

Enhancing Federated Learning with In-Cloud Unlabeled Data

Lun Wang^{1,3}, *Yang Xu^{1,3}, Hongli Xu^{1,3}, Jianchun Liu^{2,3}, Zhiyuan Wang^{1,3}, Liusheng Huang^{1,3}

¹School of Computer Science and Technology, University of Science and Technology of China

²School of Data Science, University of Science and Technology of China

³Suzhou Institute for Advanced Research, University of Science and Technology of China

{wanglun0, jsen617, cswangzy}@mail.ustc.edu.cn, {xuyangcs, xuhongli, lshuang}@ustc.edu.cn

Abstract—Federated learning (FL) has been widely applied to collaboratively train deep learning (DL) models on massive end devices (*i.e.*, clients). Due to the limited storage capacity and high labeling cost, there are always insufficient data stored and annotated on each client. Conversely, in cloud datacenters, there exist large-scale unlabeled data, which are easy to collect from public access (*e.g.*, social media). Herein, upon the federated semi-supervised learning (FSSL) technology, we propose the *Ada-FedSemi* system, which leverages both on-device labeled data and in-cloud unlabeled data to boost the performance of DL models. Given the limited communication and massive quantity of the clients, in each training round, we decide to select partial clients to participate in FL, and their local models are aggregated by the parameter server (PS) to produce pseudo-labels for the unlabeled data, which are utilized to enhance the global model. Considering that the number of participating clients and the quality of pseudo-labels will have a significant impact on the training performance (*e.g.*, efficiency and accuracy), we introduce a multi-armed bandit (MAB) based online algorithm to adaptively determine the participating fraction and confidence threshold during federated model training. Extensive experiments on benchmark models and datasets show that, given the same resource budget, the model trained by *Ada-FedSemi* achieves 3%-14.8% higher test accuracy than that of the baseline methods. Besides, when achieving the same test accuracy, *Ada-FedSemi* saves up to 48% training cost, compared with the baselines.

Index Terms—Edge Computing, Federated Learning, Semi-supervised Learning, Pseudo-labeling.

I. INTRODUCTION

With the considerable development of deep learning (DL) in recent years, more and more AI applications are penetrating our daily life, such as smart transportation [1], virtual reality [2] and intelligent assistants [3]. In order to utilize data generated at the network edge without possible leakage of personal privacy, federated learning (FL) [4] is proposed to collaboratively train DL models on massive end devices with the aid of parameter servers (PS). In FL, end devices (*i.e.*, clients) keep their data locally during training and only upload local models to the PS periodically for global aggregation [5]. Then, the PS broadcasts the global model back to clients, and the interaction procedure will last until the model converges.

The most existing works of FL concentrate on training efficiency and assume sufficient labeled data on clients. However, due to the high labeling cost or lack of expert knowledge for annotation, the scale of labeled data on each client may be small [6]. Although some simple labels on mobile phones

can be obtained automatically (*e.g.*, location of landscape photos) [7], other valuable labels need to be annotated by users or domain experts manually (*e.g.*, age and gender labels of face images or morphologic features of pathology images) [8]. Since the strong performance of DL models is largely attributed to the availability of abundant data (especially the labeled data), the insufficient labeled data on clients may introduce overfitting to the DL models [9]. On the contrary, there are various other data sources (*e.g.*, social networks [10]) continuously generating different types of unlabeled data, including text, images, videos and so on [11], which are collected and stored in cloud datacenters. For example, the large-scale WebVision Database [12] consists of 2.4 million web images crawled from the Internet, and the image labels are always missing or contain errors [13].

To fully utilize both labeled and unlabeled data in FL, a new technology of federated semi-supervised learning (FSSL) has been proposed [14]. Long *et al.* [15] and Jeong *et al.* [14] assume the (labeled and unlabeled) data are already on clients, while some works [16], [17] distribute unlabeled data from the cloud to clients and then implement FSSL using the mixed data on clients. However, limited by the storage capacity of clients, the data size on clients is much smaller than that in cloud, which restricts the performance boost of DL models. Besides, delivering additional data from the cloud to clients will incur a large amount of communication cost and also increase computation overhead for the resource-constrained clients. Instead, in works [18], [19], the PS first collects local models from all clients as teacher models, and then produces pseudo-labels for the in-cloud unlabeled data in terms of the teachers' predictions. Subsequently, the pseudo-labeled data are exploited to improve the trained model. However, since they adopt all pseudo-labeled data in the training without considering the quality of pseudo-labels, there may be many incorrect labels, which will lead to noisy training [20]. Moreover, given the massive quantity of the clients, collecting local models from all the clients will result in extremely high communication cost, which is infeasible for FL systems.

In this paper, we consider a practical FL scenario where there are limited labeled data on clients and large-scale unlabeled data in cloud, and the clients are equipped with limited computation, communication and storage resource. We fuse FL and semi-supervised learning, and leverage both the on-device labeled data and in-cloud unlabeled data to

better boost the performance of DL models, even when the labeled data are not independent and identically distributed (non-IID) across clients. During FL training, the quality of pseudo-labels (depicted as *confidence*) is low at the early stage and increases gradually as the training progresses [20]. If a large number of low-confidence pseudo-labels are adopted, the training performance will be significantly degraded, resulting in poor generalization and resource waste. Besides, the number of participating clients (referred as *participating fraction*) has an impact on training cost, including time and communication. With the increasing participating fraction, the communication cost will increase and the time will be shortened when achieving the same accuracy [21]. Note that different FL tasks usually have various preferences towards communication cost and time cost, and their preferences may change over time during the training process. For example, time-critical tasks (*e.g.*, search-and-rescue tasks) always require fast convergence while other less urgent tasks in cellular network may expect small traffic consumption. Therefore, the approaches [21], [22] with fixed strategies can barely achieve satisfactory training performance.

To improve training efficiency as well as model accuracy, we propose an adaptive FSSL system termed *Ada-FedSemi*, which employs a multi-armed bandit (MAB) based online learning algorithm to adaptively determine the participating fraction (*i.e.*, P) and confidence threshold (*i.e.*, C) in terms of the resource budgets and cost preferences. Most prior researches of semi-supervised learning concentrate on achieving the state-of-the-art model accuracy but ignore training efficiency [23]. Instead, our paper concentrates on the trade-off between model performance and training cost, which is an important issue in practical distributed systems such as an FL system [24], [25]. The main contributions of this paper are summarized as follows:

- Considering the constrained computation and communication resource of clients at the network edge, we propose to exploit limited on-device labeled data and large-scale in-cloud unlabeled data to boost the training performance of FL in a semi-supervised way.
- To adapt to different cost preferences of FL tasks, we present a multi-armed bandit based online algorithm to adaptively determine the participation fraction of clients and pseudo-labeling confidence threshold for federated model training (*i.e.*, *Ada-FedSemi*) to improve training efficiency and model accuracy.
- We develop an FL hardware prototype system and conduct extensive experiments on benchmark models and datasets. The experimental results demonstrate that, (i) given the same resource budget, *Ada-FedSemi* can improve test accuracy by 3%-14.8%, (ii) when achieving the same test accuracy, *Ada-FedSemi* saves up to 48% training cost, compared with the baseline methods.

The rest of the paper is organized as follows. Section II introduces the adaptive FSSL system and formulates the optimization problem. Section III describes the MAB based online algorithm. The experimental evaluation is presented in Section IV. Section V reviews some related works and Section VI gives the conclusions.

TABLE I: Summary of some important notations.

Notation	Meaning (at round k)
\mathcal{V}	set of clients $\mathcal{V} = \{v_1, v_2, \dots, v_M\}$
\mathcal{V}_k	set of clients participating at round k
\mathcal{D}_L	labeled dataset across clients
\mathcal{D}_U	in-cloud unlabeled dataset
$\mathcal{D}_{k,U}$	selected unlabeled data for training
\hat{y}_j	one-hot pseudo-label for the j -th unlabeled data
a_j	confidence of the pseudo-label \hat{y}_j
P_k	participating fraction of clients
C_k	confidence threshold for pseudo-labels
$t_{k,m}$	time cost of client v_m
$t_{k,p}$	time cost of the PS
t_k	total time cost
W	communication cost of delivering a local model
b_k	total communication cost
α	bias factor for cost preference
Φ_k	weighted cost

II. SYSTEM DESCRIPTION AND PROBLEM DEFINITION

In this section, we first introduce the FSSL system and the main training procedure of *Ada-FedSemi*. Then, we conduct several experiments to show the motivation of utilization of unlabeled data and further present the impact of participating fraction as well as confidence threshold. Finally, we formally describe the problem to be solved in this paper. For ease of expression, we list some important notations in Table I.

A. System Description for FSSL

An FSSL system usually includes a parameter server (PS) and a set of M distributed clients (*e.g.*, IoT devices and edge nodes) $\mathcal{V} = \{v_1, v_2, \dots, v_M\}$, which collaboratively train DL models over the networks. Each client $v_m \in \mathcal{V}$ trains a local model on its own private dataset \mathcal{D}_m with N_m labeled data, and only needs to synchronize model parameters with the PS rather than sharing the original data, which prevents the exposure of personal privacy.

Let $\mathcal{D} = \mathcal{D}_L \cup \mathcal{D}_U$ denote the whole training dataset, where $\mathcal{D}_L = \mathcal{D}_{1,L} \cup \mathcal{D}_{2,L} \cup \dots \cup \mathcal{D}_{M,L}$ is the labeled dataset distributed across clients and \mathcal{D}_U is the unlabeled dataset collected in the cloud. For the sake of description, we assume that there is no intersection between local datasets. Thus, there are $N_L = \sum_{m=1}^M N_m$ data samples in $\mathcal{D}_L = \{(x_i, y_i)\}_{i=1}^{N_L}$, where x_i is the features of the i -th data sample and $y_i = [y_{i,1}, \dots, y_{i,Q}] \in \{0, 1\}^Q$ is a one-hot label, and Q is the total number of classes. $y_{i,q} = 1, q \in [1, Q]$ means that the data sample x_i belongs to class q . For the unlabeled dataset, there are N_U data samples in $\mathcal{D}_U = \{x_j\}_{j=1}^{N_U}$, which lack the ground-truth labels. Let $F(w, x, y)$ denote the loss function over the data (x, y) and $w \in \mathbb{R}^d$ is the model parameter with d dimensions. When considering FL on labeled local data, for each client v_m , its local loss function is defined as:

$$f_m(w) = \min_{w \in \mathbb{R}^d} \mathbb{E}_{(x,y) \sim \mathcal{D}_m} F(w, x, y). \quad (1)$$

In order to utilize unlabeled data during the model training, pseudo-labeling is a general and efficient method [20], [26],

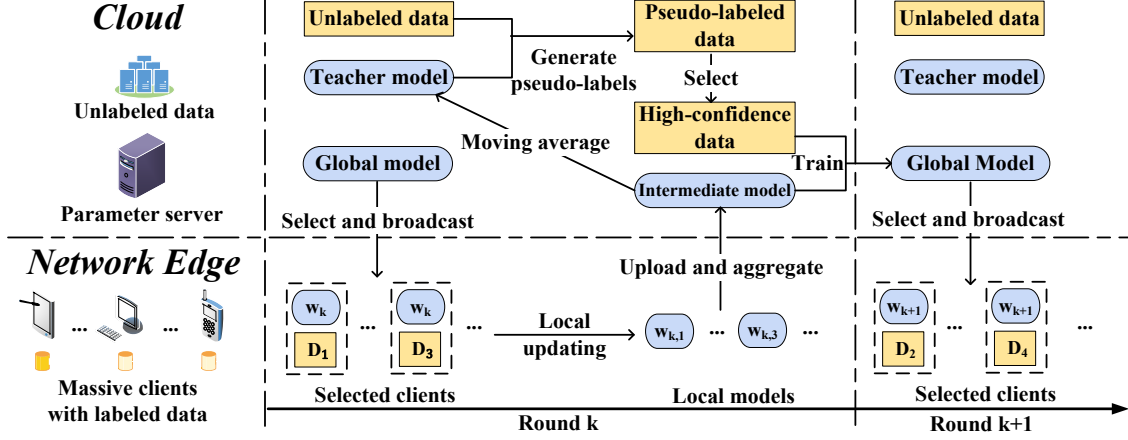


Fig. 1: System workflow of Ada-FedSemi.

in which there are two alternating steps, including training and labeling. In the training step, models are trained on both labeled and pseudo-labeled data, which is similar to traditional supervised learning but has different loss functions. In the labeling step, a trained model, also called teacher model, is used to produce predictions for unlabeled data. For a certain data sample x_j , the prediction by the teacher model is $p_j = [p_{j,1}, \dots, p_{j,Q}] \in [0, 1]^Q$ and $\sum_{q=1}^Q p_{j,q} = 1$. The pseudo-label \hat{y}_j of sample x_j is defined as:

$$\hat{y}_j = \underset{q}{\operatorname{argmax}} p_{j,q}. \quad (2)$$

When the model is trained on unlabeled data, the pseudo-labels are treated as their training targets. Therefore, the FSSL system aims to optimize the following objective function on both labeled and unlabeled datasets:

$$f^* := \min_{w \in \mathbb{R}^d} \mathbb{E}_{(x,y) \sim \mathcal{D}_L} F(w, x, y) + \mathbb{E}_{(x,\hat{y}) \sim \mathcal{D}_U} F(w, x, \hat{y}). \quad (3)$$

B. System Workflow of Ada-FedSemi

In Ada-FedSemi, the model will be trained on both on-device labeled data and in-cloud unlabeled data iteratively. Fig. 1 illustrates the workflow of the proposed Ada-FedSemi system. We consider the more effective synchronous FL training scheme, since the secure aggregation and differential privacy techniques applied in FL usually require synchronous operations [27] and the convergence of asynchronous training algorithms may not be guaranteed [28]. The process involves many training rounds, and at each round k , it mainly consists of the following four steps:

(1) **Model Broadcast and Local Updating.** In this step, the PS selects a fraction of clients (let P_k denote the participating fraction and \mathcal{V}_k denote the set of selected clients at round k), and broadcasts the global model w_k to them. Most researches [4], [29]–[31] select specific clients based on a predefined P_k . In general, larger fractions lead to faster convergence but result in more traffic consumption, and vice versa [21]. Thus, considering the properties of FL tasks and the limited resource in the distributed data system, our algorithm concentrates on

determining P_k . With a determined P_k , randomly selecting clients is an intuitive and efficient way to ensure that the global model can learn from different knowledge since data are always non-IID across clients in FL [4], [21], [29]. Nevertheless, some other advanced client selection strategies [30], [31] can be applied regarding the value of P_k derived from our algorithm, which may further improve the performance of FL but with a higher cost for selecting specific clients.

For each selected client v_m , it first initializes $w_{k,m}(0) = w_k$ and then performs local updating on its local dataset by stochastic gradient descent (SGD) [32]:

$$w_{k,m}(\tau' + 1) = w_{k,m}(\tau') - \eta_k \nabla f_m(w_{k,m}(\tau')), \quad (4)$$

where η_k is the learning rate, $\nabla f_m(w_{k,m}(\tau'))$ is the stochastic gradient, $\tau' \in [0, \tau)$ and τ is the number of local updating. Finally, client v_m gets its updated model $w_{k,m}$.

(2) **Model Uploading and Global Aggregation.** After finishing local updating, clients in \mathcal{V}_k upload their local models to the PS, and the PS aggregates these models based on the number of data samples in their local datasets as follows:

$$w_{k+\frac{1}{2}} = \frac{\sum_{v_m \in \mathcal{V}_k} N_m w_{k,m}}{\sum_{v_m \in \mathcal{V}_k} N_m}, \quad (5)$$

where $w_{k+\frac{1}{2}}$ is referred as the intermediate model. Then, the teacher model is updated by the intermediate model, which will be elaborated in Section III.

(3) **Pseudo-labels Generation and Selection.** At the PS, in terms of the teacher model, we make predictions for the unlabeled data and then generate pseudo-labels. Since \hat{y}_j may not be the ground-truth label, we need to estimate the confidence of the pseudo-labels, which indicates how likely a pseudo-label is true. Specifically, we regard the probability of label \hat{y}_j in the prediction as its confidence [23]:

$$a_j = \max_q p_{j,q}. \quad (6)$$

Generally, high-confidence pseudo-labels are more likely to be the ground-truth labels and vice versa.

To mitigate data noise introduced by pseudo-labeling, at round k , we only train the model on the high-confidence pseudo-labeled data samples, whose confidence a_j is over a threshold C_k [33].

(4) **Semi-supervised Model Training.** The intermediate

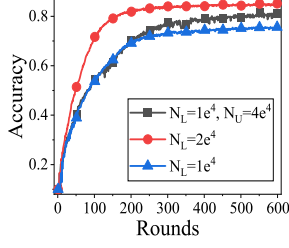


Fig. 2: Test Accuracy vs. Rounds.

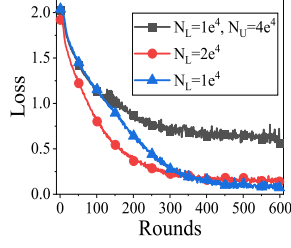


Fig. 3: Training Loss vs. Rounds.

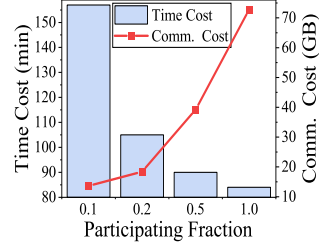


Fig. 4: Time and communication cost to achieve 60% accuracy.

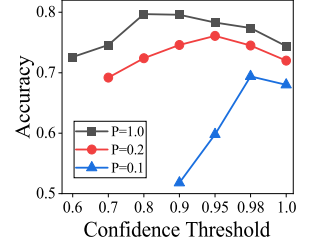


Fig. 5: Accuracy with different P and C .

model derived from the aggregation of local models may suffer from poor generalization, since the local models can easily overfit to the insufficient on-device labeled data [9]. Thus, we expect to improve the model's generalization ability by learning additional knowledge from the massive in-cloud unlabeled data. Specifically, the intermediate model in step (2) is further trained on pseudo-labeled data at the PS in a semi-supervised way:

$$w_{k+1} = w_{k+\frac{1}{2}} - \eta_k \nabla f_s(w_{k+\frac{1}{2}}), \quad (7)$$

where $f_s(w) = \mathbb{E}_{(x, \hat{y}) \sim \mathcal{D}_{k,U}} F(w, x, \hat{y})$ and $\mathcal{D}_{k,U} = \{(x_j, \hat{y}_j) | x_j \in \mathcal{D}_U \text{ and } a_j > C_k\}$ is the high-confidence pseudo-labeled dataset. As a result, the global model w_{k+1} trained on both labeled and unlabeled data is obtained. These four steps are executed repeatedly until the model converges.

C. Motivation for the Design of Ada-FedSemi

In this section, we conduct several sets of experiments to present the motivation of federated model training with the mixture of labeled and unlabeled data (Figs. 2-3), and analyze the impacts of participating fraction (*i.e.*, P) and confidence threshold (*i.e.*, C) (Figs. 4-5) on training performance. For the sake of simplicity, P and C without subscript are used to indicate that their values keep fixed during the training. We train VGG9 models for 600 rounds on the benchmark dataset CIFAR10 with 20 clients, and the details of experiments are introduced in Section IV.

We begin with experiments on labeled datasets with different scale and an unlabeled dataset. The first experiment involves 10,000 labeled data samples and 40,000 unlabeled data samples, while the second and the third experiments separately involve 20,000 and 10,000 labeled data samples without unlabeled data. The confidence threshold for pseudo-labeling is set as 0.8. We present their test accuracy and training loss in Figs. 2 and 3, respectively. In Fig. 2, with the increasing number of labeled data from 10,000 to 20,000, the performance of the trained model is improved significantly. Since obtaining the ground truth labels of data is always costly, we find that utilizing the unlabeled data can also achieve higher test accuracy, compared with training only on labeled data. This motivates us to exploit massive in-cloud unlabeled data when the scale of labeled data is limited.

However, in Fig. 3, we observe that the trend of training loss is not consistent with that of test accuracy. For example, given the same scale of labeled data, the model trained with

additional unlabeled data achieves worse training loss but higher test accuracy. Since the ultimate goal of model training is to make predictions for unseen data, training loss is not a good metric to measure the performance of a trained model. Instead, we should adopt a validation dataset to evaluate the training model and guide the decisions of P and C .

Moreover, we conduct another set of experiments on 10,000 labeled data to analyze the impacts of P and C on training performance. In Fig. 4, we present the time cost and communication cost to achieve 60% accuracy given different values of P . When achieving the same test accuracy, larger P always leads to faster convergence but results in more communication cost. Thus, we should determine P carefully to balance the trade-off between training cost and model accuracy regarding the desired cost preference.

Furthermore, we conduct experiments with different values of P and C to analyze their combined influence. Note that $C = 1.0$ means the model is only trained on labeled data since the confidence of pseudo-labels cannot exceed 1.0. The results are shown in Fig. 5, where the missing test accuracy indicates the model fails to converge given the corresponding values of P and C . We observe that, with small P (*e.g.*, 0.1), the test accuracy is more sensitive to the changes of C than that with large P (*e.g.*, 0.2 and 1.0). For example, when training with $P = 0.1$ and $C \leq 0.98$, the test accuracy degrades from 69.4% ($C = 0.98$) to 51.8% ($C = 0.9$), and the model fails to converge when $C \leq 0.9$. Conversely, when all clients participate in the training (*i.e.*, $P = 1.0$), the test accuracy is more robust to the change of C . This set of experiments shows that the values of P and C should be optimized simultaneously so as to achieve satisfactory model accuracy.

D. Problem Definition

The above experiments show the effectiveness of utilization of unlabeled data as well as impacts of two critical parameters, *i.e.*, P and C . The value of P determines how many clients will be included in the federated model training and also controls the scale of labeled dataset. Meanwhile, the value of C determines how many pseudo-labeled data are adopted to train the model, *i.e.*, the scale of unlabeled dataset. The optimization of these two parameters will have a significant impact on model performance and training cost. We formally describe the problem in the following.

In FL systems, the clients are usually equipped with limited and heterogeneous capabilities of computation and commu-

nication. Let $t_{k,m}$ denote the time cost of client v_m at round k , which includes the time for model broadcasting, updating and uploading. Since the operations of computation and communication at clients can be executed in parallel, the time cost of clients depends on the slowest participating client (*i.e.*, the straggler). Thus, the time cost at round k can be calculated as:

$$t_k = \max_{v_m \in \mathcal{V}_k} \{t_{k,m}\} + t_{k,p}, \quad (8)$$

where $t_{k,p} = t_p |\mathcal{D}_{k,U}|$ is the time cost for model training on the in-cloud pseudo-labeled data, and t_p is the time for processing a single data sample at the PS. We ignore the time cost for generating pseudo-labels, which will be elaborated in Section III.

Since the clients are usually connected with the PS via cellular network, with the increasing number of participating clients, the network may get congested and the communication cost will increase. Given the size, *i.e.*, W , of the local model (it is reasonable to assume that the sizes of local models across different clients are the same), the total communication cost can be expressed as:

$$b_k = \lceil P_k M \rceil W, \quad (9)$$

where $\lceil P_k M \rceil$ is the number of participating clients at round k . As shown in Fig. 4 and also demonstrated in the work [21], more communication cost usually leads to faster model convergence, *i.e.*, less time cost, and vice versa. Considering that different FL tasks have diverse cost preferences (*e.g.*, fast convergence or low communication cost), we consider the weighted cost of the both as in [21]:

$$\Phi_k = \alpha t_k + (1 - \alpha) b_k, \quad (10)$$

where $\alpha \in [0, 1]$ is the bias factor to adjust the preference towards time cost and communication cost. $\alpha = 0$ means that only the communication cost is taken into consideration while $\alpha = 1$ indicates that the model is expected to converge as fast as possible without considering communication cost. Note that the setting of cost preference is based on the properties of FL tasks [21]. For example, in a cellular network, traffic consumption is probably the main concern for the clients participating in FL. In contrast, in a search-and-rescue task which aims to collaboratively learn a search scheme as quickly as possible, achieving timely result would be the first priority. Thus, the cost preferences are mainly determined by the requirements of FL tasks and our algorithm can adapt to different cost preferences online, which is demonstrated in Section IV.

In the most prior works [32], [34], the optimization objective of model training aims to minimize the loss function over training data, *i.e.*, the objective defined in Eq. (3). However, as demonstrated in Section II-C, training loss fails to exactly evaluate the prediction ability of the model on unseen data (*i.e.*, generalization ability), especially when the scale of training dataset is varying. Instead, the validation dataset can be used to provide an unbiased evaluation of the model during the training [35]. The main difference between test dataset and validation dataset is that the test dataset can be dropped without affecting the model training while the validation dataset is used to guide the training process, *e.g.*, the decision of P_k and C_k in this paper. At round k , we denote the accuracy

of the global model on the validation dataset as u_k .

As stated in Section II-C, P_k and C_k have a significant impact on the performance of model training and need to be judiciously determined. As a result, the optimization problem is formulated as follows:

$$\begin{aligned} \min_{P_k, C_k, K} \quad & \sum_{k=1}^K \alpha t_k + (1 - \alpha) b_k \\ \text{s.t.} \quad & \begin{cases} u_K \geq \epsilon, \\ \sum_{k=1}^K t_k \leq T, \\ \sum_{k=1}^K b_k \leq B, \\ P_k, C_k \in [0, 1], \forall k, \end{cases} \end{aligned} \quad (11)$$

where ϵ is the target accuracy on validation dataset. T and B are the time and communication budgets for federated model training, respectively.

III. ALGORITHM DESCRIPTION

Since P_k and C_k play an important role in the model training, we propose an adaptive federated semi-supervised learning system (termed as Ada-FedSemi) to utilize both on-device labeled data and in-cloud unlabeled data efficiently. Specifically, given the desired cost preference and limited resource budgets, Ada-FedSemi employs a multi-arm bandit (MAB) based online algorithm to adaptively determine the participating fraction (*i.e.*, P_k) of clients and the confidence threshold (*i.e.*, C_k) of pseudo-labels at each round.

A. Overall Training Process of Ada-FedSemi

The overall training process of Ada-FedSemi is described in Alg. 1. Our goal is to achieve the target model accuracy while minimizing training cost. At round k , based on the value of P_k , the PS first randomly selects a subset of clients to participate in FL, and then aggregates the local models to derive the intermediate model $w_{k+\frac{1}{2}}$ at the end of local updating (Line 4-7). On the basis of the intermediate model, the teacher model \tilde{w}_k is updated (Line 8).

As suggested in [36], averaging the models across different rounds can generate a more accurate and reliable model than directly using the latest model. This is because models tend to forget past learned knowledge and fit the recent training data [37]. For example, at a certain round, if only one client is chosen to participate in the federated model training and the data on that client is highly skewed, *e.g.*, all data belong to only one class, the trained model will prefer to classify the input data as that class. Thus, we adopt the exponential moving average of the intermediate models across rounds as the teacher model, which is updated as follows and can achieve reliable performance improvement:

$$\tilde{w}_k = \gamma w_{k+\frac{1}{2}} + (1 - \gamma) \tilde{w}_{k-1}, \quad (12)$$

where $\gamma \in (0, 1]$. Unlike some existing methods, where the teacher is a well-trained model, our teacher model will be gradually improved during the training process without incurring additional training cost [36].

Subsequently, the teacher model \tilde{w}_k is used to generate pseudo-labels for unlabeled data (Line 9-10). As generating predictions for massive unlabeled data is time-consuming,

we propose two strategies to reduce the cost for pseudo-labeling. (1) Pseudo-labeling can be executed periodically (e.g., every R rounds) since the prediction ability of the teacher model will not improve significantly in several successive rounds. Furthermore, (2) pseudo-labeling can be performed in parallel with other steps like model training, broadcasting and uploading. As a result, the time cost of pseudo-labeling can be ignored. In terms of the threshold C_k , we select the high-confidence pseudo-labeled data, upon which the intermediate model is further trained to produce the global model w_{k+1} for next round (Line 11-12).

Since it is inevitable to generate incorrect pseudo-labels for unlabeled data, the model trained on these data will accumulate errors (also known as *confirmation bias*) [38]. In other words, the model keeps learning from incorrect pseudo-labels, and thereby the confidence of wrong predictions by the model continuously increases. In order to prevent error accumulation, we propose to adjust the learning rate periodically (Line 13). Specifically, we use the cosine anneal learning rate [39] to schedule the training process, which can help models jump out of local optimum and explore other regions [40]. Concretely, the learning rate is scheduled as follows:

$$\eta_k = \eta_{min} + \frac{1}{2}(\eta_{max} - \eta_{min})(1 + \cos(\frac{k \bmod \tilde{K}}{\tilde{K}}\pi)), \quad (13)$$

where η_{min} and η_{max} are the minimum and the maximum learning rates, respectively. \tilde{K} is the restart interval and k is the current training round. In each interval, the learning rate is initialized as η_{max} at the beginning and then gradually decreased to η_{min} . As a result, models trained with the trick of learning rate restart always achieve better accuracy as demonstrated in Section IV.

B. MAB based Decision Making

To adapt to system dynamics and different cost preferences, we need to adaptively determine the values of P_{k+1} and C_{k+1} (Line 14-17 of Alg. 1). As shown in Section II-C, more participating clients (i.e., larger P) always contribute to higher model accuracy but also result in more training cost. Besides, with the increasing value of C , the number of selected data decreases and the quality of pseudo-labels increases, since the high-confidence pseudo-labels are more likely to be the ground-truth labels in comparison to the low-confidence pseudo-labels. Moreover, the values of P and C are demonstrated to have an impact on each other. With more clients participating in FL (i.e., larger P), the quality of predictions by the teacher model will increase and thus the number of errors in pseudo-labels will be reduced, which will affect the decision of C . Therefore, P and C are expected to get optimized simultaneously.

However, due to the complex influence factors of federated model training (e.g., model architecture, datasets, optimizer and number of clients), it is infeasible to obtain the optimal values of P and C in advance of the training. Therefore, we propose a multi-armed bandit (MAB) based online learning algorithm to determine P and C without any prior knowledge of the FL system. In each round, the MAB algorithm chooses an action, i.e., arm, from an action set and collects a reward.

Algorithm 1 Training Process of Ada-FedSemi

Input: Client sets \mathcal{V} with their on-device data and in-cloud unlabeled data

Output: The well-trained model w

- 1: Initialize MAB agents with action sets S_P and S_C , global model w_k , participating fraction P_k , confidence threshold C_k , update interval R , target accuracy ϵ , current accuracy $u_k = 0$, $k = 1$
 - 2: **while** $u_k < \epsilon$ **do**
 - 3: **Processing at the Parameter Server**
 - 4: Select $\lceil P_k M \rceil$ clients into the set \mathcal{V}_k randomly
 - 5: Broadcast w_k to $v_m \in \mathcal{V}_k$
 - 6: Collect local models $w_{k,m}$, $\forall v_m \in \mathcal{V}_k$
 - 7: Obtain the intermediate model $w_{k+\frac{1}{2}}$ as Eq. (5)
 - 8: Update the teacher model \tilde{w}_k as Eq. (12)
 - 9: **if** $k \bmod R = 0$ **then**
 - 10: Update pseudo-labels of unlabeled data
 - 11: Select pseudo-labeled data into $\mathcal{D}_{k,U}$ based on C_k
 - 12: Train $w_{k+\frac{1}{2}}$ on $\mathcal{D}_{k,U}$ as Eq. (7), and obtain w_{k+1}
 - 13: Adjust learning rate η_k as Eq. (13)
 - 14: Validate $w_{k+\frac{1}{2}}$ and record the accuracy as $u_{k,P}$
 - 15: Validate w_{k+1} and record the accuracy as $u_{k,C}$
 - 16: Calculate the change of accuracy $\Delta u_{k,P}$ and $\Delta u_{k,C}$
 - 17: Update MAB agents, and determine P_{k+1} and C_{k+1}
 - 18: $u_k = u_{k,C}$, $k = k + 1$
 - 19: **Processing on Each Client** v_m
 - 20: **if** Receive w_k from the PS **then**
 - 21: Update the local model as Eq. (4)
 - 22: Upload the trained model $w_{k,m}$
-

Then, based on the action decisions and the corresponding rewards across different rounds, the strategies for choosing actions are updated.

We take the optimization problem in Eq. (11) as a classic MAB problem, where the values of P and C can be regarded as actions. The MAB algorithm is originally developed for discrete decision spaces. However, the values of P and C are continuous, which are within $[0, 1]$. Thus, we need to partition the continuous decision space into discrete action sets S_P and S_C , respectively. In fact, the decision space can be further zoomed in to focus on a much smaller range. For example, in Fig. 5, we find that the test accuracy of the trained model always decreases when the value of C is below 0.8. Therefore, we only consider the decision space of C in the range of 0.8 to 1.0. The decision making process of our MAB algorithm is summarized in Alg. 2. In each round, the MAB agent at the PS first makes the decision about which action is performed and then obtains a reward in response to the action (Line 3-4). According to the rewards, the MAB agent updates probabilities for the corresponding actions (Line 5-6).

We adopt a validation dataset to evaluate the accuracy of the trained models and calculate the rewards (i.e., accuracy improvement) of different actions. At round k , we denote the accuracy of the models $w_{k+\frac{1}{2}}$ and w_{k+1} as $u_{k,P}$ and $u_{k,C}$, respectively. The accuracy improvement of the two models is

$\Delta u_{k,P}$ and $\Delta u_{k,C}$:

$$\Delta u_{k,P} = u_{k,P} - u_{k-1,C}, \quad (14)$$

$$\Delta u_{k,C} = u_{k,C} - u_{k,P}. \quad (15)$$

Since the intermediate model $w_{k+\frac{1}{2}}$ is aggregated from local models, we recognize the improvement of this model, *i.e.*, $\Delta u_{k,P}$, as the outcome of the decision of P_k . Meanwhile, the model w_{k+1} is trained on the pseudo-labeled data selected by C_k , and thus we recognize $\Delta u_{k,C}$ as the outcome of the decision of C_k . Since the algorithm for determining P_k and C_k is the same, we use Δu_k and S for simplicity to introduce our algorithm in the following. We define the reward of the decision at round k as follows:

$$r_k = \begin{cases} \frac{\Delta u_k}{\Phi_k}, & \text{if } \Delta u_k \geq 0, \\ \Delta u_k \cdot \Phi_k, & \text{otherwise.} \end{cases} \quad (16)$$

As the reward design is the key to the success of MAB algorithms [41], herein, we explain the rationality of the reward function. The intuition of our reward design has two folds and the goals of the two-fold reward functions are consistent, *i.e.*, improving the model performance in a cost-efficient way. (1) When achieving the same accuracy improvement (*i.e.*, $\Delta u_k \geq 0$), the decisions which consume less training cost should be given higher rewards. In another word, we expect high accuracy improvement and small training cost. (2) While some inappropriate actions may degrade accuracy, *i.e.*, $\Delta u_k < 0$, and we still use $\Delta u_k/\Phi_k$ as the reward, a smaller training cost Φ_k will lead to a higher penalty (penalty means negative reward). This is not consistent with our design goal, *i.e.*, efficient training. Thus, in case of $\Delta u_k < 0$, we denote $\Delta u_k \cdot \Phi_k$ as the reward.

Traditional MAB algorithms estimate the actual reward of an action by averaging its received rewards across rounds. However, in this paper, the reward distribution of actions is not identical across different rounds. Firstly, the improvement speed of model accuracy is not the same during the training process. In general, the increase of model accuracy is fast at the beginning of the training and becomes slow as training progresses. Besides, the optimal decision may change over time since the quality of pseudo-labels will improve and the cost preference may vary during the training. Therefore, this is a non-stationary MAB problem [42], and it is not rational to simply average rewards of each action across rounds as traditional MAB algorithms do. Instead, we concentrate more on the recent rewards which are assigned with larger weights, and gradually decay the weights for the past rewards [43]. At round k , for each action a in the action set S , its estimated reward $\hat{r}_{k,a}$ is calculated as follows:

$$\hat{r}_{k,a} = \begin{cases} \hat{r}_{k-1,a} + \beta(r_k - \hat{r}_{k-1,a}) & , \text{ if } a_k = a, \\ \hat{r}_{k-1,a}, & \text{otherwise,} \end{cases} \quad (17)$$

where $\beta \in (0, 1]$ is the decay factor and a_k is the action chosen for round k .

Thus, the goal of our MAB algorithm is to maximize the total received rewards via a judicious trade-off between exploration and exploitation. Exploitation means pulling the best action known so far while exploration aims to explore different actions to find better solutions. Specifically, we adopt the Boltzmann exploration strategy [44], which is widely used

Algorithm 2 MAB Agent

- 1: Action set S and $\hat{r}_{k,a}$, $p_{k,a}$ for each action
 - 2: **for** round $k \in \{1, \dots, K\}$ **do**
 - 3: Select an action based on the probability $p_{k,a}$, $\forall a \in S$
 - 4: Receive the reward r_k based on Eq. (16)
 - 5: Update estimated rewards $\hat{r}_{k,a}$ as Eq. (17)
 - 6: Calculate the probability $p_{k,a}$ of action as Eq. (18)
-

for balancing exploration and exploitation. The probability of choosing action $a \in S$ at round k is calculated as follows:

$$p_{k,a} = \frac{e^{\psi \hat{r}_{k,a}}}{\sum_{a' \in S} e^{\psi \hat{r}_{k,a'}}}. \quad (18)$$

Particularly, with $\psi = 0$, the actions are uniformly chosen all the time, while $\psi \rightarrow \infty$ means that the MAB agent will always output the action with the highest reward without any exploration.

IV. EXPERIMENTAL EVALUATION

A. System Platform

We evaluate the performance of Ada-FedSemi through extensive experiments on an FL hardware prototype system. Specifically, an AMAX deep learning workstation, which is equipped with an Intel(R) Core(TM) i9-10900X CPU, 4 NVIDIA GeForce RTX 2080Ti GPUs and 128 GB RAM, is applied to serve as the PS. Besides, 20 NVIDIA Jetson TX2 developer kits are specified as the clients. The PS and clients are connected via a Wi-Fi router. The implementation for model training is based on the PyTorch deep learning framework [45], and we use the socket library of Python to build up the communication between clients and the PS.

B. Setup of Experiments

Datasets and Models: We use two benchmark datasets, *i.e.*, CIFAR10 [46] and SVHN [47], which are commonly adopted in semi-supervised learning [18]–[20], [48], to evaluate the performance of Ada-FedSemi and baselines:

- **CIFAR10:** It contains 60,000 color images labeled in 10 classes with 6,000 samples per class. By default, we split the whole dataset into four datasets: i) labeled training dataset with 10,000 samples, ii) unlabeled training dataset with 40,000 samples whose labels are discarded, iii) validation dataset with 2,000 samples, and iv) test dataset with 8,000 samples.
- **SVHN:** There are 73,257 digits for training, 26,032 digits for testing, and 531,131 additional data, which are labeled in 10 classes. By default, 5% of training data, *i.e.*, 3,660 digits, are distributed to clients as labeled data. 20,000 and 6,032 digits in testing dataset are used for testing and validation, respectively. The rest digits in the training dataset and additional dataset are all placed at the PS as unlabeled data. As a result, there are 600,728 unlabeled samples at the PS.

Since data are not always distributed uniformly across clients at the network edge, we will analyze training performance under both IID and non-IID settings. (1) In the IID

TABLE II: Performance metrics of FedSemi- P - C and FedAvg- P (i.e., FedSemi- P -1.0) on CIFAR10.

(a) Test accuracy (%)					(b) Time cost (min) for accuracy of 75%					(c) Comm. cost (GB) for accuracy of 75%				
$C \backslash P$	0.1	0.2	0.5	1.0	$C \backslash P$	0.1	0.2	0.5	1.0	$C \backslash P$	0.1	0.2	0.5	1.0
0.6	-	-	74.7	75.9	0.6	-	-	-	419	0.6	-	-	-	355
0.7	-	80.3	77.4	82.9	0.7	-	382	312	251	0.7	-	68	139	226
0.8	-	81.4	82.7	83.4	0.8	-	311	240	196	0.8	-	58	113	186
0.9	69.7	82.0	83.3	83.1	0.9	-	296	199	186	0.9	-	57	98	183
0.95	74.4	81.5	81.9	81.5	0.95	-	247	211	190	0.95	-	49	103	187
0.98	77.6	79.2	81.4	80.6	0.98	356	296	221	194	0.98	36	59	108	194
1.0	75.0	75.2	75.5	75.6	1.0	395	346	326	314	1.0	41	72	169	324

TABLE III: Performance metrics of FedSemi- P - C and FedAvg- P (i.e., FedSemi- P -1.0) on SVHN.

(a) Test accuracy (%)					(b) Time cost (min) for accuracy of 75%					(c) Comm. cost (GB) for accuracy of 75%				
$C \backslash P$	0.1	0.2	0.5	1.0	$C \backslash P$	0.1	0.2	0.5	1.0	$C \backslash P$	0.1	0.2	0.5	1.0
0.7	62.8	80.3	80.7	81.7	0.7	-	63	52	31	0.7	-	0.77	1.21	2.07
0.75	64.3	81.1	81.3	82.1	0.75	-	59	46	27	0.75	-	0.71	1.12	1.94
0.8	67.5	81.7	81.9	83.3	0.8	-	54	43	25	0.8	-	0.69	1.09	1.76
0.85	68.8	82.4	83.3	84.5	0.85	-	52	37	23	0.85	-	0.64	1.02	1.69
0.9	69.4	82.7	82.1	82.7	0.9	-	51	41	27	0.9	-	0.59	1.08	1.81
0.95	72.8	80.5	81.4	81.2	0.95	-	59	45	28	0.95	-	0.76	1.13	2.20
1.0	70.3	80.6	81.0	81.5	1.0	-	58	43	28	1.0	-	0.76	1.15	2.04

TABLE IV: Optimal combination of P and C with different cost preferences (α) on CIFAR10.

Range of α ($\times 0.1$)	[0, 1.07)	[1.07, 5.06)	[5.06, 8.68)	[8.68, 10]
Optimal P and C	(0.1, 0.98)	(0.2, 0.95)	(0.5, 0.9)	(1.0, 0.9)

setting, all labeled data are uniformly distributed to clients. (2) In the non-IID setting, as in [30], a fraction (ζ) of data samples assigned to a client belong to a certain class and the remaining data samples belong to other classes, which is denoted as non-IID- ζ . By default, the data distributions of CIFAR10 and SVHN are both non-IID-0.5.

For CIFAR10, we train a VGG9 model [49] with 3.49M parameters while a lightweight CNN model with 0.54M parameters is trained on SVHN. Besides, the SGD-momentum optimizer is adopted in our experiments to optimize models, and the momentum is set as 0.9. The restart interval for learning rate is set as 100. The maximum and minimum learning rates are set as 0.05 and 0.0001, respectively.

Baselines: We compare our proposed system with the following baselines.

- **FedSemi** [18], [19]: In FedSemi, the in-cloud unlabeled data and on-device labeled data are used to train models in a semi-supervised way. However, the two critical parameters, i.e., P and C , are fixed during the training. Given different combinations of P and C , we denote the baselines as FedSemi- P - C , e.g., FedSemi-0.5-0.9.
- **FedAvg** [4]: In FedAvg, only labeled data on clients are utilized to train models, and thus there is only one critical parameter, i.e., P . We denote the FedAvg with different P as FedAvg- P , e.g., FedAvg-0.2. Note that if the value

TABLE V: Optimal combination of P and C with different cost preferences (α) on SVHN.

Range of α ($\times 0.01$)	[0, 2.98)	[2.98, 4.57)	[4.57, 100]
Optimal P and C	(0.2, 0.9)	(0.5, 0.85)	(1.0, 0.85)

of C in FedSemi- P - C is set as 1.0 (i.e., FedSemi- P -1.0), FedSemi- P -1.0 is equivalent to FedAvg- P , since the confidence of pseudo-labels cannot exceed 1.0 and none of the unlabeled data is selected. For ease of presentation, we will use FedSemi- P -1.0 and FedAvg- P interchangeably in the later experiments.

Performance Metrics: In the experiments, we employ the following metrics to evaluate the performance of different FL systems: (1) Test accuracy. In each round, we will evaluate the global model on test dataset and record the accuracy. (2) Time cost. We will record the time to achieve the target test accuracy on different FL systems. (3) Communication cost. The communication cost for broadcasting and uploading models is also recorded when achieving the target test accuracy. (4) Weighted Cost. Based on the cost preference and Eq. (10), we combine time cost and communication cost to derive the weighted cost.

C. The Impacts of P and C

We first conduct experiments on the baselines with fixed P and C to analyze the impacts of the two parameters. The time budget of the training is set as 480min and 80min for CIFAR10 and SVHN, respectively, and the target test accuracy is set as 75%. By default, in each round of SVHN, we randomly

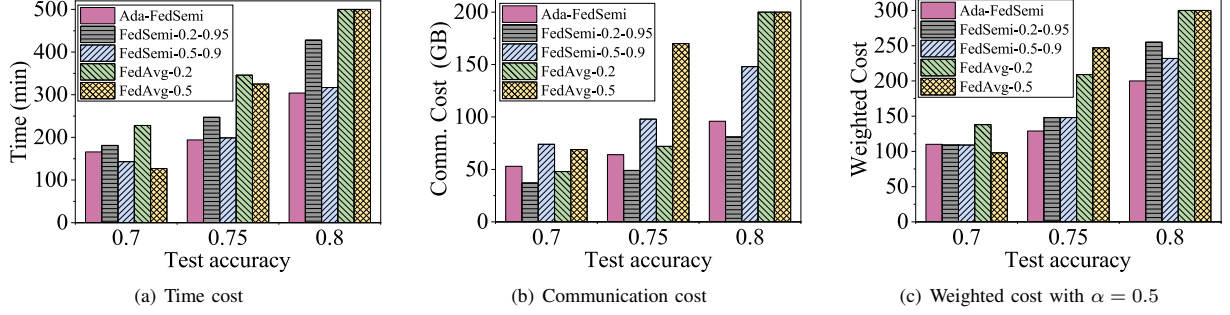


Fig. 6: Training cost on CIFAR10.

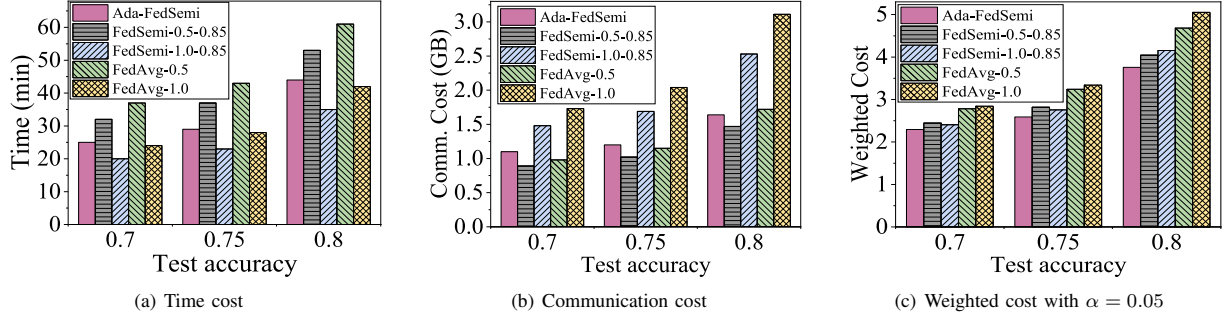


Fig. 7: Training cost on SVHN.

select 100,000 out of 600,728 unlabeled data to generate pseudo-labels for training efficiency and we will analyze the impact of the scale of unlabeled dataset in Section IV-F. The experimental results on CIFAR10 and SVHN are presented in Tables II and III, respectively.

From the perspective of test accuracy, the models trained on both labeled and unlabeled data (*i.e.*, FedSemi) can achieve up to 7.8% accuracy improvement on CIFAR10 and 3.0% accuracy improvement on SVHN, compared with the models only trained on labeled dataset (*i.e.*, FedAvg). However, we observe that FedSemi with small C and/or small P suffers from accuracy degradation (*e.g.*, $P = 0.1$ and $C \leq 0.8$ on SVHN) and even fails to converge (*e.g.*, FedSemi-0.1-0.8 on CIFAR10). Given a small P , the models will only learn knowledge from a small number of clients and labeled data, and a small C will bring in many low-confidence pseudo-labels during the training, which degrades the performance of models. Besides, the values of P have influence on the optimal values of C . For example, on CIFAR10, when all clients participate in the FL training (*i.e.*, $P = 1.0$), $C = 0.8$ achieves the highest test accuracy. However, when selecting only 10% of clients (*i.e.*, $P = 0.1$), C needs to be set as 0.98 to achieve the best accuracy. The reason lies in that the models can learn more knowledge from the labeled data with the increasing of P and thus generate pseudo-labels with higher confidence. Thus, when using the same C , we can select more samples with less errors from the pseudo-labeled data.

In terms of training cost, the time cost and communication cost are usually contradictory. With the increasing number of participating clients, the time to achieve the target accuracy

gets shorter and meanwhile the communication cost gets higher. For example, when training models on SVHN using FedAvg, with the value of P varying from 0.2 to 1.0, the time cost to achieve 75% test accuracy decreases from 58min to 28min while the communication cost increases from 0.76GB to 2.04GB. It is noteworthy that training on both labeled and unlabeled data may not always achieve better training efficiency, compared with training only on labeled dataset. For example, on CIFAR10, FedSemi-1.0-0.6 spends 33% more time cost and 10% more communication cost to achieve the same target test accuracy in comparison to FedAvg-1.0. This is because a large number of low-confidence pseudo-labels mislead the optimization of model training and result in resource waste.

In Tables II and III, we observe that a little extra time cost can help reduce the communication cost to a great extent in some cases. For example, on CIFAR10, when C is set as 0.9, the system with $P = 0.5$ spends 7% more time but saves 47% communication cost, compared with $P = 1.0$. Different FL tasks always have different cost preferences. Some tasks expect to converge fast without considering communication cost while others may prefer to perform model training in a communication-efficient way. As a result, these preferences will lead to different optimal decision of P and C . We present the optimal combination of P and C under different preferences in Tables IV and V. For example, on CIFAR10, when the task prefers saving communication cost, *i.e.*, $\alpha < 0.107$, the minimum weighted cost can be achieved with $P = 0.1$. When the task is expected to converge fast (*i.e.*, $\alpha \geq 0.868$), all clients (*i.e.*, $P = 1.0$) should participate in the training.

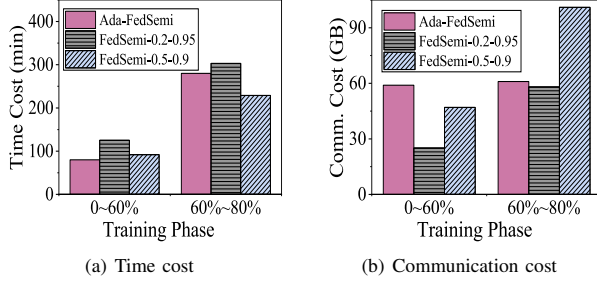


Fig. 8: Training cost on CIFAR10 with varying preference.

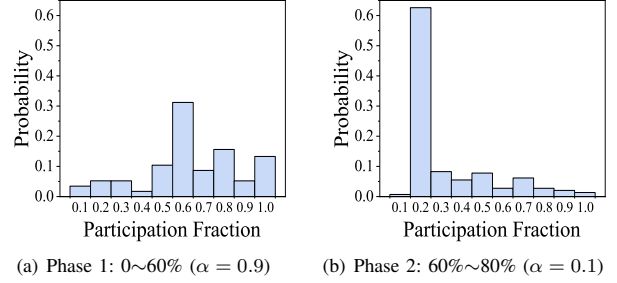


Fig. 9: Distribution of P in two training phases.

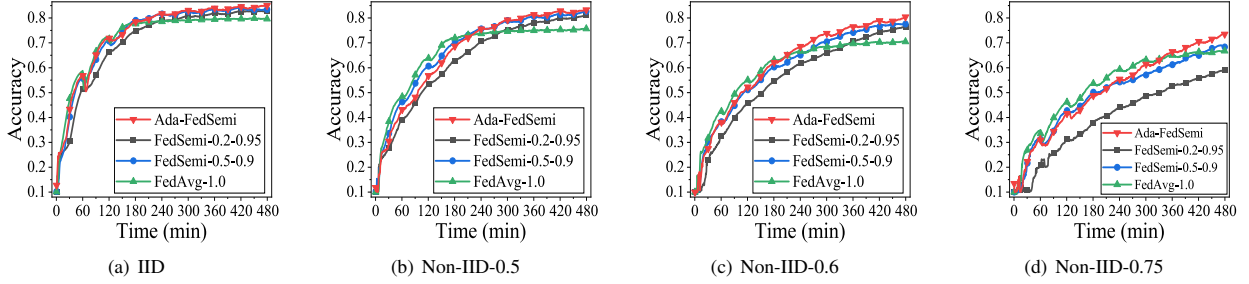


Fig. 10: Training process of Ada-FedSemi and the baselines on different data distributions.

D. Performance Comparison

In this section, we compare the performance of Ada-FedSemi and baselines. In Ada-FedSemi, the decision spaces of P and C are set as $[0.0, 1.0]$ and $[0.8, 1.0]$, respectively, and we evenly partition each of the two decision spaces into 10 discrete values. By default, we set the preference parameter α as 0.5 and 0.05 for the training on CIFAR10 and SVHN, respectively. As indicated in Table IV, on CIFAR10, $P = 0.2$ and 0.5 can achieve small weighted cost when $\alpha = 0.5$. Therefore, we choose FedAvg-0.2, FedAvg-0.5, FedSemi-0.2-0.95 and FedSemi-0.5-0.9 as baselines. Similarly, on SVHN, we take FedAvg-0.5, FedAvg-1.0, FedSemi-0.5-0.85 and FedSemi-1.0-0.85 for comparison. When a system cannot achieve the target test accuracy, its cost is set as the maximum (*i.e.*, the budgets are exhausted). For models trained on CIFAR10 and SVHN, the time budget is set as 500min and 80min, respectively.

The time cost, communication cost and weighted cost of Ada-FedSemi and the baselines for achieving the different test accuracy are presented in Figs. 6-7. Although baselines with fixed P and C may achieve the least time cost (*i.e.*, $P = 1.0$ on SVHN) or communication cost (*i.e.*, $P = 0.2$ on CIFAR10) in some cases, they cannot achieve the least weighted cost. On the contrary, Ada-FedSemi can always achieve the least weighted cost, indicating that Ada-FedSemi is able to balance time cost and communication cost given the specific cost preference. On CIFAR10, compared with FedAvg- P , Ada-FedSemi can save 35% ($P = 0.2$) and 48% ($P = 0.5$) weighted cost when achieving 75% test accuracy. Given the 80% test accuracy, Ada-FedSemi saves the weighted cost over FedSemi-0.2-0.95 and FedSemi-0.5-0.9 by 22% and 14%. However, FedAvg- P fails to achieve higher test accuracy (*i.e.*, 80%) without utilization of the in-cloud unlabeled data. On

SVHN, our algorithm saves 6%-25% weighted cost, compared with different baselines.

Moreover, we conduct another set of experiments on CIFAR10 in the scenario where the FL tasks would like to achieve an acceptable test accuracy with low time cost and then further improve model performance in a communication-efficient way. Specifically, we initially set α as 0.9 to ensure fast convergence and when the test accuracy reaches 60%, α is set as 0.1 to put more emphasis on communication cost. The training process is terminated when the test accuracy reaches 80%. The communication cost and time cost of the two training phases are presented in Fig. 8.

In the first training phase, Ada-FedSemi achieves the least time cost, resulting in the most communication cost since our goal in this phase is fast convergence. In the second training phase, our system results in similar communication cost with FedSemi-0.2-0.95 and saves about 50% communication cost in comparison to FedSemi-0.5-0.9. We further present the distribution of the values of P in the two training phases. As shown in Fig. 9, Ada-FedSemi always chooses the optimal P (*i.e.*, 0.6 in the first phase and 0.2 in the second phase) with the high probability, which indicates that our system is able to adaptively determine the optimal combination of P and C even when the cost preference is varying over time.

E. Adaptability to Data Distribution

In this section, we conduct experiments to evaluate the impact of data distributions. The models are trained on CIFAR10 with four different data distributions, *i.e.*, IID, non-IID-0.5, non-IID-0.6 and non-IID-0.75. We take FedSemi-0.2-0.95, FedSemi-0.5-0.9 and FedAvg-1.0 for comparison. The time budget is set as 480min, and the experimental results are presented in Fig. 10.

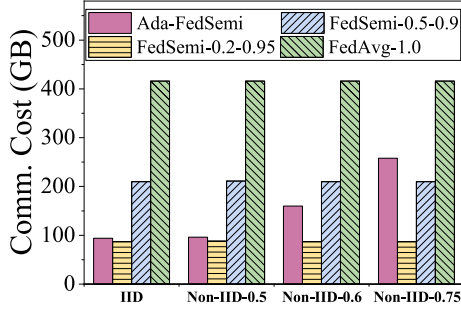


Fig. 11: Communication cost of different systems.

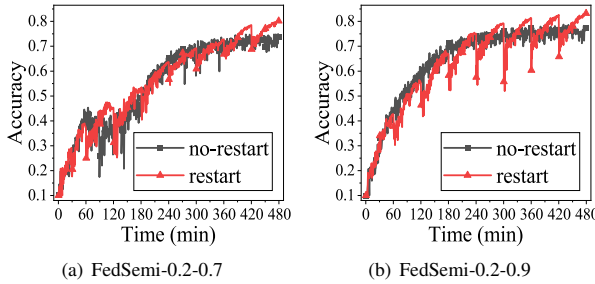


Fig. 13: Impact of learning rate restart.

Generally, the test accuracy of the models trained on all systems decreases with the increasing skewness of data distribution. The final test accuracy of Ada-FedSemi on IID, non-IID-0.5, non-IID-0.6 and non-IID-0.75 is 85.0%, 83.3%, 80.4%, and 73.5%, respectively. For the same data distribution, Ada-FedSemi can always achieve the highest test accuracy and outperform the three baselines by 1.5% to 14.8%. We find that FedAvg-1.0 always achieves the best test accuracy at the beginning of the training. This is because, in FedAvg, models are optimized without perturbation of errors from low-confidence pseudo-labels and thus can converge fast. For FedSemi, as the training progresses, the prediction ability of models is improved, and thus the teacher model generates more and more high-confidence pseudo-labels for the unlabeled data. As a result, the models trained on these high-confidence pseudo-labeled data can achieve higher test accuracy, compared with FedAvg. Nevertheless, in Fig. 10(d), the test accuracy of FedSemi-0.2-0.95 and FedAvg-1.0 is separately 59.0% and 66.7%, indicating that with highly skewed data, a small number of clients may fail to generate high-confidence pseudo-labels, and thus the incorrect pseudo-labels may mislead the model optimization.

In Fig. 11, we also present the communication cost of different systems on the four data distributions. Except Ada-FedSemi, the communication cost of other three baselines is almost the same across different data distributions since they always use fixed P and C . As the skewness of training data increases, our system can adaptively increase communication cost to ensure the best model performance. For example, on non-IID-0.75 dataset, although Ada-FedSemi consumes 20% more communication cost than FedSemi-0.5-0.9, it improves

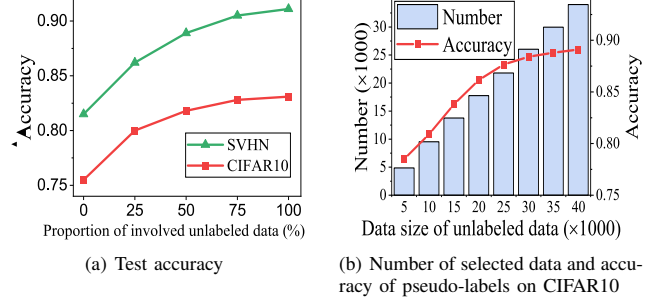


Fig. 12: Performance of models trained with different scale of unlabeled data.

the final test accuracy from 66.7% to 73.5%. The results of the experiments demonstrate that Ada-FedSemi has the ability of adapting to different data distributions.

F. The Impact of the Scale of Unlabeled Dataset

As mentioned in Section IV-B, the number of unlabeled data in CIFAR10 and SVHN is 40,000 and 600,728, respectively. To explore the impact of the scale of unlabeled dataset, we fix the size of labeled dataset and conduct experiments by changing the proportion of in-cloud unlabeled data involved in the training, and the results are shown in Fig. 12(a). When the proportion of involved unlabeled data varies from 0% to 100%, the accuracy of models on CIFAR10 increases from 75.5% to 83.1% and that on SVHN increases from 81.5% to 91.1%. Furthermore, we use the trained model with $C = 0.9$ to select pseudo-labeled data of CIFAR10 and the results are presented in Fig. 12(b). With the increasing of scale of unlabeled data, the number of selected data increases from 4,868 to 34,008, and the accuracy of pseudo-labels increases from 78.5% to 89.1%. This set of experiments demonstrates that in Ada-FedSemi, the final test accuracy of trained models is positively correlated to the scale of unlabeled data. Therefore, when the scale of labeled data on clients is small, it is an effective way to collect and exploit large-scale in-cloud unlabeled data to boost the model performance, and it will not incur additional training cost for resource-constrained clients.

G. Comparison of Methods for Online Optimization

In fact, after formulating the problem as Eq. (11) and designing the reward function in Eq. (16), we can adopt different online optimization methods to find the optimal P and C . MAB-based algorithms [50] and Bayesian optimization [51] are powerful tools to make decisions online under uncertainty. Herein, we compare training cost, model accuracy and the time for decision making when Boltzmann exploration, UCB, GP-UCB and Bayesian optimization are adopted in our system. The experimental results are presented in Table VI. Training time and traffic consumption are compared when models achieve 80% accuracy and the time for decision making is accumulative during the whole training process. Four methods achieve similar training performance while it takes much more time to make decisions by Bayesian optimization than others. Nevertheless, the time for decision making can be

TABLE VI: Comparison of different methods for adaptive model training on CIFAR10 with $\alpha = 0.5$. (BE: Boltzmann exploration; BO: Bayesian optimization.)

Metrics	Model Training			Decision Making
	Time (min)	Traffic (GB)	Accuracy (%)	Time (min)
BE	304	96	83.3	0.03
UCB	310	95	82.9	0.02
GP-UCB	302	99	83.0	0.33
BO	305	93	83.5	2.65

always ignored, compared with the time for model training. Thus, Ada-FedSemi is compatible with different optimization methods and we adopt Boltzmann exploration for its efficiency and ease of implementation.

H. The Impact of Learning Rate Restart

Finally, we also conduct experiments to evaluate the impact of learning rate restart. The constant learning rate 0.05, denoted as no-restart, is taken as comparison. We perform model training on CIFAR10 with FedSemi-0.2-0.9 and FedSemi-0.2-0.7, and the corresponding results are presented in Figs. 13(a) and 13(b), respectively. It shows that the test accuracy first degrades at each moment of learning rate restart and then resumes quickly. Although the models trained with the constant learning rate can achieve continuous improvement, it converges earlier and fails to reach higher test accuracy, compared with the models trained with learning rate restart. This is because the constant learning rate may make the models get trapped in local minimum, especially when there exists noise in pseudo-labels. Instead, restarting learning rate helps the models jump out of local minimum and converge to better solutions.

V. RELATED WORKS

A. Federated Learning with Labeled and Unlabeled Data

The significant improvement of AI in recent years is largely attributed to the utilization of large scale labeled dataset. However, obtaining labels of data is often very costly and time-consuming in practice [6]. Therefore, semi-supervised learning (SSL) [52] is proposed to train models on both small scale of labeled dataset and large scale of unlabeled dataset.

There are two main methods in SSL. The first one is consistency regularization based algorithms [48], [53]. These algorithms require that the predictions of unlabeled data are invariant to different data augmentations of a single data sample, which significantly increases training overhead [26]. The other method is pseudo-labeling based algorithms [20], [26], which regard the high-confidence predictions of unlabeled data as their pseudo-labels. However, all above methods only care about the final test accuracy of the trained model but don't take the features (*e.g.*, limited capacity of communication and computation) of FL into consideration and thus fail to achieve training efficiency.

Recently, several works try to perform SSL under FL settings. Some works try to exploit on-device unlabeled data. For example, Jeong *et al.* [14] propose to select other clients' local models for each client to help exploit local unlabeled data. Besides, Long *et al.* [15] adopt two networks (teacher

and student) at each client to train models on both labeled and unlabeled data. Moreover, considering the limited scale of on-device unlabeled data, works [16], [17] first distribute in-cloud unlabeled data to clients and then perform SSL algorithms. However, all above methods try to utilize unlabeled data at clients, which will increase training cost for clients. Since end devices are always resource-constrained and there are large-scale in-cloud public unlabeled data, it is more efficient to utilize these data at the PS rather than on end devices.

B. Resource-efficient Federated Learning

To achieve efficient training of FL, many algorithms are proposed to reduce time cost and communication cost. Some recent works [22], [30] aim to optimize the training time by utilizing deep reinforcement learning (DRL) to schedule clients for federated model training. However, these works all employ a fixed participating fraction of clients and mainly concentrate on optimizing a single objective (*e.g.*, training time), which cannot satisfy the various cost preferences (*e.g.*, fast convergence or low communication cost) for different FL tasks. As stated in the work [21], a large participating fraction can lead to reduction of training time while a small fraction contributes to saving communication cost. However, the work [21] is designed to determine the (offline) optimal participating fraction before performing the FL tasks, which fails to adapt to the dynamic changes of cost preferences online.

In addition, to reduce the volume of transmitted data, many works propose various compression techniques for distributed model training, *e.g.*, *Quantization* [54] and *Sparsification* [55]. Moreover, Wang *et al.* [32] propose to determine the optimal local updating steps adaptively with the constraint of available resource. Note that our work, which determines the optimal participating fraction of clients and confidence threshold of pseudo-labels adaptively, is orthogonal to compression techniques and adaptive frequency of local updating. Thus, other communication-efficient techniques can be adopted in our system to further reduce training cost.

VI. CONCLUSIONS

To fully utilize the on-device labeled and in-cloud unlabeled data in FL, we propose an adaptive FSSL system called Ada-FedSemi. It employs an MAB based online learning algorithm to adaptively determine the fraction of participating clients and confidence threshold for pseudo-labeling during the federated model training. The dynamic optimization of participating fraction and confidence threshold can contribute to the trade-off between model accuracy and training efficiency given the limited resource budgets. The extensive experimental results demonstrate that Ada-FedSemi significantly outperforms the existing baselines, including FedAvg and FedSemi.

ACKNOWLEDGMENTS

The corresponding author of this paper is Yang Xu. This research is supported in part by the National Science Foundation of China (NSFC) under Grants 62132019, 61936015, 62102391 and U1709217, and the National Key Research and Development Program of China (Grant No. 2021YFB3301501).

REFERENCES

- [1] F. Zantalis, G. Koulouras, S. Karabetsos, and D. Kandris, "A review of machine learning and iot in smart transportation," *Future Internet*, vol. 11, no. 4, p. 94, 2019.
- [2] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, P. Martinez-Gonzalez, and J. Garcia-Rodriguez, "A survey on deep learning techniques for image and video semantic segmentation," *Applied Soft Computing*, vol. 70, pp. 41–65, 2018.
- [3] A. S. Tulshan and S. N. Dhage, "Survey on virtual assistant: Google assistant, siri, cortana, alexa," in *International symposium on signal processing and intelligent recognition systems*. Springer, 2018, pp. 190–201.
- [4] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial Intelligence and Statistics*, 2017, pp. 1273–1282.
- [5] Y. Wang, Y. Tong, D. Shi, and K. Xu, "An efficient approach for cross-silo federated learning to rank," in *2021 IEEE 37th International Conference on Data Engineering (ICDE)*. IEEE, 2021, pp. 1128–1139.
- [6] K. Li, G. Li, Y. Wang, Y. Huang, Z. Liu, and Z. Wu, "Crowdrl: An end-to-end reinforcement learning framework for data labelling," in *2021 IEEE 37th International Conference on Data Engineering (ICDE)*. IEEE, 2021, pp. 289–300.
- [7] T. Weyand, I. Kostrikov, and J. Philbin, "Planet-photo geolocation with convolutional neural networks," in *European Conference on Computer Vision*. Springer, 2016, pp. 37–55.
- [8] R. Rothe, R. Timofte, and L. Van Gool, "Deep expectation of real and apparent age from a single image without facial landmarks," *International Journal of Computer Vision*, vol. 126, no. 2, pp. 144–157, 2018.
- [9] B. Neyshabur, S. Bhojanapalli, D. McAllester, and N. Srebro, "Exploring generalization in deep learning," *arXiv preprint arXiv:1706.08947*, 2017.
- [10] C. Yang, Q. Huang, Z. Li, K. Liu, and F. Hu, "Big data and cloud computing: innovation opportunities and challenges," *International Journal of Digital Earth*, vol. 10, no. 1, pp. 13–53, 2017.
- [11] N. A. Ghani, S. Hamid, I. A. T. Hashem, and E. Ahmed, "Social media big data analytics: A survey," *Computers in Human Behavior*, vol. 101, pp. 417–428, 2019.
- [12] W. Li, L. Wang, W. Li, E. Agustsson, and L. Van Gool, "Webvision database: Visual learning and understanding from web data," *arXiv preprint arXiv:1708.02862*, 2017.
- [13] S. Guo, W. Huang, H. Zhang, C. Zhuang, D. Dong, M. R. Scott, and D. Huang, "Curriculumnet: Weakly supervised learning from large-scale web images," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 135–150.
- [14] W. Jeong, J. Yoon, E. Yang, and S. J. Hwang, "Federated semi-supervised learning with inter-client consistency & disjoint learning," in *International Conference on Learning Representations*, 2021.
- [15] Z. Long, L. Che, Y. Wang, M. Ye, J. Luo, J. Wu, H. Xiao, and F. Ma, "Fedsemi: An adaptive federated semi-supervised learning framework," *arXiv preprint arXiv:2012.03292*, 2020.
- [16] S. Itahara, T. Nishio, Y. Koda, M. Morikura, and K. Yamamoto, "Distillation-based semi-supervised federated learning for communication-efficient collaborative training with non-iid private data," *IEEE Transactions on Mobile Computing*, pp. 1–1, 2021.
- [17] I. Bistriz, A. Mann, and N. Bambos, "Distributed distillation for on-device learning," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [18] N. Papernot, M. Abadi, U. Erlingsson, I. Goodfellow, and K. Talwar, "Semi-supervised knowledge transfer for deep learning from private training data," in *International Conference on Learning Representations*, 2017.
- [19] N. Papernot, S. Song, I. Mironov, A. Raghunathan, K. Talwar, and U. Erlingsson, "Scalable private learning with pate," in *International Conference on Learning Representations*, 2018.
- [20] P. Cascante-Bonilla, F. Tan, Y. Qi, and V. Ordóñez, "Curriculum labeling: Revisiting pseudo-labeling for semi-supervised learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- [21] B. Luo, X. Li, S. Wang, J. Huang, and L. Tassiulas, "Cost-effective federated learning design," in *IEEE INFOCOM 2021-IEEE Conference on Computer Communications*, 2021.
- [22] W. Xia, T. Q. Quek, K. Guo, W. Wen, H. H. Yang, and H. Zhu, "Multi-armed bandit-based client scheduling for federated learning," *IEEE Transactions on Wireless Communications*, vol. 19, no. 11, pp. 7108–7123, 2020.
- [23] K. Sohn, D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. A. Raffel, E. D. Cubuk, A. Kurakin, and C.-L. Li, "Fixmatch: Simplifying semi-supervised learning with consistency and confidence," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 596–608.
- [24] P. Zhou, Q. Lin, D. Loghin, B. C. Ooi, Y. Wu, and H. Yu, "Communication-efficient decentralized machine learning over heterogeneous networks," in *2021 IEEE 37th International Conference on Data Engineering (ICDE)*. IEEE, 2021, pp. 384–395.
- [25] A. Li, L. Zhang, J. Wang, J. Tan, F. Han, Y. Qin, N. M. Freris, and X.-Y. Li, "Efficient federated-learning model debugging," in *2021 IEEE 37th International Conference on Data Engineering (ICDE)*. IEEE, 2021, pp. 372–383.
- [26] M. N. Rizve, K. Duarte, Y. S. Rawat, and M. Shah, "In defense of pseudo-labeling: An uncertainty-aware pseudo-label selection framework for semi-supervised learning," in *International Conference on Learning Representations*, 2021.
- [27] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, et al., "Advances and open problems in federated learning," *Foundations and Trends® in Machine Learning*, vol. 14, no. 1–2, pp. 1–210, 2021.
- [28] D. Verma, S. Julier, and G. Cirincione, "Federated ai for building ai solutions across multiple agencies," *arXiv preprint arXiv:1809.10036*, 2018.
- [29] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of fedavg on non-iid data," in *International Conference on Learning Representations*, 2020.
- [30] H. Wang, Z. Kaplan, D. Niu, and B. Li, "Optimizing federated learning on non-iid data with reinforcement learning," in *IEEE INFOCOM 2020-IEEE Conference on Computer Communications*. IEEE, 2020, pp. 1698–1707.
- [31] M. Tang and V. W. Wong, "An incentive mechanism for cross-silo federated learning: A public goods perspective," in *IEEE INFOCOM 2021-IEEE Conference on Computer Communications*. IEEE, 2021, pp. 1–10.
- [32] 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 Journal on Selected Areas in Communications*, vol. 37, no. 6, pp. 1205–1221, 2019.
- [33] X. J. Zhu, "Semi-supervised learning literature survey," 2005.
- [34] P. Han, S. Wang, and K. K. Leung, "Adaptive gradient sparsification for efficient federated learning: An online learning approach," in *2020 IEEE 40th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2020, pp. 300–310.
- [35] J. Gareth, W. Daniela, H. Trevor, and T. Robert, *An introduction to statistical learning: with applications in R*. Springer, 2013.
- [36] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [37] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning (still) requires rethinking generalization," *Communications of the ACM*, vol. 64, no. 3, pp. 107–115, 2021.
- [38] E. Arazo, D. Ortego, P. Albert, N. E. O'Connor, and K. McGuinness, "Pseudo-labeling and confirmation bias in deep semi-supervised learning," in *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2020, pp. 1–8.
- [39] I. Loshchilov and F. Hutter, "Sgdr: Stochastic gradient descent with warm restarts," in *International Conference on Learning Representations*, 2017.
- [40] A. Gotmare, N. S. Keskar, C. Xiong, and R. Socher, "A closer look at deep learning heuristics: Learning rate restarts, warmup and distillation," in *International Conference on Learning Representations*, 2019.
- [41] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [42] O. Besbes, Y. Gur, and A. Zeevi, "Stochastic multi-armed-bandit problem with non-stationary rewards," *Advances in neural information processing systems*, vol. 27, pp. 199–207, 2014.
- [43] R. S. Sutton and A. G. Barto, "Reinforcement learning: An introduction," 2011.
- [44] N. Cesa-Bianchi, C. Gentile, G. Lugosi, and G. Neu, "Boltzmann exploration done right," in *Advances in Neural Information Processing Systems*, vol. 30. Curran Associates, Inc., 2017.
- [45] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga et al., "Pytorch: An imperative style, high-performance deep learning library," in *Advances in neural information processing systems*, 2019, pp. 8026–8037.
- [46] A. Krizhevsky, G. Hinton et al., "Learning multiple layers of features from tiny images," 2009.

- [47] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," 2011.
- [48] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. A. Raffel, "Mixmatch: A holistic approach to semi-supervised learning," in *Advances in Neural Information Processing Systems*, vol. 32. Curran Associates, Inc., 2019.
- [49] H. Wang, M. Yurochkin, Y. Sun, D. Papailiopoulos, and Y. Khazaeni, "Federated learning with matched averaging," in *International Conference on Learning Representations*, 2020.
- [50] V. Kuleshov and D. Precup, "Algorithms for multi-armed bandit problems," *arXiv preprint arXiv:1402.6028*, 2014.
- [51] B. Shahriari, K. Swersky, Z. Wang, R. P. Adams, and N. De Freitas, "Taking the human out of the loop: A review of bayesian optimization," *Proceedings of the IEEE*, vol. 104, no. 1, pp. 148–175, 2015.
- [52] X. Zhu and A. B. Goldberg, "Introduction to semi-supervised learning," *Synthesis lectures on artificial intelligence and machine learning*, vol. 3, no. 1, pp. 1–130, 2009.
- [53] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," in *International Conference on Learning Representations*, 2018.
- [54] Z. Tang, S. Shi, and X. Chu, "Communication-efficient decentralized learning with sparsification and adaptive peer selection," *arXiv preprint arXiv:2002.09692*, 2020.
- [55] T. Dettmers, "8-bit approximations for parallelism in deep learning," in *International Conference on Learning Representations*, 2016.