

# HASFL: Heterogeneity-aware Split Federated Learning over Edge Computing Systems

Zheng Lin, Zhe Chen *Member, IEEE*, Xianhao Chen, *Member, IEEE*, Wei Ni, *Fellow, IEEE*, and Yue Gao, *Fellow, IEEE*

**Abstract**—Split federated learning (SFL) has emerged as a promising paradigm to democratize machine learning (ML) on edge devices by enabling layer-wise model partitioning. However, existing SFL approaches suffer significantly from the straggler effect due to the heterogeneous capabilities of edge devices. To address the fundamental challenge, we propose adaptively controlling batch sizes (BSs) and model splitting (MS) for edge devices to overcome resource heterogeneity. We first derive a tight convergence bound of SFL that quantifies the impact of varied BSs and MS on learning performance. Based on the convergence bound, we propose HASFL, a heterogeneity-aware SFL framework capable of adaptively controlling BS and MS to balance communication-computing latency and training convergence in heterogeneous edge networks. Extensive experiments with various datasets validate the effectiveness of HASFL and demonstrate its superiority over state-of-the-art benchmarks.

**Index Terms**—Federated learning, split federated learning, batch size, model splitting, mobile edge computing.

## I. INTRODUCTION

Conventional machine learning (ML) frameworks predominantly rely on centralized learning (CL), where raw data is gathered and processed at a central server for model training. However, CL is often impractical due to its high communication latency, increased backbone traffic, and privacy risks [1]–[4]. To address these limitations, federated learning (FL) [5], [6] has emerged as a promising alternative that allows participating devices to collaboratively train a shared model via exchanging model parameters (e.g., gradients) rather than raw data, thereby protecting data privacy and reducing communication costs [7], [8]. Despite its advantage, on-device training of FL poses a significant challenge for its deployment on resource-constrained edge devices as ML models scale up [9], [10]. Considering the fact that training costs significantly more latency and memory space than model inference, training large-sized models, such as Gemini Nano-2 with 3.25 billion parameters (3GB for 32-bit floats), is computationally prohibitive for resource-constrained edge devices [11].

Split learning (SL) [12] has emerged as a viable solution to overcoming the weaknesses of FL, which offloads the

Z. Lin and X. Chen are with the Department of Electrical and Electronic Engineering, University of Hong Kong, Pok Fu Lam, Hong Kong, China (e-mail: linzheng@eee.hku.hk; xchen@eee.hku.hk).

Z. Chen and Y. Gao are with the Institute of Space Internet, Fudan University, Shanghai 200438, China, and the School of Computer Science, Fudan University, Shanghai 200438, China (e-mail: zhechen@fudan.edu.cn; gao.yue@fudan.edu.cn).

W. Ni is with Data61, CSIRO, Marsfield, NSW 2122, Australia, and the School of Computing Science and Engineering, and the University of New South Wales, Kensington, NSW 2052, Australia (e-mail: wei.ni@ieee.org).

(Corresponding author: Xianhao Chen)

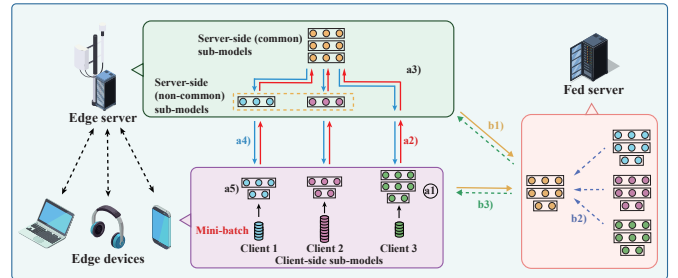


Fig. 1. The illustration of HASFL over edge computing systems, where a1) and a5) denote client-side forward propagation (FP) and backward pass (BP), a3) represents server-side model FP and BP, a2) and a4) are activations and activations' gradients transmissions, b1) and b3) denote sub-model uploading and downloading, and b2) represents client-side model aggregation.

major workload from edge devices to a more powerful central server via layer-wise model splitting, alleviating computing workload and memory requirements on edge devices [13]–[15]. The original SL framework, named vanilla SL [12], employs an inter-device sequential training where the central server collaborates with one edge device after another for co-training, however, at the expense of excessive training latency. To address this, a prominent variant of SL, split federated learning (SFL) [16], integrates the strengths of FL and SL to enable parallel training across edge devices. As shown in Fig. 1, SFL not only involves model splitting but also adheres to FL principles, requiring the fed server to periodically aggregate client sub-models from edge devices in a synchronous manner, akin to FedAvg [5].

Unfortunately, deploying SFL in heterogeneous edge systems poses significant challenges. The SFL framework enforces a synchronous aggregation protocol that necessitates every device to complete its local training before model aggregation. However, the heterogeneous computing and communication resources across edge devices lead to significantly varied training delays, causing a severe straggler effect<sup>1</sup> [17]–[19]. Although existing studies [20]–[23] attempt to mitigate this issue by determining varied cut layers for different devices, optimizing model splitting (MS) alone is insufficient as MS *lacks the flexibility* of adjusting client-side workload. For instance, while splitting a convolutional neural network (CNN) at a shallow point leads to lower computing workload

<sup>1</sup>The straggler effect is a primary bottleneck in distributed learning, referring to the problem of edge devices with slower training speeds (i.e., stragglers) impeding the overall model training progress. Specifically, all edge devices must wait for the slowest one to complete model training before proceeding.

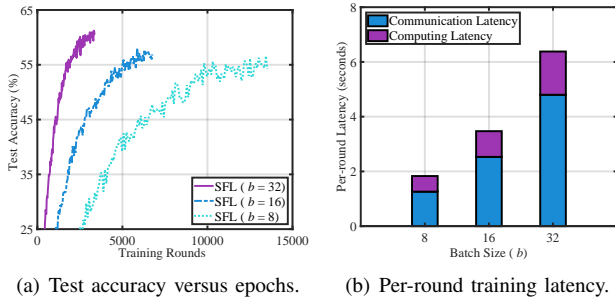


Fig. 2. The impact of BS on training performance and per-round training latency. Fig. 2(a) shows the performance for test accuracy versus number of epochs, and Fig. 2(b) illustrates the effect of BS on per-round training latency. The experiment is conducted on the CIFAR-10 dataset under the non-IID setting with  $L_c = 8$  and  $I = 15$ .

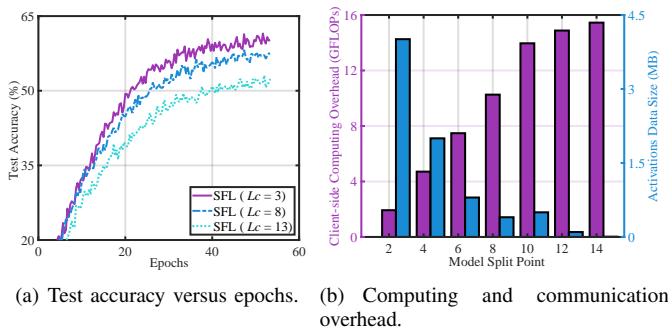


Fig. 3. The impact of MS on computing and communication overhead as well as training performance. Fig. 3(b) shows the computing and communication overhead of SFL at different model split points. Fig. 3(a) presents the performance for test accuracy versus the number of epochs. The experiment is conducted on the CIFAR-10 dataset under the non-IID setting with  $b = 16$  and  $I = 15$ .

on clients, it also increases communication overhead due to a larger output size in early layers [21].

To effectively alleviate the straggler effect, BS serves as another key control variable. Our key insight is that *both* the communication and computing workload of an edge device scales approximately linearly with its BS, making BS a flexible lever for customizing client-side workload. Moreover, BS and MS significantly influence the performance of SFL not only in latency but also in training accuracy: i) **Impact of BS:** As shown in Fig. 2, increasing parameter  $b$  (batch size) expedites model convergence, but this improvement comes at the cost of increased per-round communication-computing latency. In particular, smashed data (activations) communication latency is linearly proportional to BSs of edge devices – a phenomenon that does not exist in conventional FL settings. ii) **Impact of MS:** As illustrated in Fig. 3, MS introduces a complicated trade-off among computing, communications, and training accuracy. Selecting a shallower cut layer (i.e., smaller  $L_c$ ) enhances training accuracy by allocating a larger portion of the model as the server-side sub-model<sup>2</sup>. However, it also

<sup>2</sup>In SFL, the aggregation frequencies of client-side and server-side sub-models may differ. Specifically, server-side sub-models co-located on edge servers can be synchronized in each training round without introducing communication overhead. In contrast, client-side sub-model aggregation incurs communication costs and is therefore performed periodically over several training rounds. A shallower cut layer indicates a larger portion of the model being aggregated at each round (i.e., more frequent model aggregation), leading to better training convergence [24].

incurs higher communication overhead due to the larger output dimensions of shallower layers in CNNs.

Optimizing BS and MS necessitates a fundamental trade-off analysis of SFL under resource constraints, which is not yet studied. To fill the void, in this paper, we propose HASFL, a heterogeneity-aware SFL framework with adaptive BS and MS to jointly determine varied BS and MS under heterogeneous SFL systems. We first derive an analytical convergence upper bound that quantifies the impact of BS and MS on training convergence, and then formulate the problem to minimize total training latency by jointly optimizing BS and MS. The challenges addressed by HASFL are twofold: i) From the convergence analysis perspective, unlike FL with varied batch sizes [25]–[27], the server-side sub-model training in SFL is equivalent to concatenating the entire batch from all clients, whereas the client-side sub-model training runs on its local batch. The discrepancy also makes BS and MS jointly affect training convergence, thus resulting in intricate convergence behavior. We highlight that this paper provides the first convergence bound of SFL with *varied* BSs and cut layers for clients. ii) From a system optimization perspective, both BS and MS balance the tradeoff between training convergence and per-round latency, resulting in a complicated convergence-latency tradeoff.

The main contributions of this paper are summarized as follows.

- We propose HASFL, a heterogeneity-aware SFL framework, which controls BS and MS to accelerate SFL in heterogeneous resource-constrained edge computing systems.
- We establish the *first* convergence bound that rigorously quantifies the impact of varied BSs and MS on SFL, and formulate a new problem to jointly optimize BS and MS.
- We decompose the new problem into two tractable sub-problems of optimizing BS and MS, and develop an efficient optimization method to solve the sub-problems alternately.
- We conduct extensive simulations across various datasets to validate our analysis and demonstrate the superiority of the proposed solution over the state-of-the-art in training accuracy and convergence speed.

The remainder of this paper is organized as follows. Section II elaborates on related work and Section III introduces the system model and HASFL framework. Section IV provides the convergence analysis of HASFL. We formulate the optimization problem in Section V and offer the corresponding solution approach in Section VI. Section VII provides the simulation results. Finally, concluding remarks are presented in Section VIII.

## II. RELATED WORK

Some studies have been conducted to enhance the training performance of FL via BS control [25]–[30]. Shi *et al.* [28] proposed a dynamic BS assisted FL framework that dynamically controls BS to minimize the total latency of FL over mobile devices. Ma *et al.* [26] developed an efficient FL algorithm that adaptively adjusts BS with scaled learning rate for heterogeneous devices to reduce their waiting time and

improve the model convergence rate. Liu *et al.* [25] jointly optimized the BS, compression ratio, and spectrum allocation to maximize the convergence rate under the given training latency constraint. Liu *et al.* [27] first established the convergence bound for training error considering heterogeneous datasets across devices and then derived closed-form solutions for co-optimized BS and aggregation frequency configuration. Smith *et al.* [29] and Yu *et al.* [30] empirically designed dynamic BS schemes to establish a more communication-efficient FL framework. Nevertheless, these existing BS control schemes cannot be directly applied to SFL, as aggregation discrepancies between client-side and server-side sub-models lead to distinct convergence behaviors from conventional FL. This phenomenon, as alluded to earlier, causes the intricate relationship between BS and MS in SFL.

MS heavily impacts the training latency of SL, as it determines the computing load distribution between devices and server, the size of smashed data, and model convergence. Tremendous research efforts have been made to determine the optimal model split point for SL [20], [22], [31]. Wu *et al.* [20] developed a cluster-based parallel SL framework and a joint MS, device clustering, and subchannel allocation optimization algorithm to minimize the training latency under heterogeneous device capabilities and dynamic network conditions. Lee *et al.* [22] proposed a hierarchical Stackelberg game framework for SFL that jointly optimizes MS and incentive mechanisms to balance the training load between the server and clients, ensuring the necessary level of privacy for the clients. Lin *et al.* [31] designed an efficient parallel SL framework that aggregates last-layer gradients to reduce the dimensionality of activations' gradients and devised a joint MS, subchannel allocation, and power control to reduce the overall training latency. Unfortunately, these works fail to consider the impact of MS on training convergence.

The prior work [32] conducted a convergence analysis for SFL, but does not consider the impact of MS, which plays a critical role in the system performance of SFL. To the best of our knowledge, our recent works [21], [23] were the first to provide a theoretical convergence analysis for SFL by considering the impact of MS. Since these studies assumed that all edge devices employ the same BS, neither convergence analysis nor system optimization of SFL can be directly applied to HASFL, both of which are non-trivial.

### III. THE HETEROGENEITY-AWARE SFL FRAMEWORK

This section presents the system model in Section III-A and details the proposed heterogeneity-aware SFL framework, HASFL, in Section III-B.

#### A. System Model

As illustrated in Fig. 1, we consider a typical edge computing scenario of HASFL, which comprises three components:

- **Edge devices:** We consider that each client possesses an edge device capable of executing local computations, i.e., the client-side forward propagation (FP) and backward pass (BP). The set of participating edge devices is represented as  $\mathcal{N} = \{1, 2, \dots, N\}$ , where  $N$  is the total number of edge devices. For the  $i$ -th edge device,

local dataset is  $\mathcal{D}_i = \{\mathbf{x}_{i,k}, y_{i,k}\}_{k=1}^{|\mathcal{D}_i|}$ , where  $\mathbf{x}_{i,k}$  and  $y_{i,k}$  are the  $k$ -th input data sample and its corresponding label, respectively. The client-side sub-model of the  $i$ -th edge device is represented as  $\mathbf{w}_{c,i}$ .

- **Edge server:** The edge server is a computing entity with powerful computing capability, responsible for executing the server-side model training. The server-side sub-model of the  $i$ -th edge device is denoted by  $\mathbf{w}_{s,i} = [\mathbf{h}_s; \mathbf{h}_{m,i}]$ , where  $\mathbf{h}_s$  and  $\mathbf{h}_{m,i}$  are the server-side common and server-side non-common sub-models, respectively. The common sub-model is the shared component across all clients and is synchronized in every training round. The non-common sub-model arises from the discrepancies in the number of layers maintained on the server across clients, and it is periodically aggregated with client-side sub-model every several training rounds. Moreover, the edge server is also tasked with network information collection, such as device computing capabilities and channel conditions, to implement optimization decisions of HASFL.
- **Fed server:** The fed server takes charge of the synchronization of client-side sub-models and server-side non-common sub-models, periodically aggregating the client-side sub-models together with server-side non-common sub-models across all participating edge devices. For privacy concerns, fed and edge servers usually belong to separate parties, as possessing both the client-side sub-models and smashed data could potentially enable a malicious server to reconstruct the original user data [33].

The global model is represented as  $\mathbf{w} = [\mathbf{w}_{s,i}; \mathbf{w}_{c,i}]$ . The goal of SL is to obtain the optimal global model  $\mathbf{w}^*$  by minimizing the following finite-sum non-convex global loss function:

$$\min_{\mathbf{w}} f(\mathbf{w}) \triangleq \min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{w}), \quad (1)$$

where  $f_i(\mathbf{w}) \triangleq \mathbb{E}_{\xi_i \sim \mathcal{D}_i} [F_i(\mathbf{w}; \xi_i)]$  denotes the local loss function of the  $i$ -th edge device and  $\xi_i$  is training randomness<sup>3</sup> from the local dataset  $\mathcal{D}_i$ . For conciseness, we consider that all edge devices have local datasets with uniform size. Consistent with the standard setting [34]–[36], it is assumed that the stochastic gradient is an unbiased estimate of the true gradient, i.e.,  $\mathbb{E}_{\xi_i^t \sim \mathcal{D}_i} [\nabla F_i(\mathbf{w}_i^{t-1}; \xi_i^t) | \xi^{[t-1]}] = \nabla f_i(\mathbf{w}_i^{t-1})$ , where  $\nabla_{\mathbf{w}} F(\mathbf{w}; \xi)$  is the gradient of the function  $F(\mathbf{w}; \xi)$  with respect to model parameter  $\mathbf{w}$  and  $\xi^{[t-1]} \triangleq [\xi_i^\tau]_{i \in \{1, 2, \dots, N\}, \tau \in \{1, \dots, t-1\}}$  represents all training randomness up to the  $(t-1)$ -th training round.

To solve problem (1), conventional SFL frameworks employ a uniform BS and fixed MS across edge devices throughout model training. However, this scheme leads to severe straggler effects in heterogeneous edge computing systems, significantly slowing down model convergence. Motivated by this, we propose a HASFL featuring heterogeneity-aware BS and MS, as will be detailed in the following section.

<sup>3</sup>Training randomness refers to the stochasticity introduced during model training, primarily due to mini-batch data sampling and random data shuffling within the local dataset [34]–[36].

## B. HASFL Procedure

This section presents the workflow of the proposed HASFL framework. The salient feature of HASFL lies in heterogeneity-aware BS and MS. By jointly optimizing BS and MS, HASFL can significantly reduce the overall training latency while retaining the desired learning performance.

Before model training starts, the edge server initializes the ML model. Then, the optimal BS is determined (see Section VI), and the global model is partitioned into client-side and server-side sub-models via optimal MS (see Section VI). Afterwards, HASFL aggregates client-side sub-models and updates BS and MS based on device computing and communication resources every  $I$  training rounds. This process loops until the model converges. The training procedure of HASFL consists of two primary stages: split training and client-side model aggregation. The split training is executed in every training round, while the client-side model aggregation is triggered every  $I$  training rounds. As illustrated in Fig. 1, for any training round  $t$ , i.e.,  $t \in \mathcal{R} = \{1, 2, \dots, R\}$ , the training workflow of HASFL is presented as follows.

*a. The split training stage:* This stage involves client-side and server-side model updates and smashed data exchange in each training round, comprising the following five steps.

*a1) Client-side model forward propagation:* This step involves the parallel execution of client-side FP of edge devices. For the arbitrary  $t$ -th training round, each edge device  $i$  randomly samples a mini-batch  $\mathcal{B}_i^t \subseteq \mathcal{D}_i$  containing  $b_i$  data samples from its local dataset for model training. For the  $t$ -th training round, the input data and corresponding labels of the mini-batch in training round  $t$  are represented as  $\mathbf{x}_i^t$  and  $\mathbf{y}_i^t$ , respectively. The client-side sub-model of the  $i$ -th edge device trained up to the  $(t-1)$ -th round, is denoted by  $\mathbf{w}_{c,i}^{t-1}$ . After feeding a mini-batch into the client-side sub-model, activations are generated at the cut layer. The activations of the  $i$ -th edge device are expressed as

$$\mathbf{a}_i^t = \varphi(\mathbf{x}_i^t; \mathbf{w}_{c,i}^{t-1}), \quad (2)$$

where  $\varphi(\mathbf{x}; \mathbf{w})$  is the mapping function between input data  $\mathbf{x}$  and its predicted value given model parameter  $\mathbf{w}$ .

*a2) Activations transmissions:* After the client-side FP is completed, each edge device transmits its generated activations and corresponding labels to the edge server (usually over wireless channels).

*a3) Server-side model forward propagation and backward pass:* The edge server gathers activations from participating edge devices and then feeds these activations into the server-side sub-models to perform server-side FP. For the  $i$ -th edge device, the predicted value of the server-side sub-model is represented as

$$\hat{\mathbf{y}}_i^t = \varphi(\mathbf{a}_i^t; \mathbf{w}_{s,i}^{t-1}), \quad (3)$$

where  $\mathbf{w}_{s,i}^{t-1} = [\mathbf{h}_s^{t-1}; \mathbf{h}_{m,i}^{t-1}]$ ;  $\mathbf{h}_s^{t-1}$  and  $\mathbf{h}_{m,i}^{t-1}$  are the server-side common and server-side non-common sub-models, respectively. The predicted value and labels are utilized to compute the loss function and derive the server-side sub-model's gradients.

Due to heterogeneous cut layers across edge devices, the server-side sub-model consists of two parts: the common sub-model, shared by all edge devices, and the non-common sub-model, formed by the extra layers from edge devices with deeper cut layers. These structural disparities necessitate distinct update mechanisms for each component.

For the server-side common sub-model, the edge server updates it<sup>4</sup> every training round by averaging the common sub-model updates from all edge devices, i.e.,

$$\mathbf{h}_s^t = \frac{1}{N} \sum_{i=1}^N \mathbf{h}_{s,i}^t, \quad (4)$$

where  $\mathbf{h}_{s,i}^t \leftarrow \mathbf{h}_{s,i}^{t-1} - \gamma \nabla_{\mathbf{h}_s} F_i(\mathbf{h}_{s,i}^{t-1}; \xi_i^t)$  is the server-side common sub-model of the  $i$ -th edge device,  $\nabla_{\mathbf{h}_s} F_i(\mathbf{h}_{s,i}^{t-1}; \xi_i^t)$  is the stochastic gradients of server-side common sub-models of the  $i$ -th edge device, and  $\gamma$  represents the learning rate. Since the aggregation step in Eqn. (4) does not incur any communication overhead, it can be executed at every training round to expedite training convergence [21], [23]. Consequently, the update process for the server-side common sub-model is equivalent to centralized training with stochastic gradient descent.

In contrast, the server-side non-common sub-models, specific to the edge devices with deeper cut layers, are updated individually on their respective devices without cross-device aggregation. The update procedure for the server-side non-common sub-model of the  $i$ -th edge device is given by

$$\mathbf{h}_{m,i}^t \leftarrow \mathbf{h}_{m,i}^{t-1} - \gamma \nabla_{\mathbf{h}_m} F_i(\mathbf{h}_{m,i}^{t-1}; \xi_i^t), \quad (5)$$

where  $\nabla_{\mathbf{h}_m} F_i(\mathbf{h}_{m,i}^{t-1}; \xi_i^t)$  gives the stochastic gradients of server-side non-common sub-models of the  $i$ -th edge device.

*a4) Activations' gradients transmissions:* After the completion of server-side BP, the edge server sends the activations' gradients to the respective edge devices.

*a5) Client-side model backward pass:* Each edge device updates its client-side sub-model based on the received activations' gradients. For the  $i$ -th edge device, the client-side sub-model is updated through

$$\mathbf{w}_{c,i}^t \leftarrow \mathbf{w}_{c,i}^{t-1} - \gamma \nabla_{\mathbf{w}_c} F_i(\mathbf{w}_{c,i}^{t-1}; \xi_i^t), \quad (6)$$

where  $\nabla_{\mathbf{w}_c} F_i(\mathbf{w}_{c,i}^{t-1}; \xi_i^t)$  represent the stochastic gradients of client-side sub-models of the  $i$ -th edge device.

*b. The client-side model aggregation stage:* The client-side model aggregation stage aims to aggregate forged client-specific sub-models (including client-side sub-models and server-side non-common sub-models specific to each client) on the fed server. This stage is executed every  $I$  training rounds and comprises the following three steps.

*b1) Sub-model uploading:* In this step, edge devices and the edge server simultaneously upload the client-side sub-models and corresponding server-side non-common sub-models to the fed server, serving as the basis for the subsequent forged client-specific sub-model aggregation.

<sup>4</sup>The edge server can update the model for multiple edge devices in either a serial or parallel manner, which does not impact training performance as aggregation occurs every round. Here, we formulate the parallel manner.

---

**Algorithm 1** The HASFL Training Framework
 

---

**Input:**  $I, \gamma, E, \mathcal{N}$  and  $\mathcal{D}$ .

**Output:**  $\mathbf{w}^*$ .

```

1: Initialization:  $\mathbf{w}^0, b^0, \mathbf{w}_i^0 \leftarrow \mathbf{w}^0, b_i^0 \leftarrow b^0, \tau \leftarrow 0$ , and  $\rho \leftarrow 0$ .
2: while  $\sum_{\eta=0}^{\tau} \sum_{i=1}^N b_i^\eta \leq E$  do
3:
4:   /** Runs on edge devices **/
5:   for all edge device  $i \in \mathcal{N}$  in parallel do
6:      $\mathbf{a}_i^\tau \leftarrow \varphi(\mathbf{x}_i^\tau; \mathbf{w}_{c,i}^{\tau-1})$ 
7:     Send  $(\mathbf{a}_i^\tau, \mathbf{y}_i^\tau)$  to the edge server
8:   end for
9:
10:  /** Runs on edge server **/
11:   $\hat{\mathbf{y}}_i^\tau = \varphi(\mathbf{a}_i^\tau; \mathbf{w}_{s,i}^{\tau-1})$ 
12:  Calculate loss function value  $f(\mathbf{w}^{\tau-1})$ 
13:   $\mathbf{h}_s^\tau \leftarrow \mathbf{h}_s^{\tau-1} - \gamma \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_s} F_i(\mathbf{h}_s^{\tau-1}; \xi_i^\tau)$ ,
14:   $\mathbf{h}_{m,i}^\tau \leftarrow \mathbf{h}_{m,i}^{\tau-1} - \gamma \nabla_{\mathbf{h}_m} F_i(\mathbf{h}_{m,i}^{\tau-1}; \xi_i^\tau)$ 
15:  Send activations' gradients to corresponding edge devices
16:
17:  /** Runs on edge devices **/
18:  for all edge device  $i \in \mathcal{N}$  in parallel do
19:     $\mathbf{w}_{c,i}^\tau \leftarrow \mathbf{w}_{c,i}^{\tau-1} - \gamma \nabla_{\mathbf{w}_c} F_i(\mathbf{w}_{c,i}^{\tau-1}; \xi_i^\tau)$ 
20:  end for
21:
22:  /** Runs on the fed server **/
23:  if  $\tau \bmod I = 0$  then
24:    Determine  $\mathbf{b}^\tau = [b_1^\tau, b_2^\tau, \dots, b_N^\tau]$  and  $\boldsymbol{\mu}^\tau$  based on Algo-
rithm 2
25:    Forge client-side sub-models  $\mathbf{h}_{c,i}^\tau = [\mathbf{h}_{m,i}^\tau; \mathbf{w}_{c,i}^\tau]$ 
26:     $\mathbf{h}_c^\tau = \frac{1}{N} \sum_{i=1}^N \mathbf{h}_{c,i}^\tau$ 
27:     $\mathbf{h}_{c,i}^\tau \leftarrow \mathbf{h}_c^\tau$ 
28:  else if  $\tau \bmod I \neq 0$  then
29:    Determine  $\mathbf{b}^\tau = [b_1^\tau, b_2^\tau, \dots, b_N^\tau]$  based on Theorem 1
30:  end if
31:
32:  Update  $\tau \leftarrow \tau + 1$ 
33: end while
    
```

---

*b2) Client-side model aggregation:* The fed server pairs and assembles the received client-side sub-models and server-side non-common sub-models to forge the client-specific models. Then, these forged client-specific models are aggregated across all participating devices, as given by

$$\mathbf{h}_c^t = \frac{1}{N} \sum_{i=1}^N \mathbf{h}_{c,i}^t, \quad (7)$$

where  $\mathbf{h}_{c,i}^t = [\mathbf{h}_{m,i}^t; \mathbf{w}_{c,i}^t]$  is the forged client-specific model of the  $i$ -th edge device.

*b3) Sub-model downloading:* After the client-side model aggregation is completed, the fed server transmits updated client-side sub-models and server-side non-common sub-models back to corresponding edge devices and the edge server, respectively.

The workflow of the HASFL training framework is presented in **Algorithm 1**, where  $E$  denotes the number of stochastic gradients involved throughout model training.

#### IV. CONVERGENCE ANALYSIS OF HASFL

This section presents the convergence analysis of SFL, characterizing the effect of BS and MS on training convergence.

This analysis serves as the theoretical foundation for designing an efficient iterative optimization algorithm in Section VI.

We define the global model of the  $i$ -th edge device at the  $t$ -th training round as  $\mathbf{w}_i^t = [\mathbf{h}_{s,i}^t; \mathbf{h}_{c,i}^t]$ , which comprises server-side common sub-model  $\mathbf{h}_{s,i}^t$  and forged client-specific sub-model  $\mathbf{h}_{c,i}^t$ . The server-side common sub-model is updated and synchronized across all devices every training round, whereas the forged client-specific sub-models are aggregated every  $I$  training rounds.

Due to non-uniform model partitioning, different edge devices retain varying numbers of layers on the client side. Since the client-side and server-side sub-models follow different aggregation schedules, with only the client-side sub-models being aggregated periodically every  $I$  rounds, we define the model split point  $L_c$  as the maximum depth (i.e., number of layers) of client-specific sub-models among all participating edge devices. The gradients of forged client-specific and server-side common sub-models are given by

$$\mathbf{g}_{c,i}^t = [\nabla_{\mathbf{h}_m} F_i(\mathbf{h}_{m,i}^{t-1}; \xi_i^t); \nabla_{\mathbf{w}_c} F_i(\mathbf{w}_{c,i}^{t-1}; \xi_i^t)] \quad (8)$$

and

$$\mathbf{g}_{s,i}^t = \sum_{i=1}^N \nabla_{\mathbf{h}_s} F_i(\mathbf{h}_s^{t-1}; \xi_i^t), \quad (9)$$

where  $\mathbf{h}_s^t = \frac{1}{N} \sum_{i=1}^N \mathbf{h}_{s,i}^t$ . Therefore, the gradient of the global model of the  $i$ -th edge device is  $\mathbf{g}_i^t = [\mathbf{g}_{s,i}^t; \mathbf{g}_{c,i}^t]$ .

In line with the seminal studies in distributed stochastic optimization [37]–[40] that analyze convergence bound on the aggregated version of individual solutions, we conduct the convergence analysis of  $\mathbf{w}^t = \frac{1}{N} \sum_{i=1}^N \mathbf{w}_i^t$ . To analyze the convergence upper bound of HASFL, we consider the following two standard assumptions on loss functions [37]–[40]:

**Assumption 1** (*Smoothness of the local loss functions*). Each local loss function  $f_i(\mathbf{w})$  is differentiable and  $\beta$ -smooth, i.e., for any  $\mathbf{w}$  and  $\mathbf{w}'$ , we have

$$\|\nabla_{\mathbf{w}} f_i(\mathbf{w}) - \nabla_{\mathbf{w}} f_i(\mathbf{w}')\| \leq \beta \|\mathbf{w} - \mathbf{w}'\|, \quad \forall i. \quad (10)$$

**Assumption 2** (*Bounded variance and second-order moments of stochastic gradients*). The stochastic gradients computed via mini-batch sampling have bounded variance and second-order moments for each layer:

$$\mathbb{E}_{\xi_i \sim \mathcal{D}_i} \|\nabla_{\mathbf{w}} F_i(\mathbf{w}; \xi_i) - \nabla_{\mathbf{w}} f_i(\mathbf{w})\|^2 \leq \sum_{j=1}^l \frac{\sigma_j^2}{b}, \quad \forall \mathbf{w}, \quad \forall i, \quad (11)$$

$$\mathbb{E}_{\xi_i \sim \mathcal{D}_i} \|\nabla_{\mathbf{w}} F_i(\mathbf{w}; \xi_i)\|^2 \leq \sum_{j=1}^l G_j^2, \quad \forall \mathbf{w}, \quad \forall i, \quad (12)$$

where  $b$  is the mini-batch size,  $l$  denotes the number of layers for model  $\mathbf{w}$ ,  $\sigma_j^2$  is the constant,  $\frac{\sigma_j^2}{b}$  and  $G_j^2$  represent the bounded variance and second-order moments for the  $j$ -th layer of model  $\mathbf{w}$ , respectively.

**Lemma 1.** *Under Assumption 1, Algorithm 1 ensures*

$$\mathbb{E}[\|\mathbf{h}_c^t - \mathbf{h}_{c,i}^t\|^2] \leq \mathbb{1}_{\{I>1\}} 4\gamma^2 I^2 \sum_{j=1}^{L_c} G_j^2, \forall i, \forall t, \quad (13)$$

where  $I$  denotes the client-side model aggregation interval, and  $\mathbf{h}_c^t$  is defined in Eqn. (7) as the virtual aggregated version of client-side model at the  $t$ -th training round.

*Proof.* Consider an arbitrary training round  $t \geq 1$ . Let  $t_0 \leq t$  be the most recent round at which client-side model aggregation took place, i.e., the largest  $t_0$  satisfying  $t_0 \bmod I = 0$  (Such  $t_0$  must exist and  $t - t_0 \leq I$ ). Recalling Eqn. (5), Eqn. (6) and Eqn. (8), the forged client-specific sub-model of the  $i$ -th edge device can be expressed as

$$\mathbf{h}_{c,i}^t = \mathbf{h}_c^{t_0} - \gamma \sum_{\tau=t_0+1}^t \mathbf{g}_{c,i}^\tau. \quad (14)$$

By (7), we have

$$\mathbf{h}_c^t = \mathbf{h}_c^{t_0} - \gamma \sum_{\tau=t_0+1}^t \frac{1}{N} \sum_{i=1}^N \mathbf{g}_{c,i}^\tau. \quad (15)$$

Subtracting Eqn. (15) from Eqn. (14) and taking the squared norm yields:

$$\begin{aligned} & \mathbb{E}[\|\mathbf{h}_c^t - \mathbf{h}_{c,i}^t\|^2] \\ &= \mathbb{1}_{\{I>1\}} \mathbb{E}[\|\gamma \sum_{\tau=t_0+1}^t \frac{1}{N} \sum_{i=1}^N \mathbf{g}_{c,i}^\tau - \gamma \sum_{\tau=t_0+1}^t \mathbf{g}_{c,i}^\tau\|^2] \\ &= \mathbb{1}_{\{I>1\}} \gamma^2 \mathbb{E}[\|\sum_{\tau=t_0+1}^t \frac{1}{N} \sum_{i=1}^N \mathbf{g}_{c,i}^\tau - \sum_{\tau=t_0+1}^t \mathbf{g}_{c,i}^\tau\|^2] \\ &\stackrel{(a)}{\leq} \mathbb{1}_{\{I>1\}} 2\gamma^2 \mathbb{E}[\|\sum_{\tau=t_0+1}^t \frac{1}{N} \sum_{i=1}^N \mathbf{g}_{c,i}^\tau\|^2 + \|\sum_{\tau=t_0+1}^t \mathbf{g}_{c,i}^\tau\|^2] \\ &\stackrel{(b)}{\leq} \mathbb{1}_{\{I>1\}} 2\gamma^2 (t - t_0) \mathbb{E}[\sum_{\tau=t_0+1}^t \|\frac{1}{N} \sum_{i=1}^N \mathbf{g}_{c,i}^\tau\|^2 + \sum_{\tau=t_0+1}^t \|\mathbf{g}_{c,i}^\tau\|^2] \\ &\stackrel{(c)}{\leq} \mathbb{1}_{\{I>1\}} 2\gamma^2 (t - t_0) \mathbb{E}[\sum_{\tau=t_0+1}^t (\frac{1}{N} \sum_{i=1}^N \|\mathbf{g}_{c,i}^\tau\|^2) + \sum_{\tau=t_0+1}^t \|\mathbf{g}_{c,i}^\tau\|^2] \\ &\stackrel{(d)}{\leq} \mathbb{1}_{\{I>1\}} 4\gamma^2 I^2 \sum_{j=1}^{L_c} G_j^2, \end{aligned}$$

where (a)-(c) follow by the inequality  $\|\sum_{i=1}^n \mathbf{z}_i\|^2 \leq n \sum_{i=1}^n \|\mathbf{z}_i\|^2$  for any vectors  $\mathbf{z}_i$  and any positive integer  $n$  (using  $n = 2$  in (a),  $n = t - t_0$  in (b), and  $n = N$  in (c)); and (d) follows from **Assumption 2**.  $\square$

**Theorem 1.** *Under Assumption 1 and Assumption 2, if  $0 < \gamma \leq \frac{1}{\beta}$  in Algorithm 1, then for any total number of training rounds  $R \geq 1$ , we have*

$$\begin{aligned} & \frac{1}{R} \sum_{t=1}^R \mathbb{E}[\|\nabla_{\mathbf{w}} f(\mathbf{w}^{t-1})\|^2] \\ & \leq \frac{2\vartheta}{\gamma R} + \frac{\beta\gamma \sum_{i=1}^N \sum_{j=1}^L \frac{\sigma_j^2}{b_i}}{N^2} + \mathbb{1}_{\{I>1\}} 4\beta^2 \gamma^2 I^2 \sum_{j=1}^{L_c} G_j^2, \end{aligned} \quad (16)$$

where  $\vartheta = f(\mathbf{w}^0) - f^*$ ,  $b_i$  denotes the batch size of the  $i$ -th edge device,  $L$  and  $f^*$  represent the total number of global model layers and global minimum value of problem (1), respectively.

*Proof.* We begin by analyzing the expected decrease in the global loss function across one training round. By invoking the smoothness of the loss function  $f(\cdot)$ , we have

$$\begin{aligned} \mathbb{E}[f(\mathbf{w}^t)] &\leq \mathbb{E}[f(\mathbf{w}^{t-1})] + \mathbb{E}[\langle \nabla_{\mathbf{w}} f(\mathbf{w}^{t-1}), \mathbf{w}^t - \mathbf{w}^{t-1} \rangle] \\ &\quad + \frac{\beta}{2} \mathbb{E}[\|\mathbf{w}^t - \mathbf{w}^{t-1}\|^2]. \end{aligned} \quad (17)$$

Note that

$$\begin{aligned} & \mathbb{E}[\|\mathbf{w}^t - \mathbf{w}^{t-1}\|^2] \\ &= \mathbb{E}[\|[\mathbf{h}_c^t; \mathbf{h}_s^t] - [\mathbf{h}_c^{t-1}; \mathbf{h}_s^{t-1}]\|^2] \\ &= \mathbb{E}[\|[\mathbf{h}_c^t - \mathbf{h}_c^{t-1}; \mathbf{h}_s^t - \mathbf{h}_s^{t-1}]\|^2] \\ &= \mathbb{E}[\|\mathbf{h}_c^t - \mathbf{h}_c^{t-1}\|^2] + \mathbb{E}[\|\mathbf{h}_s^t - \mathbf{h}_s^{t-1}\|^2], \end{aligned} \quad (18)$$

where  $\mathbb{E}[\|\mathbf{h}_c^t - \mathbf{h}_c^{t-1}\|^2]$  can be bounded, as given by

$$\begin{aligned} & \mathbb{E}[\|\mathbf{h}_c^t - \mathbf{h}_c^{t-1}\|^2] \stackrel{(a)}{=} \gamma^2 \mathbb{E}[\|\frac{1}{N} \sum_{i=1}^N \mathbf{g}_{c,i}^t\|^2] \\ & \stackrel{(b)}{=} \gamma^2 \mathbb{E}[\|\frac{1}{N} \sum_{i=1}^N (\mathbf{g}_{c,i}^t - \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1}))\|^2] \\ & \quad + \gamma^2 \mathbb{E}[\|\frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1})\|^2] \\ & \stackrel{(c)}{=} \frac{\gamma^2}{N^2} \sum_{i=1}^N \mathbb{E}[\|\mathbf{g}_{c,i}^t - \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1})\|^2] \\ & \quad + \gamma^2 \mathbb{E}[\|\frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1})\|^2] \\ & \stackrel{(d)}{\leq} \frac{\gamma^2}{N^2} \sum_{i=1}^N \sum_{j=1}^{L_c} \frac{\sigma_j^2}{b_i} + \gamma^2 \mathbb{E}[\|\frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1})\|^2], \end{aligned} \quad (19)$$

where (a) follows from Eqn. (7) and Eqn. (14); (b) follows that  $\mathbb{E}[\mathbf{g}_{c,i}^t] = \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1})$  and  $\mathbb{E}[\|\mathbf{z}\|^2] = \mathbb{E}[\|\mathbf{z} - \mathbb{E}[\mathbf{z}]\|^2] + \|\mathbb{E}[\mathbf{z}]\|^2$  that holds for any random vector  $\mathbf{z}$ ; (c) is because  $\mathbf{g}_{c,i}^t - \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1})$  has zero mean and is independent between edge devices; and (d) follows from **Assumption 2**.

Similarly,  $\mathbb{E}[\|\mathbf{h}_s^t - \mathbf{h}_s^{t-1}\|^2]$  has an upper bound:

$$\begin{aligned} & \mathbb{E}[\|\mathbf{h}_s^t - \mathbf{h}_s^{t-1}\|^2] \\ & \leq \frac{\gamma^2}{N^2} \sum_{i=L_c+1}^N \sum_{j=1}^L \frac{\sigma_j^2}{b_i} + \gamma^2 \mathbb{E}[\|\frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_s} f_i(\mathbf{h}_{s,i}^{t-1})\|^2]. \end{aligned} \quad (20)$$

Substituting Eqn. (19) and Eqn. (20) into Eqn. (18) yields

$$\begin{aligned} & \mathbb{E}[\|\mathbf{w}^t - \mathbf{w}^{t-1}\|^2] \\ & \leq \frac{\gamma^2}{N^2} \sum_{i=1}^N \sum_{j=1}^L \frac{\sigma_j^2}{b_i} + \gamma^2 \mathbb{E}[\|\frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1})\|^2] \\ & \quad + \gamma^2 \mathbb{E}[\|\frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_s} f_i(\mathbf{h}_{s,i}^{t-1})\|^2] \end{aligned}$$

$$= \frac{\gamma^2}{N^2} \sum_{i=1}^N \sum_{j=1}^L \frac{\sigma_j^2}{b_i} + \gamma^2 \mathbb{E}[\|\frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} f_i(\mathbf{w}_i^{t-1})\|^2]. \quad (21)$$

We further note that

$$\begin{aligned} & \mathbb{E}[\langle \nabla_{\mathbf{w}} f(\mathbf{w}^{t-1}), \mathbf{w}^t - \mathbf{w}^{t-1} \rangle] \\ & \stackrel{(a)}{=} -\gamma \mathbb{E}[\langle \nabla_{\mathbf{w}} f(\mathbf{w}^{t-1}), \frac{1}{N} \sum_{i=1}^N \mathbf{g}_i^t \rangle] \\ & \stackrel{(b)}{=} -\gamma \mathbb{E}[\langle \nabla_{\mathbf{w}} f(\mathbf{w}^{t-1}), \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} f_i(\mathbf{w}_i^{t-1}) \rangle] \\ & \stackrel{(c)}{=} -\frac{\gamma}{2} \mathbb{E}[\|\nabla_{\mathbf{w}} f(\mathbf{w}^{t-1})\|^2 + \|\frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} f_i(\mathbf{w}_i^{t-1})\|^2 \\ & \quad - \|\nabla_{\mathbf{w}} f(\mathbf{w}^{t-1}) - \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} f_i(\mathbf{w}_i^{t-1})\|^2], \quad (22) \end{aligned}$$

where (a) follows from  $\mathbf{w}^t = \frac{1}{N} \sum_{i=1}^N \mathbf{w}_i^t$ ; (c) follows from the identity  $\langle \mathbf{z}_1, \mathbf{z}_2 \rangle = \frac{1}{2} (\|\mathbf{z}_1\|^2 + \|\mathbf{z}_2\|^2 - \|\mathbf{z}_1 - \mathbf{z}_2\|^2)$  for any two vectors  $\mathbf{z}_1, \mathbf{z}_2$  of the same length; (b) follows from

$$\begin{aligned} & \mathbb{E}[\langle \nabla_{\mathbf{w}} f(\mathbf{w}^{t-1}), \frac{1}{N} \sum_{i=1}^N \mathbf{g}_i^t \rangle] \\ & = \mathbb{E}[\mathbb{E}[\langle \nabla_{\mathbf{w}} f(\mathbf{w}^{t-1}), \frac{1}{N} \sum_{i=1}^N \mathbf{g}_i^t | \boldsymbol{\xi}^{[t-1]} \rangle]] \\ & = \mathbb{E}[\langle \nabla_{\mathbf{w}} f(\mathbf{w}^{t-1}), \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{g}_i^t | \boldsymbol{\xi}^{[t-1]}] \rangle] \\ & = \mathbb{E}[\langle \nabla_{\mathbf{w}} f(\mathbf{w}^{t-1}), \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} f_i(\mathbf{w}_i^{t-1}) \rangle], \end{aligned}$$

where the first equality follows by the law of expectations, the second equality holds because  $\mathbf{w}^{t-1}$  is determined by  $\boldsymbol{\xi}^{[t-1]} = [\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^{t-1}]$ , and the third equality follows from  $\mathbb{E}[\mathbf{g}_i^t | \boldsymbol{\xi}^{[t-1]}] = \mathbb{E}[\nabla F_i(\mathbf{w}_i^{t-1}; \boldsymbol{\xi}_i^t) | \boldsymbol{\xi}^{[t-1]}] = \nabla f_i(\mathbf{w}_i^{t-1})$ .

Substituting Eqn. (21) and Eqn. (22) into Eqn. (17), we have

$$\begin{aligned} & \mathbb{E}[f(\mathbf{w}^t)] \\ & \leq \mathbb{E}[f(\mathbf{w}^{t-1})] - \frac{\gamma - \gamma^2 \beta}{2} \mathbb{E}[\|\frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} f_i(\mathbf{w}_i^{t-1})\|^2] \\ & \quad - \frac{\gamma}{2} \mathbb{E}[\|\nabla_{\mathbf{w}} f(\mathbf{w}^{t-1})\|^2] + \frac{\beta \gamma^2}{2N^2} \sum_{i=1}^N \sum_{j=1}^L \frac{\sigma_j^2}{b_i} \\ & \quad + \frac{\gamma}{2} \mathbb{E}[\|\nabla_{\mathbf{w}} f(\mathbf{w}^{t-1}) - \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} f_i(\mathbf{w}_i^{t-1})\|^2] \\ & \stackrel{(a)}{\leq} \mathbb{E}[f(\mathbf{w}^{t-1})] - \frac{\gamma}{2} \mathbb{E}[\|\nabla_{\mathbf{w}} f(\mathbf{w}^{t-1})\|^2] + \frac{\beta \gamma^2}{2N^2} \sum_{i=1}^N \sum_{j=1}^L \frac{\sigma_j^2}{b_i} \\ & \quad + \frac{\gamma}{2} \mathbb{E}[\|\nabla_{\mathbf{h}_c} f(\mathbf{h}_c^{t-1}) - \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1})\|^2] \\ & \quad + \frac{\gamma}{2} \mathbb{E}[\|\nabla_{\mathbf{h}_s} f(\mathbf{h}_s^{t-1}) - \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_s} f_i(\mathbf{h}_{s,i}^{t-1})\|^2] \end{aligned}$$

$$\begin{aligned} & \stackrel{(b)}{\leq} \mathbb{E}[f(\mathbf{w}^{t-1})] - \frac{\gamma}{2} \mathbb{E}[\|\nabla_{\mathbf{w}} f(\mathbf{w}^{t-1})\|^2] + \frac{\beta \gamma^2}{2N^2} \sum_{i=1}^N \sum_{j=1}^L \frac{\sigma_j^2}{b_i} \\ & \quad + \mathbb{1}_{\{I>1\}} 2\beta^2 \gamma^3 I^2 \sum_{j=1}^{L_c} G_j^2, \quad (23) \end{aligned}$$

where (a) follows from  $0 < \gamma \leq \frac{1}{\beta}$  and (b) holds because of the following inequality (24) and (25)

$$\begin{aligned} & \mathbb{E}[\|\nabla_{\mathbf{h}_c} f(\mathbf{h}_c^{t-1}) - \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1})\|^2] \\ & = \mathbb{E}[\|\frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_c^{t-1}) - \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1})\|^2] \\ & = \frac{1}{N^2} \mathbb{E}[\|\sum_{i=1}^N (\nabla_{\mathbf{h}_c} f_i(\mathbf{h}_c^{t-1}) - \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1}))\|^2] \\ & \leq \frac{1}{N} \mathbb{E}[\sum_{i=1}^N \|\nabla_{\mathbf{h}_c} f_i(\mathbf{h}_c^{t-1}) - \nabla_{\mathbf{h}_c} f_i(\mathbf{h}_{c,i}^{t-1})\|^2] \\ & \leq \beta^2 \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\|\mathbf{h}_c^{t-1} - \mathbf{h}_{c,i}^{t-1}\|^2] \\ & \leq \mathbb{1}_{\{I>1\}} 4\beta^2 \gamma^2 I^2 \sum_{j=1}^{L_c} G_j^2, \quad (24) \end{aligned}$$

where the first inequality follows from  $\|\sum_{i=1}^N \mathbf{z}_i\|^2 \leq N \sum_{i=1}^N \|\mathbf{z}_i\|^2$  for any vectors  $\mathbf{z}_i$ ; the second inequality follows from the smoothness of  $f_i$  under **Assumption 1**; and the third inequality follows from **Lemma 1**. Moreover, we have

$$\begin{aligned} & \mathbb{E}[\|\nabla_{\mathbf{h}_s} f(\mathbf{h}_s^{t-1}) - \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{h}_s} f_i(\mathbf{h}_{s,i}^{t-1})\|^2] \\ & \leq \beta^2 \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\|\mathbf{h}_s^{t-1} - \mathbf{h}_{s,i}^{t-1}\|^2] \stackrel{(a)}{=} 0, \quad (25) \end{aligned}$$

where (a) holds because the server-side common sub-models are aggregated in each training round (i.e.,  $I = 1$ ). Therefore, at any training round  $t$ , the server-side common sub-model of each edge device is the aggregated version.

Dividing both sides of Eqn. (23) by  $\frac{\gamma}{2}$  and rearranging terms yields

$$\begin{aligned} & \mathbb{E}[\|\nabla_{\mathbf{w}} f(\mathbf{w}^{t-1})\|^2] \\ & \leq \frac{2}{\gamma} (\mathbb{E}[f(\mathbf{w}^{t-1})] - \mathbb{E}[f(\mathbf{w}^t)]) + \frac{\beta \gamma \sum_{i=1}^N \sum_{j=1}^L \frac{\sigma_j^2}{b_i}}{N^2} + \mathbb{1}_{\{I>1\}} 4\beta^2 \gamma^2 I^2 \sum_{j=1}^{L_c} G_j^2. \end{aligned}$$

Summing over  $t \in \{1, \dots, R\}$  and dividing both sides by  $R$  yields

$$\begin{aligned} & \frac{1}{R} \sum_{t=1}^R \mathbb{E}[\|\nabla_{\mathbf{w}} f(\mathbf{w}^{t-1})\|^2] \\ & \leq \frac{2}{\gamma R} (f(\mathbf{w}^0) - \mathbb{E}[f(\mathbf{w}^R)]) + \frac{\beta \gamma \sum_{i=1}^N \sum_{j=1}^L \frac{\sigma_j^2}{b_i}}{N^2} + \mathbb{1}_{\{I>1\}} 4\beta^2 \gamma^2 I^2 \sum_{j=1}^{L_c} G_j^2 \\ & \stackrel{(a)}{\leq} \frac{2}{\gamma R} (f(\mathbf{w}^0) - f^*) + \frac{\beta \gamma \sum_{i=1}^N \sum_{j=1}^L \frac{\sigma_j^2}{b_i}}{N^2} + \mathbb{1}_{\{I>1\}} 4\beta^2 \gamma^2 I^2 \sum_{j=1}^{L_c} G_j^2, \end{aligned}$$



where (a) is because  $f^*$  is the global minimum value of problem (1).  $\square$

Substituting Eqn. (16) into Eqn. (26) yields **Corollary 1**, revealing a lower bound on the number of training rounds for achieving target convergence accuracy.

**Corollary 1.** *The number  $R$  of training rounds for achieving target convergence accuracy  $\varepsilon$ , i.e., satisfying*

$$\frac{1}{R} \sum_{t=1}^R \mathbb{E}[\|\nabla_{\mathbf{w}} f(\mathbf{w}^{t-1})\|^2] \leq \varepsilon, \quad (26)$$

is lower bounded by

$$R \geq \frac{2\vartheta}{\gamma \left( \varepsilon - \frac{\beta\gamma \sum_{i=1}^N \sum_{j=1}^L \frac{\sigma_j^2}{b_i}}{N^2} - \mathbb{1}_{\{I>1\}} 4\beta^2\gamma^2 I^2 \sum_{j=1}^{L_c} G_j^2 \right)}. \quad (27)$$

**Insight 1:** Eqn. (27) reveals that the number  $R$  of training rounds for achieving target convergence accuracy  $\varepsilon$  decreases with the increase of  $b_i$ . This is because a larger  $b_i$  reduces the variance of stochastic gradients, thereby enhancing the training convergence. For a fixed number of training rounds  $R$ , increasing  $b_i$  leads to higher converged accuracy, i.e., smaller  $\varepsilon$ . It is interesting to note that the BSs of clients can compensate for each other: A stronger client can take a larger BS, whereas a weaker client takes a smaller BS, thereby mitigating the straggler effect without affecting the convergence upper bound.

**Insight 2:** Eqn. (27) also shows that  $L_c$  plays a critical role in model convergence when the client-side aggregation interval  $I > 1$ . A smaller  $L_c$  (i.e., placing more layers on the edge server) enables more frequent aggregations of a larger portion of the model for expediting convergence, whereas a larger  $L_c$  retains more layers on the edge devices, slowing down model convergence. Conversely, when  $I = 1$ ,  $L_c$  has no impact on convergence since all layers are aggregated synchronously.

These observations align with the experimental results described in Fig. 2 and Fig. 3. While a larger BS can expedite training convergence by reducing gradient variance, it also incurs higher per-round computation and communication costs, which may prolong training latency on resource-constrained edge devices. On the other hand, MS determines not only the computing workload between edge devices and the server, but also the volume of smashed data exchanged. Therefore, optimizing BS and MS, i.e., assigning different BS and MS to diverse edge devices, is essential for accelerating SFL at heterogeneous edge devices.

## V. PROBLEM FORMULATION

The convergence analysis quantifies the impact of BS and MS on model convergence. Building upon the convergence upper bound in Section IV, we formulate a joint BS and MS optimization problem to minimize the training latency of achieving convergence for HASFL over edge networks with heterogeneous participating edge devices. Then, we propose a heterogeneity-aware BS and MS strategy in Section VI. For clarity, the decision variables and definitions are listed below.

- $\mathbf{b}$ :  $b_i \in \mathbb{N}^+$  represents the BS decision, which indicates the number of data samples utilized by the  $i$ -th edge device for one training round.  $\mathbf{b} = [b_1, b_2, \dots, b_N]$  denotes the collection of BS decisions.
- $\boldsymbol{\mu}$ :  $\mu_{i,j} \in \{0, 1\}$  is the MS decision, where  $\mu_{i,j} = 1$  indicates that the  $j$ -th neural network layer is selected as the cut layer for the  $i$ -th edge device, and  $\mu_{i,j} = 0$ , otherwise.  $\boldsymbol{\mu} = [\mu_{1,1}, \mu_{1,2}, \dots, \mu_{N,L}]$  represents the collection of MS decisions.

### A. Training Latency Analysis

This section presents a detailed analysis for the training latency of HASFL. Without loss of generality, we focus on one training round for analysis. For notational brevity, the training round number index  $t$  is omitted. We start by analyzing the latency of the split training stage for one training round as follows.

*a1) Client-side model forward propagation latency:* In this step, each edge device performs the client-side FP with a mini-batch sampled from its local dataset. The computing workload (in float-point operations (FLOPs) [21], [31], [41]) of the client-side FP for the  $i$ -th edge device per data sample is represented as  $\Phi_{c,i}^F(\boldsymbol{\mu}) = \sum_{j=1}^L \mu_{i,j} \rho_j$ , where  $\rho_j$  denotes the FP computing workload of propagating the first  $j$  layer of neural network per data sample. The  $i$ -th edge device utilizes a mini-batch containing  $b_i$  data samples to execute the client-side FP. The client-side FP latency of the  $i$ -th edge device is given by [20], [31]

$$T_i^F = \frac{b_i \Phi_{c,i}^F(\boldsymbol{\mu})}{f_i}, \quad \forall i \in \mathcal{N}, \quad (28)$$

where  $f_i$  denotes the computing capability (in float-point operations per second (FLOPS)) of the  $i$ -th edge device.

*a2) Activation uploading latency:* After the client-side FP is completed, each edge device transmits the activations to the edge server. Let  $\Gamma_{a,i}(\boldsymbol{\mu}) = \sum_{j=1}^L \mu_{i,j} \psi_j$  represent the data size (in bits) of activations for the  $i$ -th edge device, where  $\psi_j$  denotes the data size of activations at the cut layer  $j$ . For the  $i$ -th edge device, the activation uploading latency is given by

$$T_{a,i}^U = \frac{b_i \Gamma_{a,i}(\boldsymbol{\mu})}{r_i^U}, \quad \forall i \in \mathcal{N}, \quad (29)$$

where  $r_i^U$  represents the uplink data rate from the  $i$ -th edge device to the edge server.

*a3) Server-side model forward propagation and backward pass latency:* In this step, the edge server performs server-side FP and BP with activations received from all participating edge devices. Let  $\Phi_s^F(\mathbf{b}, \boldsymbol{\mu}) = \sum_{i=1}^N \sum_{j=1}^L b_i \mu_{i,j} (\rho_L - \rho_j)$  and  $\Phi_s^B(\mathbf{b}, \boldsymbol{\mu}) = \sum_{i=1}^N \sum_{j=1}^L b_i \mu_{i,j} (\varpi_L - \varpi_j)$  denote the computing workload of the server-side FP and BP, respectively, where  $\varpi_j$  denotes the BP computing workload of propagating the first  $j$



layer of neural network per data sample. Thus, the server-side model FP and BP latency can be determined as

$$T_s^F = \frac{\Phi_s^F(\mathbf{b}, \boldsymbol{\mu})}{f_s} \quad (30)$$

and

$$T_s^B = \frac{\Phi_s^B(\mathbf{b}, \boldsymbol{\mu})}{f_s}, \quad (31)$$

where  $f_s$  is the computing capability of the edge server.

*a4) Downloading latency of activations' gradients:* After the completion of server-side BP, the edge server sends the activations' gradients back to the respective edge devices. Let  $\Gamma_{g,i}(\boldsymbol{\mu}) = \sum_{j=1}^L \mu_{i,j} \chi_j$  represent the data size of activations' gradients for the  $i$ -th edge device per data sample, where  $\chi_j$  is the data size of activations' gradients at the cut layer  $j$ . Downloading latency of activations' gradients of the  $i$ -th edge device can be calculated as

$$T_{g,i}^D = \frac{b_i \Gamma_{g,i}(\boldsymbol{\mu})}{r_i^D}, \quad \forall i \in \mathcal{N}, \quad (32)$$

where  $r_i^D$  is the downlink data rate from edge server to the  $i$ -th edge device.

*a5) Client-side model backward pass latency:* In this step, each edge device executes client-side BP based on the received activations' gradients. Let  $\Phi_{c,i}^B(\boldsymbol{\mu}) = \sum_{j=1}^L \mu_{i,j} \varpi_j$  represent the computing workload of the client-side BP for the  $i$ -th edge device per data sample. For the  $i$ -th edge device, the client-side BP latency can be obtained as

$$T_i^B = \frac{b_i \Phi_{c,i}^B(\boldsymbol{\mu})}{f_i}, \quad \forall i \in \mathcal{N}. \quad (33)$$

Next, we analyze the latency of the client-side model aggregation stage every  $I$  training rounds as follows.

*b1) Sub-model uploading latency:* In this step, each edge device sends its client-side sub-model to the fed server, while the edge server uploads the server-side non-common sub-models to the fed server. Let  $\Lambda_{c,i}(\boldsymbol{\mu}) = \sum_{j=1}^L \mu_{i,j} \delta_j$  denote the data size of client-side sub-model for the  $i$ -th edge device, where  $\delta_j$  is the data size of client-side sub-model with the cut layer  $j$ . The total data size of the exchanged server-side non-common sub-models is denoted by  $\Lambda_s(\boldsymbol{\mu}) = N \max_i \left\{ \sum_{j=1}^L \mu_{i,j} \delta_j \right\} - \sum_{i=1}^N \sum_{j=1}^L \mu_{i,j} \delta_j$ . The uploading latency of client-side sub-model for the  $i$ -th edge device and server-side non-common sub-models are expressed as

$$T_{c,i}^U = \frac{\Lambda_{c,i}(\boldsymbol{\mu})}{r_{i,f}^U}, \quad \forall i \in \mathcal{N} \quad (34)$$

and

$$T_s^U = \frac{\Lambda_s(\boldsymbol{\mu})}{r_{s,f}}, \quad (35)$$

where  $r_{i,f}^U$  and  $r_{s,f}$  denote the uplink data rate for transferring the sub-model from the  $i$ -th edge device and edge server to the fed server, respectively.

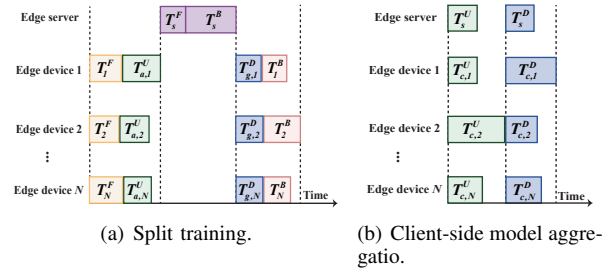


Fig. 4. An illustration of split training and client-side model aggregation stages.

*b2) Client-side model aggregation:* The fed server assembles the received client-side sub-models and server-side non-common sub-models into client-specific models and then aggregates them. For simplicity, the sub-model aggregation latency for this part is ignored, as it is negligible compared to the latency from other stages [42], [43].

*b3) Sub-model downloading latency:* After completing the client-side model aggregation, the fed server sends the updated client-side sub-models to the respective edge devices, along with the server-side non-common sub-models to the edge server. The downloading latency of the client-side sub-model for the  $i$ -th edge device and server-side non-common sub-models are given by

$$T_{c,i}^D = \frac{\Lambda_{c,i}(\boldsymbol{\mu})}{r_{i,f}^D}, \quad \forall i \in \mathcal{N} \quad (36)$$

and

$$T_s^D = \frac{\Lambda_s(\boldsymbol{\mu})}{r_{f,s}}, \quad (37)$$

where  $r_{i,f}^D$  and  $r_{f,s}$  represent the downlink data rates for transmitting sub-models from the fed server to the  $i$ -th edge device and edge server, respectively.

## B. Problem Formulation

As illustrated in Fig. 4, the per-round latency of split training can be expressed as

$$T_S(\mathbf{b}, \boldsymbol{\mu}) = \max_i \{T_i^F + T_{a,i}^U\} + T_s^F + T_s^B + \max_i \{T_{g,i}^D + T_i^B\}, \quad (38)$$

and client-side model aggregation latency is calculated as

$$T_A(\mathbf{b}, \boldsymbol{\mu}) = \max_i \{T_{c,i}^U, T_s^U\} + \max_i \{T_{c,i}^D, T_s^D\}. \quad (39)$$

Since split training is conducted in each training round and client-side model aggregation occurs every  $I$  training rounds, the total latency for  $R$  training rounds is given by

$$T(\mathbf{b}, \boldsymbol{\mu}) = RT_S(\mathbf{b}, \boldsymbol{\mu}) + \left\lceil \frac{R}{I} \right\rceil T_A(\mathbf{b}, \boldsymbol{\mu}). \quad (40)$$

As alluded in Section I, BS balances the tradeoffs between per-round latency and training convergence, while MS significantly impacts communication-computing overhead and model convergence. Consequently, a joint optimization of BS and MS is crucial for expediting the training process. To this end,

we formulate the following problem to minimize the training latency for model convergence:

$$\begin{aligned}
 \mathcal{P} : \min_{\mathbf{b}, \boldsymbol{\mu}} T(\mathbf{b}, \boldsymbol{\mu}) \quad (41) \\
 \text{s.t. C1 : } \frac{1}{R} \sum_{t=1}^R \mathbb{E}[\|\nabla_{\mathbf{w}} f(\mathbf{w}^{t-1})\|^2] \leq \varepsilon, \\
 \text{C2 : } \mu_{i,j} \in \{0, 1\}, \quad \forall i \in \mathcal{N}, j = 1, 2, \dots, L, \\
 \text{C3 : } \sum_{j=1}^L \mu_{i,j} = 1, \quad \forall i \in \mathcal{N}, \\
 \text{C4 : } \sum_{j=1}^L \mu_{i,j} (b_i \tilde{\psi}_j + b_i \tilde{\chi}_j + \tilde{\vartheta}_j + \delta_j) < v_{c,i}, \quad \forall i \in \mathcal{N}, \\
 \text{C5 : } b_i \geq 1, b_i \in \mathbb{Z}, \quad \forall i \in \mathcal{N},
 \end{aligned}$$

where  $\tilde{\psi}_j = \sum_{k=1}^j \psi_k$  and  $\tilde{\chi}_l = \sum_{k=1}^j \chi_k$  are the cumulative sum of data size (in bits) of activations and activations' gradients for the first  $j$  layers of the neural network per data sample;  $\tilde{\vartheta}_l$  represents the accumulated data size of the optimizer state for the first  $j$  layers of the neural network, depending on the choice of the optimizer (e.g. Momentum, SGD, and Adam);  $v_{c,i}$  denotes the memory limitation of the  $i$ -th edge device. Constrain C1 guarantees model convergence loss; C2 and C3 ensure the uniqueness of the cut layer for each edge device, indicating that the global model is partitioned into client-side and server-side sub-models; C4 represents the memory limitation of edge devices [44]; C5 ensures that the BS decision variable is a positive integer.

Problem (41) is a combinatorial problem with a non-convex mixed-integer non-linear objective function, which is generally NP-hard. Consequently, obtaining an optimal solution via polynomial-time algorithms is impractical.

## VI. SOLUTION APPROACH

In this section, we design an efficient iterative algorithm by decoupling problem (41) into two tractable BS and MS sub-problems and then solving them alternately.

We first derive the explicit expression of  $R$  by utilizing **Corollary 1**. Given that  $R$  is proportional to the objective function, the objective function is minimized if and only if (27) holds as an equality. In general,  $R$  is substantially larger than  $I$ , and therefore we can approximate  $\lfloor \frac{R}{I} \rfloor \approx \frac{R}{I}$  in Eqn. (40). By substituting (27) into (40), problem (41) can be converted into

$$\begin{aligned}
 \mathcal{P}' : \min_{\mathbf{b}, \boldsymbol{\mu}} \Theta(\mathbf{b}, \boldsymbol{\mu}) \quad (42) \\
 \text{s.t. C2} - \text{C4},
 \end{aligned}$$

where

$$\Theta(\mathbf{b}, \boldsymbol{\mu}) = \frac{2\vartheta \left( T_S(\mathbf{b}, \boldsymbol{\mu}) + \frac{T_A(\mathbf{b}, \boldsymbol{\mu})}{I} \right)}{\gamma \left( \varepsilon - \frac{\beta \gamma \sum_{i=1}^N \sum_{j=1}^L \frac{\sigma_i^2}{b_i}}{N^2} - \mathbb{1}_{\{I>1\}} 4\beta^2 \gamma^2 I^2 \sum_{j=1}^{L_c} G_j^2 \right)}. \quad (43)$$

Problem (42) is a mixed-integer non-linear programming problem, and the non-convexity and non-smoothness of the objective function renders it still intractable. To tackle this

issue, we define a set of constants  $\tilde{\mathbf{G}} = [\tilde{G}_1^2, \tilde{G}_2^2, \dots, \tilde{G}_L^2]$ , where  $\tilde{G}_j^2$  denotes the cumulative sum of the bounded second-order moments for the first  $j$  layers of neural network, defined as  $\tilde{G}_j^2 = \sum_{k=1}^j G_k^2$ . Thus, the term  $\sum_{j=1}^{L_c} G_j^2$  in the objective function (43) can be reformulated as  $\max_i \left\{ \sum_{j=1}^L \mu_{i,j} \tilde{G}_j^2 \right\}$ . To linearize the objective function, we introduce a set of auxiliary variables  $\mathbf{T} = [T_1, T_2, \dots, T_6]$ , i.e.,  $\max_i \left\{ \sum_{j=1}^L \mu_{i,j} \tilde{G}_j^2 \right\} \leq T_1$ ,  $\max_i \left\{ \sum_{j=1}^L \mu_{i,j} \delta_j \right\} \leq T_2$ ,  $\max_i \left\{ T_i^F + T_a^U \right\} \leq T_3$ ,  $\max_i \left\{ T_{g,i}^D + T_i^B \right\} \leq T_4$ ,  $\max_i \left\{ T_{c,i}^U, T_s^U \right\} \leq T_5$ , and  $\max_i \left\{ T_{c,i}^D, T_s^D \right\} \leq T_6$ . As a consequence, problem (42) can be transformed into

$$\begin{aligned}
 \mathcal{P}'' : \min_{\mathbf{b}, \boldsymbol{\mu}, \mathbf{T}} \Theta'(\mathbf{b}, \boldsymbol{\mu}, \mathbf{T}) \quad (44) \\
 \text{s.t. C2} - \text{C5}, \\
 \text{R1 : } \sum_{j=1}^L \mu_{i,j} \tilde{G}_j^2 \leq T_1, \quad \forall i \in \mathcal{N}, \\
 \text{R2 : } \sum_{j=1}^L \mu_{i,j} \delta_j \leq T_2, \quad \forall i \in \mathcal{N}, \\
 \text{R3 : } \frac{b_i \sum_{j=1}^L \mu_{i,j} \rho_j}{f_i} + \frac{b_i \sum_{j=1}^L \mu_{i,j} \psi_j}{r_i^U} \leq T_3, \quad \forall i \in \mathcal{N}, \\
 \text{R4 : } \frac{b_i \sum_{j=1}^L \mu_{i,j} \chi_j}{r_i^D} + \frac{b_i \sum_{j=1}^L \mu_{i,j} \varpi_j}{f_i} \leq T_4, \quad \forall i \in \mathcal{N}, \\
 \text{R5 : } \frac{\sum_{j=1}^L \mu_{i,j} \delta_j}{r_{i,f}^U} \leq T_5, \quad \forall i \in \mathcal{N}, \\
 \text{R6 : } \frac{NT_2 - \sum_{i=1}^N \sum_{j=1}^L \mu_{i,j} \delta_j}{r_{s,f}} \leq T_5, \\
 \text{R7 : } \frac{\sum_{j=1}^L \mu_{i,j} \delta_j}{r_{i,f}^D} \leq T_6 \quad \forall i \in \mathcal{N}, \\
 \text{R8 : } \frac{NT_2 - \sum_{i=1}^N \sum_{j=1}^L \mu_{i,j} \delta_j}{r_{f,s}} \leq T_6,
 \end{aligned}$$

where

$$\begin{aligned}
 \Theta'(\mathbf{b}, \boldsymbol{\mu}, \mathbf{T}) = \frac{2\vartheta \left( T_3 + T_s^F + T_s^B + T_4 + \frac{T_5 + T_6}{I} \right)}{\gamma \left( \varepsilon - \frac{\beta \gamma \sum_{i=1}^N \sum_{j=1}^L \frac{\sigma_i^2}{b_i}}{N^2} - \mathbb{1}_{\{I>1\}} 4\beta^2 \gamma^2 I^2 T_1 \right)}. \quad (45)
 \end{aligned}$$

As seen from problem (44), the introduced auxiliary variables are tightly coupled with the original decision variables, posing significant barriers for solving this problem. To combat

this, we decompose problem (44) into two tractable sub-problems according to decision variables and devise efficient algorithms to solve them iteratively.

*BS sub-problem.* We fix the variables  $\boldsymbol{\mu}$  and  $\mathbf{T}$  to investigate the sub-problem involving BS, which is expressed as

$$\mathcal{P}_1 : \min_{\mathbf{b}} \Theta'(\mathbf{b}) \quad (46)$$

s.t. C4 – C5, R3 – R4.

Then, we can derive the following proposition:

**Proposition 1.** *The optimal BS decision is given by*

$$\mathbf{b}^* = [b_1^*, b_2^*, \dots, b_N^*], \quad (47)$$

where

$$b_i^* = \begin{cases} 1 & \hat{b}_i \leq 1 \\ \arg \min_{b_i \in \{\lfloor \hat{b}_i \rfloor, \lceil \hat{b}_i \rceil\}} \Theta'(\mathbf{b}) & 1 < \hat{b}_i < \kappa_i \\ \lfloor \kappa_i \rfloor & \hat{b}_i \geq \kappa_i, \end{cases} \quad (48)$$

where  $\kappa_i = \min \left\{ \frac{v_{c,i} - \sum_{j=1}^L \mu_{i,j}(\vartheta_j + \delta_j)}{\sum_{j=1}^L \mu_{i,j}(\psi_j + \chi_j)}, \frac{T_3}{\sum_{j=1}^L \mu_{i,j} \left( \frac{\rho_j}{T_i} + \frac{\psi_j}{T_i} \right)}, \frac{T_4}{\sum_{j=1}^L \mu_{i,j} \left( \frac{x_j}{T_i} + \frac{\varpi_j}{T_i} \right)} \right\}$ ,

$\hat{\mathbf{b}} = \{\hat{b}_i \mid i \in \mathcal{N}\}$  can be easily obtained by solving  $\frac{\partial \Theta'(\mathbf{b})}{\partial \mathbf{b}} = \mathbf{0}$  with Newton-Jacobi method [45];  $\lfloor \cdot \rfloor$  and  $\lceil \cdot \rceil$  represent floor and ceiling operations, respectively.

*Proof.* We first conduct a functional analysis for the objective

function in problem (46). Let  $\Theta'(\mathbf{b}) = \frac{2\vartheta \left( \sum_{i=1}^N b_i C_i + D \right)}{\gamma \left( A - \sum_{i=1}^N \frac{B}{b_i} \right)}$ ,

where  $A = \varepsilon - \mathbb{1}_{\{I>1\}} 4\beta^2 \gamma^2 I^2 T_1$ ,  $B = \frac{\beta \gamma}{N^2} \sum_{j=1}^L \sigma_j^2$ ,  $C_i =$

$\frac{\sum_{j=1}^L \mu_{i,j}(\rho_L - \rho_j + \varpi_L - \varpi_j)}{f_s}$ , and  $D = T_3 + T_4 + \frac{T_5 + T_6}{I}$ . Without

loss of generality, we take the first-order derivative for the BS of the arbitrary  $i'$ -th edge device  $b_{i'}$ , yielding

$$\frac{\partial \Theta'(\mathbf{b})}{\partial b_{i'}} = \frac{2\vartheta}{\gamma} \frac{\Xi(\mathbf{b})}{\left( A - \sum_{i=1}^N \frac{B}{b_i} \right)^2}, \quad (49)$$

where

$$\Xi(\mathbf{b}) = C_{i'} \left( A - \sum_{i=1}^N \frac{B}{b_i} \right) - \left( \sum_{i=1}^N b_i C_i + D \right) \frac{B}{b_{i'}}. \quad (50)$$

Since  $\frac{\partial \Xi(\mathbf{b})}{\partial b_{i'}} = \left( \sum_{i=1}^N b_i C_i + D \right) \frac{2B}{b_{i'}^3} > 0$ ,  $\Xi(\mathbf{b})$  is an increasing function of  $b_{i'}$ . Considering that  $\lim_{b_{i'} \rightarrow +\infty} \Xi(\mathbf{b}) = \frac{2\vartheta}{\gamma} \frac{C_{i'}}{\left( A - \sum_{k=1, k \neq i'}^N \frac{B}{b_k} \right)} > 0$  and

$\lim_{b_{i'} \rightarrow 0^+} \Xi(\mathbf{b}) = -\frac{2\vartheta}{\gamma B} \left( \sum_{k=1, k \neq i'}^N b_k C_k + D \right) < 0$ , there

must exist  $\hat{b}_{i'} \in (0, +\infty)$  satisfying  $\Xi(\mathbf{b})|_{b_{i'} = \hat{b}_{i'}} = 0$ , indicating that the objective function first decreases and then increases with respect to  $b_{i'}$  and  $\cdot$ . Therefore, the minimum is taken at  $b_{i'} = \hat{b}_{i'}$ . Considering  $N$  edge devices with diverse BSs, we generalize (50) into a system of  $N$  equations, i.e.,  $\frac{\partial \Theta'(\mathbf{b})}{\partial \mathbf{b}} = \left[ \frac{\partial \Theta'(\mathbf{b})}{\partial b_1}, \frac{\partial \Theta'(\mathbf{b})}{\partial b_2}, \dots, \frac{\partial \Theta'(\mathbf{b})}{\partial b_N} \right] = \mathbf{0}$ . To solve this system

**Algorithm 2** BCD-based Algorithm.

**Input:**  $\mathbf{b}^{(0)}$ ,  $\boldsymbol{\mu}^{(0)}$ ,  $\mathbf{T}^{(0)}$  and  $\varepsilon$ .

**Output:**  $\mathbf{b}^*$  and  $\boldsymbol{\mu}^*$

1: Initialization:  $\tau \leftarrow 0$ .

2: **repeat**

3:  $\tau \leftarrow \tau + 1$

4: Update  $\mathbf{b}^{(\tau)}$  by solving problem (46)

5: Update  $\boldsymbol{\mu}^{(\tau)}$  and  $\mathbf{T}^{(\tau)}$  by solving problem (53)

6: **until**  $|\Theta'(\mathbf{b}^{(\tau)}, \boldsymbol{\mu}^{(\tau)}, \mathbf{T}^{(\tau)}) - \Theta'(\mathbf{b}^{(\tau-1)}, \boldsymbol{\mu}^{(\tau-1)}, \mathbf{T}^{(\tau-1)})| \leq \varepsilon$

of equations, we can utilize the Newton-Jacobi method [45], yielding the solution  $\hat{\mathbf{b}} = [\hat{b}_1, \hat{b}_2, \dots, \hat{b}_N]$ . According to constraint C5, which requires the BS to be a positive integer, the optimal batch size  $b_{i'}^*$  for the  $i'$ -th edge device must be one of the two integers closest to the continuous solution  $\hat{b}_{i'}$ . Combining this insight with constraints C4, R3 and R4, the optimal solution is given by

$$\mathbf{b}^* = [b_1^*, b_2^*, \dots, b_N^*], \quad (51)$$

where  $b_i^*$  is determined as

$$b_i^* = \begin{cases} 1 & \hat{b}_i \leq 1 \\ \arg \min_{b_i \in \{\lfloor \hat{b}_i \rfloor, \lceil \hat{b}_i \rceil\}} \Theta'(\mathbf{b}) & 1 < \hat{b}_i < \kappa_i \\ \lfloor \kappa_i \rfloor & \hat{b}_i \geq \kappa_i, \end{cases} \quad (52)$$

Here,  $\kappa_i = \min \left\{ \frac{v_{c,i} - \sum_{j=1}^L \mu_{i,j}(\vartheta_j + \delta_j)}{\sum_{j=1}^L \mu_{i,j}(\psi_j + \chi_j)}, \frac{T_3}{\sum_{j=1}^L \mu_{i,j} \left( \frac{\rho_j}{T_i} + \frac{\psi_j}{T_i} \right)}, \frac{T_4}{\sum_{j=1}^L \mu_{i,j} \left( \frac{x_j}{T_i} + \frac{\varpi_j}{T_i} \right)} \right\}$ .  $\square$

**Remark:** Due to constraints C4, R3 and R4, the BS for each edge device has three potential solutions, i.e., 1,  $\lfloor \hat{b}_i \rfloor$  or  $\lceil \hat{b}_i \rceil$ , and  $\lfloor \kappa_i \rfloor$ . Based on **Proposition 1**, the optimal solution to problem (46) can be obtained via exhaustive search over all  $3^N$  possible combinations – a reduction from the original  $B^N$  combinations, where  $B$  is the maximum BS. The objective function is calculated for each combination, and the one with the minimum value is identified as the global optimum. While exhaustive search is feasible for small-scale systems, it becomes computationally prohibitive as  $N$  increases. To overcome this scalability challenge, an efficient sub-optimal solution can be derived using **Proposition 1**. Specifically, by solving the first-order condition  $\frac{\partial \Theta'(\mathbf{b})}{\partial \mathbf{b}} = \mathbf{0}$ ,  $\hat{\mathbf{b}} = \{\hat{b}_i \mid i \in \mathcal{N}\}$  is obtained. Each  $\hat{b}_i$  is then discretized and adjusted using Eqn. (48) to yield the final integer BS  $b_i^*$ . This approach only requires solving a nonlinear system of  $N$  variables followed by a one-time correction step, significantly reducing computational complexity.

Therefore, we obtain an efficient solution to problem (46).

*MS sub-problem.* By fixing the decision variable  $\mathbf{b}$ , we transform problem (44) into a standard mixed-integer linear fractional programming with respect to  $\boldsymbol{\mu}$  and  $\mathbf{T}$ , which is expressed as

$$\mathcal{P}_2 : \min_{\boldsymbol{\mu}, \mathbf{T}} \Theta'(\boldsymbol{\mu}, \mathbf{T}) \quad (53)$$

s.t. C2 – C4, R1 – R8.

TABLE I  
 SIMULATION PARAMETERS.

Parameter	Value	Parameter	Value
$f_s$	20 TFLOPS	$f_i$	[1, 2] TFLOPS
$N$	20	$r_i^U/r_{i,f}^U$	[75, 80] Mbps
$r_i^D/r_{i,f}^D$	[360, 380] Mbps	$r_s^U/r_s^D$	[360, 380] Mbps
$\gamma$	$5 \times 10^{-4}$	$I$	15

We leverage the Dinkelbach algorithm [46] to solve problem (53) by transforming it into mixed-integer linear programming, which guarantees the optimal solution [47], [48].

As mentioned earlier, we decompose the original problem (41) into two tractable sub-problems  $\mathcal{P}1$  and  $\mathcal{P}2$  according to decision variables and develop efficient algorithms to solve each sub-problem. We design an iterative block-coordinate descent (BCD)-based algorithm [49] to solve problem (41), as described in **Algorithm 2**. The key parameters required for executing the algorithm (e.g.,  $\beta$ ,  $G_j^2$  and  $\sigma_j^2$ ) are estimated following the approach in [24].

## VII. PERFORMANCE EVALUATION

This section provides numerical results to evaluate the performance of HASFL from three aspects: i) comparisons with four benchmarks to demonstrate the superiority of HASFL; ii) investigating the robustness of HASFL to varying network computing and communication resources, and the number of edge devices; iii) conducting the ablation study to validate the necessity of each meticulously designed component in HASFL, including BS and MS.

### A. Experimental Setup

**Implementation and hyper-parameters.** HASFL is implemented using Python 3.7 and PyTorch 1.9.1. All training procedures are conducted on a ThinkPad P17 Gen1 workstation, which is configured with an NVIDIA Quadro RTX 3000 GPU, an Intel i9-10885H processor, and a 4TB solid-state drive. We deploy  $N$  edge devices, and  $N$  is set to 20 by default unless specified otherwise. The computing capability of each edge device is uniformly distributed within [1, 2] TFLOPS, and the edge server is provisioned with a computing capability of 20 TFLOPS. The uplink data rates from the  $i$ -th edge device to the edge server and fed server follow uniform distribution within [75, 80] Mbps, and the corresponding downlink data rates are uniformly distributed within [360, 380] Mbps. For convenience, the inter-server data rate, namely,  $r_s^U$  and  $r_s^D$ , also identically follow uniform distribution within [360, 380] Mbps. The client-side sub-model aggregation interval and learning rate are set to 15 and  $5 \times 10^{-4}$ , respectively. For readers' convenience, the detailed experiment parameters are summarized in Table I.

**Dataset and model.** To evaluate the learning performance of HASFL, we adopt two widely used image classification datasets, CIFAR-10 and CIFAR-100 [50]. The CIFAR-10 dataset consists of 10 distinct categories of object images, such as airplanes and trucks, and contains 50000 training samples and 10000 test samples. The CIFAR-100 dataset comprises object images from 100 categories, with each cat-

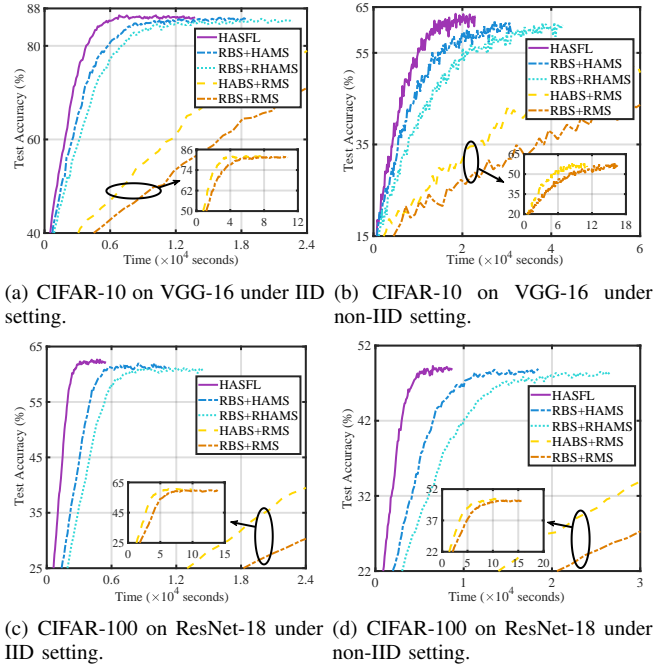


Fig. 5. The training performance for CIFAR-10 and CIFAR-100 datasets under IID and non-IID settings using VGG-16 and ResNet-18.

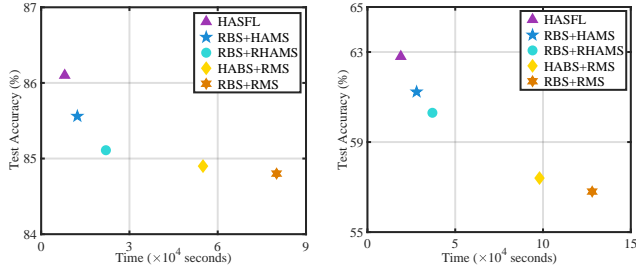
egory containing 500 training samples and 100 test samples. The experiments are conducted under IID and non-IID data settings. In the IID setting, the dataset is randomly shuffled and evenly allocated across all edge devices. In the non-IID setting [7], [51], [52], the dataset is first sorted based on class labels, and then partitioned into 40 shards, with each of 20 edge devices receiving two randomly distributed shards. Additionally, we employ the well-known ML models, VGG-16 [53] and ResNet-18 [54] for performance evaluation. VGG-16 is a classical deep convolutional neural network composed of 13 convolution layers and 3 fully connected layers, while ResNet-18 is a residual neural network consisting of 17 convolutional layers and 1 fully connected layer.

**Benchmarks.** To comprehensively evaluate the performance of HASFL, we compare it against the following alternatives:

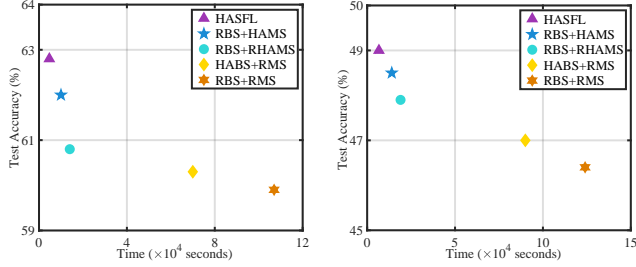
- **RBS+HAMS:** This benchmark utilizes a random BS strategy (i.e., randomly drawing BS from 1 to 64 during model training), and employs the tailored heterogeneity-aware MS scheme in Section VI.
- **HABS+RMS:** This benchmark adopts the heterogeneity-aware BS strategy in Section VI and deploys a random MS scheme (i.e., randomly selecting model split points during model training).
- **RBS+RMS:** This benchmark employs the random BS and MS strategy.
- **RBS+RHAMS:** This benchmark utilizes the RBS scheme and adopts a resource-heterogeneity-aware MS strategy [55].

### B. Superiority of HASFL

Fig. 5 presents the training performance of HASFL and four benchmarks on the CIFAR-10 and CIFAR-100 datasets. HASFL demonstrates superior convergence speed and im-

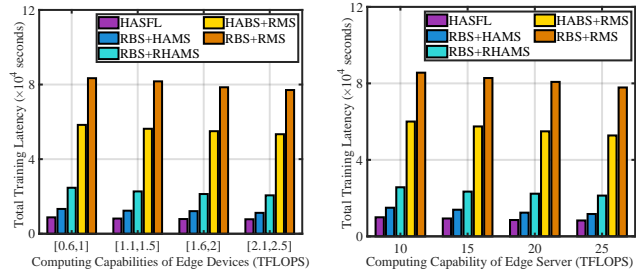


(a) CIFAR-10 on VGG-16 under IID setting. (b) CIFAR-10 on VGG-16 under non-IID setting.



(c) CIFAR-100 on ResNet-18 under IID setting. (d) CIFAR-100 on ResNet-18 under non-IID setting.

Fig. 6. The converged accuracy and time for CIFAR-10 and CIFAR-100 datasets under IID and non-IID settings using VGG-16 and ResNet-18.

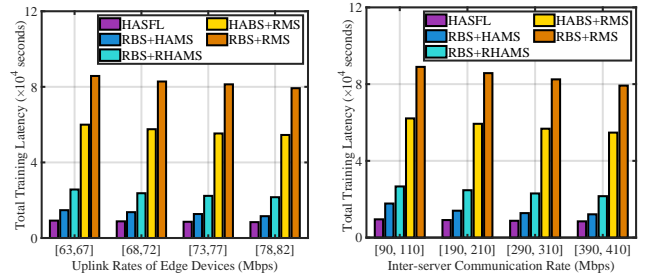


(a) Computing capabilities of edge devices. (b) Computing capabilities of edge server.

Fig. 7. The converged time versus network computing resources on CIFAR-10 under IID setting using VGG-16.

proved test accuracy compared to other benchmarks. Notably, HASFL, RBS+HAMS, and RBS+RHAMS achieve significantly faster model convergence than HABS+RMS and RBS+RMS, primarily due to heterogeneity-aware MS, which strikes a balance between communication-computing overhead and training convergence. Furthermore, the performance comparison between HASFL and RBS+HAMS reveals the advantage of the tailored BS scheme, which can expedite the model training while retaining training accuracy. The performance discrepancy between RBS+HAMS and RBS+RHAMS is mainly attributed to the heuristic design of RBS+RHAMS, which lacks systematic optimization capturing the interplay between model convergence and MS. By comparing Fig. 5(a) with Fig. 5(b), and Fig. 5(c) with Fig. 5(d), we show that HASFL and the other four benchmarks consistently converge more slowly under the non-IID setting than IID setting.

Fig. 6 illustrates the converged accuracy and time (i.e., the test accuracy increases by less than 0.02% across five consecutive training rounds) on the CIFAR-10 and CIFAR-100 datasets. Comparing HASFL with RBS+HAMS and



(a) Uplink rates of edge devices. (b) Inter-server communication rate.

Fig. 8. The converged time versus network communication resources on CIFAR-10 under IID setting using VGG-16.

HABS+RMS reveals a significant impact of BS and MS on both converged accuracy and time. Notably, MS plays a more decisive role in model training than BS. This is because the selected model split point directly affects the overall model aggregation frequency. For instance, a shallower model split point (i.e., fewer layers placed on the edge device) demonstrates that larger portions of the global model are aggregated at each training round. In the IID setting, the non-optimized MS scheme leads to an accuracy drop by 1.3% (resp. 2.4%) and nearly a sixfold (resp. fifteenfold) increase in the converged time on CIFAR-10 (resp. CIFAR-100) dataset, which is 0.7% and 4.4 times (resp. 1.6% and 7 times) that of the non-optimized MS strategy. In the non-IID setting, HASFL still significantly outperforms the non-optimized MS scheme in both converged accuracy and time, achieving approximately 5.6% (resp. 2%) accuracy increase and 4.8 (resp. 13.4) times faster model convergence acceleration on the CIFAR-10 (resp. CIFAR-100) dataset. The reason behind this is the model bias introduced by the heterogeneous nature of local data distributions, consistent with FL. By comparing HASFL with RBS+RMS, HASFL achieves a faster model convergence of 9.4 (resp. 6.4) times over its counterpart without optimization on the CIFAR-10 (resp. CIFAR-100) dataset while guaranteeing training accuracy, showcasing the superior performance of HASFL. Though RBS+HAMS, HABS+RMS, and RBS+RHAMS exhibit better training performance over RBS+RMS, we observe that the performance gains remain relatively limited because they lack a unified optimization of BS and MS.

### C. Impact of Varying Network Resources

Figs. 7 and 8 show the converged time versus network computing and communication resources on CIFAR-10 under the IID setting. HASFL consistently achieves faster model convergence than the other benchmarks across varying network resources. The training performance of RBS+RMS deteriorates significantly as network resources decline. This degradation arises from the lack of optimization of BS and MS, leading to sub-optimal trade-offs between communication-computing overhead and model convergence. The comparison of HASFL with RBS+HAMS and HABS+RMS reveals that optimizing either BS or MS can compensate for the performance decline caused by the reduced network resources. Though RBS+RHAMS achieves faster model convergence



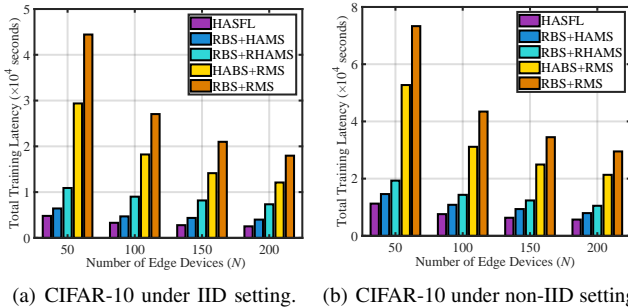


Fig. 9. The converged time versus number of edge devices on CIFAR-10 under IID and non-IID setting using VGG-16.

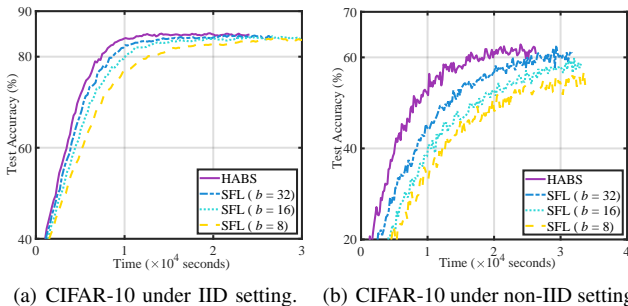


Fig. 10. Ablation experiments for HABS strategy on the CIFAR-10 dataset under IID and non-IID setting using VGG-16 with  $L_c = 8$ .

than RBS+RMS, its performance gains remain limited by the lack of joint optimization of BS and MS. In contrast, HASFL demonstrates better robustness over other benchmarks, with the converged time of HASFL exhibiting only a marginal increase as network resources diminish, owing to its heterogeneity-aware BS and MS design. Specifically, HASFL dynamically adapts BS and MS based on network resource conditions to achieve the minimum rounds required for model convergence to expedite model training, demonstrating the superior adaptability to varying network resources.

#### D. Impact of Number of Edge Devices

Fig. 9 presents the converged time versus the number of edge devices on CIFAR-10 under the IID and non-IID settings. HASFL consistently outperforms the other benchmarks by achieving significantly lower training latency across all configurations of  $N$ . Notably, the performance gap between HASFL and the other benchmarks becomes increasingly pronounced with the increasing number of edge devices, underscoring HASFL's superior scalability and efficiency in large-scale edge networks. The converged time of HASFL and all benchmarks is noticeably longer under the non-IID setting than under IID setting. This degradation in performance is primarily attributed to the statistical heterogeneity of non-IID data, which often results in inconsistent local updates across edge devices and hinders model convergence. HASFL maintains its advantage over the other benchmarks, underscoring its robustness to diverse and imbalanced data distributions.

#### E. Ablation Study of The HASFL Framework

Fig. 10 illustrates the impact of BS on training performance for the CIFAR-10 dataset. The proposed HABS scheme

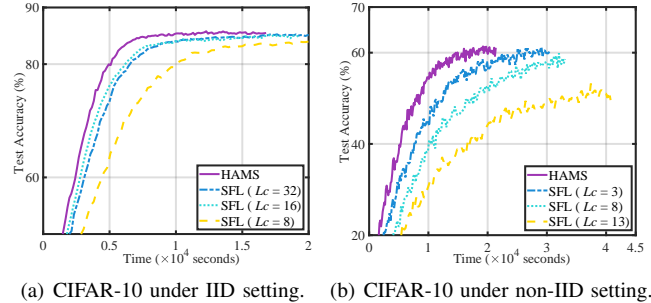


Fig. 11. Ablation experiments for HAMS scheme on the CIFAR-10 dataset under IID and non-IID setting using VGG-16 with  $b = 16$ .

achieves faster model convergence and higher accuracy than fixed BS benchmarks, where all edge devices utilize the same BS ( $b = 8, 16$ , and  $32$ ). Specifically, HABS achieves at least 0.9% and 1% accuracy improvements and reduces the convergence time by a factor of 1.4 and 1.6 compared to fixed BS benchmarks under the IID and non-IID settings, respectively. This performance gain primarily stems from the principled design of HABS, which leverages a theoretical convergence bound to determine BS. By dynamically assigning the optimal BS to each edge device based on network resource conditions and training convergence, HABS effectively mitigates the straggler effect to expedite model training under heterogeneous edge systems. Hence, the effectiveness of HABS is validated.

Fig. 11 presents the impact of MS on training performance for the CIFAR-10 dataset. Both converged time and accuracy degrade as the model split point  $L_c$  becomes deeper, consistent with the derived convergence bound in Eqn. (16). This performance decline occurs because a deeper model split point places fewer layers on the server-side sub-model, reducing the proportion of the global model that is aggregated at each round. Consequently, the overall update frequency of the global model decreases, slowing model convergence and impairing the model's generalization. Furthermore, the impact of MS is more significant under the non-IID than IID setting, as the heterogeneity of local datasets exacerbates the sensitivity of MS under the non-IID setting. The comparisons between HAMS and the other three benchmarks reveal that HAMS can expedite model convergence and improve training performance. This demonstrates the effectiveness of the HAMS design in adapting to heterogeneous network resources.

## VIII. CONCLUSIONS

In this paper, we have proposed a heterogeneity-aware SFL framework, named HASFL, to minimize the training latency of SFL under heterogeneous resource-constrained edge computing systems. We have observed that both BS and MS significantly affect training latency, control of which can mitigate the straggler issue to a great extent. To guide the system optimization, we have first derived the convergence bound of HASFL to quantify the impact of BS and MS on training convergence. Then, we have formulated a joint optimization problem for BS and MS based on the derived convergence bound, and developed efficient algorithms to solve it. Extensive experiments on various datasets have validated the effectiveness of HASFL, demonstrating at least  $4\times$

faster convergence and 1% higher training and testing accuracy compared to state-of-the-art benchmarks.

## REFERENCES

- [1] X. Liu, Z. Yan, Y. Zhou, D. Wu, X. Chen, and J. H. Wang, "Optimizing Parameter Mixing Under Constrained Communications in Parallel Federated Learning," *IEEE/ACM Trans. Networking*, vol. 31, no. 6, pp. 2640–2652, Dec. 2023.
- [2] Z. Lin, Y. Zhang, Z. Chen, Z. Fang, X. Chen, P. Vepakomma, W. Ni, J. Luo, and Y. Gao, "HSplitLoRA: A Heterogeneous Split Parameter-Efficient Fine-Tuning Framework for Large Language Models," *arXiv preprint arXiv:2505.02795*, 2025.
- [3] Y. Deng, X. Chen, G. Zhu, Y. Fang, Z. Chen, and X. Deng, "Actions at the Edge: Jointly Optimizing the Resources in Multi-access Edge Computing," *IEEE Wireless Commun.*, vol. 29, no. 2, pp. 192–198, Apr. 2022.
- [4] X. Chen, G. Zhu, Y. Deng, and Y. Fang, "Federated Learning over Multihop Wireless Networks with In-network Aggregation," *IEEE Trans. Wireless Commun.*, vol. 21, no. 6, pp. 4622–4634, Apr. 2022.
- [5] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. Arcas, "Communication-efficient Learning of Deep Networks From Decentralized Data," in *Proc. AISTATS*, Apr. 2017.
- [6] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated Learning: Strategies for Improving Communication Efficiency," *arXiv preprint arXiv:1610.05492*, Oct. 2016.
- [7] Z. Lin, Z. Chen, Z. Fang, X. Chen, X. Wang, and Y. Gao, "FedSN: A Federated Learning Framework over Heterogeneous LEO Satellite Networks," *IEEE Trans. Mobile Comput.*, 2024.
- [8] M. Hu, J. Zhang, X. Wang, S. Liu, and Z. Lin, "Accelerating Federated Learning with Model Segmentation for Edge Networks," *IEEE Trans. Green Commun. Netw.*, 2024.
- [9] Z. Fang, Z. Lin, Z. Chen, X. Chen, Y. Gao, and Y. Fang, "Automated Federated Pipeline for Parameter-efficient Fine-tuning of Large Language Models," *arXiv preprint arXiv:2404.06448*, 2024.
- [10] Z. Lin, Y. Zhang, Z. Chen, Z. Fang, C. Wu, X. Chen, Y. Gao, and J. Luo, "Leo-Split: A Semi-Supervised Split Learning Framework over LEO Satellite Networks," *arXiv preprint arXiv:2501.01293*, 2025.
- [11] G. Team, R. Anil, S. Borgeaud, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, K. Millican *et al.*, "Gemini: A Family of Highly Capable Multimodal Models," *arXiv preprint arXiv:2312.11805*, Dec. 2023.
- [12] 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*, Dec. 2018.
- [13] W. Wei, Z. Lin, T. Li, X. Li, and X. Chen, "Pipelining Split Learning in Multi-hop Edge Networks," *arXiv preprint arXiv:2505.04368*, 2025.
- [14] Z. Lin, G. Qu, Q. Chen, X. Chen, Z. Chen, and K. Huang, "Pushing Large Language Models to the 6G Edge: Vision, Challenges, and Opportunities," *IEEE Commun. Mag.*, 2023.
- [15] S. Lyu, Z. Lin, G. Qu, X. Chen, X. Huang, and P. Li, "Optimal Resource Allocation for U-Shaped Parallel Split Learning," in *Proc. IEEE Globecom Wkshps*, 2023, pp. 197–202.
- [16] C. Thapa, P. C. M. Arachchige, S. Camepe, and L. Sun, "Splitfed: When Federated Learning Meets Split Learning," in *Proc. AAAI*, Feb. 2022.
- [17] M. Chen, D. Gündüz, K. Huang, W. Saad, M. Bennis, A. V. Feljan, and H. V. Poor, "Distributed Learning in Wireless Networks: Recent Progress and Future Challenges," *IEEE J. Sel. Areas Commun.*, Dec. 2021.
- [18] Y. Xiao, X. Zhang, Y. Li, G. Shi, M. Krunz, D. N. Nguyen, and D. T. Hoang, "Time-sensitive Learning For Heterogeneous Federated Edge Intelligence," *IEEE Trans. Mobile Comput.*, vol. 23, no. 2, pp. 1382–1400, 2023.
- [19] K. Lee, M. Lam, R. Pedarsani, D. Papailiopoulos, and K. Ramchandran, "Speeding Up Distributed Machine Learning Using Codes," *IEEE Trans. Inf. Theory*, vol. 64, no. 3, pp. 1514–1529, Mar. 2017.
- [20] W. Wu, M. Li, K. Qu, C. Zhou, X. Shen, W. Zhuang, X. Li, and W. Shi, "Split learning over Wireless Networks: Parallel Design and Resource Management," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 4, pp. 1051–1066, Feb. 2023.
- [21] Z. Lin, G. Qu, W. Wei, X. Chen, and K. K. Leung, "AdaptsFL: Adaptive Split Federated Learning in Resource-Constrained Edge Networks," *arXiv preprint arXiv:2403.13101*, 2024.
- [22] J. Lee, J. Cho, W. Lee, M. Seif, and H. V. Poor, "Game-Theoretic Joint Incentive and Cut Layer Selection Mechanism in Split Federated Learning," *arXiv preprint arXiv:2412.07813*, 2024.
- [23] Z. Lin, W. Wei, Z. Chen, C.-T. Lam, X. Chen, Y. Gao, and J. Luo, "Hierarchical Split Federated Learning: Convergence Analysis and System Optimization," *IEEE Trans. Mobile Comput.*, Apr. 2025.
- [24] S. Wang, T. Tuor, T. Salonidis, K. K. Leung, C. Makaya, T. He, and K. Chan, "Adaptive Federated Learning in Resource Constrained Edge Computing Systems," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 6, pp. 1205–1221, Mar. 2019.
- [25] S. Liu, G. Yu, R. Yin, J. Yuan, and F. Qu, "Adaptive Batchsize Selection and Gradient Compression For Wireless Federated Learning," in *Proc. GLOBECOM*, 2020.
- [26] Z. Ma, Y. Xu, H. Xu, Z. Meng, L. Huang, and Y. Xue, "Adaptive Batch Size For Federated Learning in Resource-constrained Edge Computing," *IEEE Trans. Mobile Comput.*, vol. 22, no. 1, pp. 37–53, Jan. 2021.
- [27] W. Liu, X. Zhang, J. Duan, C. Joe-Wong, Z. Zhou, and X. Chen, "DYNAMITE: Dynamic Interplay of Mini-batch Size and Aggregation Frequency For Federated Learning with Static and Streaming Datasets," *IEEE Trans. Mobile Comput.*, vol. 23, no. 7, pp. 7664–7679, 2023.
- [28] D. Shi, L. Li, M. Wu, M. Shu, R. Yu, M. Pan, and Z. Han, "To Talk or to Work: Dynamic Batch Sizes Assisted Time Efficient Federated Learning over Future Mobile Edge Devices," *IEEE Trans. Wireless Commun.*, vol. 21, no. 12, pp. 11 038–11 050, Dec. 2022.
- [29] S. L. Smith, P.-J. Kindermans, and Q. V. Le, "Don't Decay the Learning Rate, Increase the Batch Size," in *Proc. ICLR*, Feb. 2018.
- [30] H. Yu and R. Jin, "On the Computation and Communication Complexity of Parallel SGD with Dynamic Batch Sizes for Stochastic Non-convex Optimization," in *Proc. ICML*, Jun. 2019.
- [31] Z. Lin, G. Zhu, Y. Deng, X. Chen, Y. Gao, K. Huang, and Y. Fang, "Efficient Parallel Split Learning over Resource-constrained Wireless Edge Networks," *IEEE Trans. Mobile Comput.*, 2024.
- [32] P. Han, C. Huang, G. Tian, M. Tang, and X. Liu, "Convergence Analysis of Split Federated Learning on Heterogeneous Data," in *Proc. NIPS*, 2025.
- [33] D. Pasquini, G. Ateniese, and M. Bernaschi, "Unleashing the Tiger: Inference Attacks on Split Learning," in *Proc. CCS*, Nov. 2021.
- [34] Karimireddy, Sai Praneeth and Kale, Satyen and Mohri, Mehryar and Reddi, Sashank and Stich, Sebastian and Suresh, Ananda Theertha, "On the Convergence Properties of A K-step Averaging Stochastic Gradient Descent Algorithm for Nonconvex Optimization," in *Proc. IJCAI*, Jul. 2018.
- [35] H. Yu, R. Jin, and S. Yang, "On the Linear Speedup Analysis of Communication Efficient Momentum SGD For Distributed Non-convex Optimization," in *Proc. ICML*, Jun. 2019, pp. 7184–7193.
- [36] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "Scaffold: Stochastic Controlled Averaging for Federated Learning," in *Proc. ICLR*, Apr. 2020.
- [37] Y. Zhang, M. J. Wainwright, and J. C. Duchi, "Communication-efficient Algorithms for Statistical Optimization," in *Proc. NIPS*, Jun. 2012.
- [38] X. Lian, C. Zhang, H. Zhang, C.-J. Hsieh, W. Zhang, and J. Liu, "Can Decentralized Algorithms Outperform Centralized Algorithms? A Case Study for Decentralized Parallel Stochastic Gradient Descent," in *Proc. NIPS*, Jun. 2017.
- [39] H. Mania, X. Pan, D. Papailiopoulos, B. Recht, K. Ramchandran, and M. I. Jordan, "Perturbed Iterate Analysis for Asynchronous Stochastic Optimization," *SIAM J. Optim.*, vol. 27, no. 4, pp. 2202–2229, Jan. 2017.
- [40] T. Lin, S. U. Stich, K. K. Patel, and M. Jaggi, "Don't Use Large Mini-batches, Use Local SGD," in *Proc. ICLR*, Dec. 2019.
- [41] L. Yi, G. Wang, X. Wang, and X. Liu, "QSFL: Two-Level Communication-Efficient Federated Learning on Mobile Edge Devices," *IEEE Trans. Serv. Comput.*, 2024.
- [42] W. Shi, S. Zhou, Z. Niu, M. Jiang, and L. Geng, "Joint Device Scheduling and Resource Allocation for Latency Constrained Wireless Federated Learning," *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 453–467, Sep. 2020.
- [43] W. Xia, W. Wen, K.-K. Wong, T. Q. Quek, J. Zhang, and H. Zhu, "Federated-learning-based Client Scheduling for Low-latency Wireless Communications," *IEEE Wireless Commun.*, vol. 28, no. 2, pp. 32–38, Apr. 2021.
- [44] G. Yeung, D. Borowiec, R. Yang, A. Friday, R. Harper, and P. Garraghan, "Horus: Interference-aware and Prediction-based Scheduling in Deep Learning Systems," *IEEE Trans. Parallel Distrib. Syst.*, vol. 33, no. 1, pp. 88–100, May. 2021.
- [45] Z. Geng, C. Wei, Y. Han, Q. Wei, and Z. Ouyang, "Pipeline Network Simulation Calculation based on Improved Newton Jacobian Iterative Method," in *Proc. AICS*, Jul. 2019.
- [46] W. Dinkelbach, "On Nonlinear Fractional Programming," *Manage. Sci.*, vol. 13, no. 7, pp. 492–498, Mar. 1967.



- [47] D. Yue and F. You, "A Reformulation-linearization Method for the Global Optimization of Large-scale Mixed-Integer Linear Fractional Programming Problems and Cyclic Scheduling aApplication," in *Proc. ACC*, Jun. 2013.
- [48] R. G. Ródenas, M. L. López, and D. Verastegui, "Extensions of Dinkelbach's Algorithm for Solving Non-linear Fractional Programming Problems," *Top*, vol. 7, pp. 33–70, Jun. 1999.
- [49] P. Tseng, "Convergence of A Block Coordinate Descent Method for Nondifferentiable Minimization," *J Optim Theory Appl*, vol. 109, pp. 475–494, Jun. 2001.
- [50] A. Krizhevsky, G. Hinton *et al.*, "Learning Multiple Layers of Features From Tiny Images," *Tech. Rep.*, Apr. 2009.
- [51] G. Zhu, Y. Wang, and K. Huang, "Broadband analog aggregation for low-latency federated edge learning," *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 491–506, Oct. 2019.
- [52] Z. Yang, M. Chen, W. Saad, C. S. Hong, and M. Shikh-Bahaei, "Energy Efficient Federated Learning over Wireless Communication Networks," *IEEE Trans. Wireless Commun.*, vol. 20, no. 3, pp. 1935–1949, Nov. 2020.
- [53] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-scale Image Recognition," in *Proc. ICLR*, 2015.
- [54] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *Proc. CVPR*, Jun. 2016, pp. 770–778.
- [55] Z. Wang, H. Xu, Y. Xu, Z. Jiang, and J. Liu, "CoopFL: Accelerating Federated Learning with DNN Partitioning and Offloading in Heterogeneous Edge Computing," *Comput. Netw.*, vol. 220, p. 109490, Jan. 2023.