [14] M. Zhang, G. Zhu, S. Wang, J. Jiang, Q. Liao, C. Zhong, and S. Cui, "Communication-efficient federated edge learning via optimal probabilistic device scheduling," *IEEE Transactions on Wireless Communications*, vol. 21, no. 10, pp. 8536–8551, 2022.
- [15] N. H. Tran, W. Bao, A. Zomaya, M. N. H. Nguyen, and C. S. Hong, "Federated Learning over Wireless Networks: Optimization Model Design and Analysis," in *Proc. IEEE INFOCOM 2019 - IEEE Conference on Computer Communications*, 2019, pp. 1387–1395.
- [16] Z. Yang, M. Chen, W. Saad, C. S. Hong, and M. Shikh-Bahaei, "Energy efficient federated learning over wireless communication networks," *IEEE Transactions on Wireless Communications*, vol. 20, no. 3, pp. 1935–1949, 2020.
- [17] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, "Tackling the Objective Inconsistency Problem in Heterogeneous Federated Optimization," in *Proc. Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 7611–7623.
- [18] Y. Deng, F. Lyu, J. Ren, Y. Zhang, Y. Zhou, Y. Zhang, and Y. Yang, "SHARE: Shaping Data Distribution at Edge for Communication-Efficient Hierarchical Federated Learning," in *Proc. International Conference on Distributed Computing Systems*, 2021, pp. 24–34.
- [19] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep Learning with Differential Privacy," in *Proc. ACM SIGSAC Conference on Computer and Communications Security*, 2016, pp. 308–318.
- [20] K. Wei, J. Li, M. Ding, C. Ma, H. Su, B. Zhang, and H. V. Poor, "User-level privacy-preserving federated learning: Analysis and performance optimization," *IEEE Transactions on Mobile Computing*, vol. 21, no. 9, pp. 3388–3401, 2021.
- [21] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. S. Quek, and H. V. Poor, "Federated Learning With Differential Privacy: Algorithms and Performance Analysis," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 3454–3469, 2020.
- [22] D. Liu and O. Simeone, "Privacy for Free: Wireless Federated Learning via Uncoded Transmission With Adaptive Power Control," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 1, pp. 170–185, 2021.
- [23] S. Truex, N. Baracaldo, A. Anwar, T. Steinke, H. Ludwig, R. Zhang, and Y. Zhou, "A Hybrid Approach to Privacy-Preserving Federated Learning," in *Proc. ACM Workshop on Artificial Intelligence and Security*, 2019, pp. 1–11.
- [24] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for privacy-preserving machine learning," in *Proc. ACM SIGSAC Conference on Computer and Communications Security*, 2017, pp. 1175–1191.
- [25] B. Zhao, X. Liu, W.-N. Chen, and R. Deng, "CrowdFL: privacy-preserving mobile crowdsensing system via federated learning," *IEEE Transactions on Mobile Computing*, 2022.
- [26] R. Shokri and V. Shmatikov, "Privacy-Preserving Deep Learning," in *Proc. ACM SIGSAC Conference on Computer and Communications Security*, 2015, pp. 1310–1321.
- [27] A. Triastcyn and B. Faltings, "Federated Generative Privacy," *IEEE Intelligent Systems*, vol. 35, no. 4, pp. 50–57, 2020.
- [28] Y. Lu, X. Huang, Y. Dai, S. Maharjan, and Y. Zhang, "Blockchain and Federated Learning for Privacy-Preserved Data Sharing in Industrial IoT," *IEEE Transactions on Industrial Informatics*, vol. 16, no. 6, pp. 4177–4186, 2020.
- [29] Y. Qu, L. Gao, T. H. Luan, Y. Xiang, S. Yu, B. Li, and G. Zheng, "Decentralized Privacy Using Blockchain-Enabled Federated Learning in Fog Computing," *IEEE Internet of Things Journal*, vol. 7, no. 6, pp. 5171–5183, 2020.
- [30] O. Gupta and R. Raskar, "Distributed learning of deep neural network over multiple agents," *Journal of Network and Computer Applications*, vol. 116, pp. 1–8, 2018.
- [31] P. Vepakomma, O. Gupta, T. Swedish, and R. Raskar, "Split learning for health: Distributed deep learning without sharing raw patient data," *arXiv preprint arXiv:1812.00564*, 2018.
- [32] O. Li, J. Sun, X. Yang, W. Gao, H. Zhang, J. Xie, V. Smith, and C. Wang, "Label Leakage and Protection in Two-party Split Learning," in *International Conference on Learning Representations*, 2021.
- [33] S. Pal, M. Uniyal, J. Park, P. Vepakomma, R. Raskar, M. Bennis, M. Jeon, and J. Choi, "Server-Side Local Gradient Averaging and Learning Rate Acceleration for Scalable Split Learning," *arXiv preprint arXiv:2112.05929*, 2021.
- [34] C. Thapa, P. C. M. Arachchige, S. Camtepe, and L. Sun, "Splitfed: When federated learning meets split learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 8, 2022, pp. 8485–8493.
- [35] D.-J. Han, H. I. Bhatti, J. Lee, and J. Moon, "Accelerating federated learning with split learning on locally generated losses," in *ICML 2021 Workshop on Federated Learning for User Privacy and Data Confidentiality*. ICML Board, 2021.
- [36] S. Oh, J. Park, P. Vepakomma, S. Baek, R. Raskar, M. Bennis, and S.-L. Kim, "LocFedMix-SL: Localize, Federate, and Mix for Improved Scalability, Convergence, and Latency in Split Learning," in *Proceedings of the ACM Web Conference 2022*, 2022, pp. 3347–3357.
- [37] X. Liu, Y. Deng, and T. Mahmoodi, "Wireless distributed learning: A new hybrid split and federated learning approach," *IEEE Transactions on Wireless Communications*, 2022.
- [38] A. Abedi and S. S. Khan, "FedSL: Federated Split Learning on Distributed Sequential Data in Recurrent Neural Networks," *arXiv preprint arXiv:2011.03180*, 2020.
- [39] V. Turina, Z. Zhang, F. Esposito, and I. Matta, "Federated or Split? A Performance and Privacy Analysis of Hybrid Split and Federated Learning Architectures," in *2021 IEEE 14th International Conference on Cloud Computing (CLOUD)*. IEEE, 2021, pp. 250–260.
- [40] Y. Tian, Y. Wan, L. Lyu, D. Yao, H. Jin, and L. Sun, "FedBERT: When Federated Learning Meets Pre-Training," *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2022.
- [41] D. Romanini, A. J. Hall, P. Papadopoulos, T. Titcombe, A. Ismail, T. Cebere, R. Sandmann, R. Roehm, and M. A. Hoeh, "Pyvertical: A vertical federated learning framework for multi-headed splitnn," *arXiv preprint arXiv:2104.00489*, 2021.
- [42] F. Jameel, Z. Hamid, F. Jabeen, S. Zeedally, and M. A. Javed, "A Survey of Device-to-Device Communications: Research Issues and Challenges," *IEEE Communications Surveys Tutorials*, vol. 20, no. 3, pp. 2133–2168, 2018.
- [43] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [44] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.

- [45] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [46] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [47] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.
- [48] Y. Zhan, J. Zhang, Z. Hong, L. Wu, P. Li, and S. Guo, "A Survey of Incentive Mechanism Design for Federated Learning," *IEEE Transactions on Emerging Topics in Computing*, pp. 1–1, 2021.



Michael Shell Biography text here.

John Doe Biography text here.

Jane Doe Biography text here.