

Communication-Efficient Federated Distillation with Active Data Sampling

Lumin Liu, Jun Zhang, S. H. Song, and Khaled B. Letaief, *Fellow, IEEE*

Dept. of ECE, The Hong Kong University of Science and Technology, Hong Kong
Email: lliubb@ust.hk, eejzhang@ust.hk, eeshsong@ust.hk, eekhaled@ust.hk

Abstract—Federated learning (FL) is a promising paradigm to enable privacy-preserving deep learning from distributed data. Most previous works are based on federated average (FedAvg), which, however, faces several critical issues, including a high communication overhead and the difficulty in dealing with heterogeneous model architectures. Federated Distillation (FD) is a recently proposed alternative to enable communication-efficient and robust FL, which achieves orders of magnitude reduction of the communication overhead compared with FedAvg and is flexible to handle heterogeneous models at the clients. However, so far there is no unified algorithmic framework or theoretical analysis for FD-based methods. In this paper, we first present a generic meta-algorithm for FD and investigate the influence of key parameters through empirical experiments. Then, we verify the empirical observations theoretically. Based on the empirical results and theory, we propose a communication-efficient FD algorithm with active data sampling to improve the model performance and reduce the communication overhead. Empirical simulations on benchmark datasets will demonstrate that our proposed algorithm effectively and significantly reduces the communication overhead while achieving a satisfactory performance.

I. INTRODUCTION

Federated Learning (FL) has recently attracted considerable attention due to its ability to collaboratively and effectively train machine learning models while preserving users' privacy [1]. A popular FL algorithm is Federated Average (FedAvg) [2], which aggregates models trained by different clients via weight averaging. FedAvg has been successfully implemented on real-world applications [3] and has inspired tremendous research interests in designing efficient and robust FL algorithms [4].

Nevertheless, weight-averaging-based methods have many limitations. For example, the local neural network architectures at different clients have to be the same, and the communication overhead is proportional to the local model size. The communication issue has been partially addressed by adopting model compression techniques to reduce the communication cost [5], while the restrictions of model architectures have been largely ignored. In a realistic FL system, clients have heterogeneous computational and communication resources. Hence, it would be highly ineffective to require all the local models to be of the same architecture.

To allow heterogeneous models and reduce the communication overhead, knowledge distillation (KD) was introduced to enable effective low-cost information exchange in FL. KD [6]

is an effective mechanism to transfer knowledge from a large teacher model to a small student model, where the student model mimics the teacher model's output, i.e., logits, on the same training data. Thus, the model architecture of the student can be different and the communication cost only depends on the logits size rather than the model weights. However, since KD is data-dependent, the training data were assumed to be universally accessible in classic KD methods. Considering the privacy regulation in FL, Federated Distillation (FD) needs to achieve distillation without sharing the local private data.

In [7], distillation was achieved by transmitting and aggregating label-wise logits of the local training data. In [8], an auxiliary distillation dataset was generated with a linear mixture of the local training data. However, the learning performance of these two approaches degrades noticeably compared with FedAvg. In [9], it was assumed that there exists a public unlabeled dataset at both the server and the clients for the distillation process. An entropy reduction technique was proposed to improve the model performance under non-iid data. In [10], delta-coding on the logits was proposed to further reduce the communication cost and the knowledge was distilled at the server side. In [11], distillation was introduced as an additional technique after weight averaging at the server side. In [12], fully distributed distillation in a connected network was considered and the gradient of the training loss function was proved to converge to zero asymptotically. These approaches showed comparable or even better performance than the weight-averaging method with a much less communication cost and even in heterogeneous model architectures.

Existing FD algorithms, while sharing similar key steps, are proposed from different perspectives, which makes it difficult to characterize and improve their performance. For FedAvg, systematic and theoretical understandings have been developed [13], which enables further design and optimization for the FL system with weight-averaging-based methods. However, for these FD algorithms, despite the empirical success, there lacks a clear understanding, either experimentally or theoretically, of the key components, i.e., 1) the auxiliary data distribution; 2) the logits aggregation strategy; and 3) the size of the upload logits.

In this paper, we endeavor to fill this important gap and answer these questions. We will first propose a generic meta-algorithm for FD, and investigate the effects of key parameters. Our results will show that in order to achieve a good

training performance, the public auxiliary data distribution should be close to the local training data, the logits aggregation strategy should reduce the logits entropy, and the size of the upload logits size should be sufficiently large. To verify and better understand these observations, we will provide a theoretical characterization of the FD meta-algorithm with a binary classification problem and Gaussian mixture models.

Inspired by the findings from these empirical observations and theoretical results, i.e., the logit entropy should be low and the distillation set size should be large, we will propose a communication-efficient FD algorithm, named, Federated distillation with Active data Sampling (FAS). In the proposed algorithm, each client only uploads a subset of the logits with low entropy. Accordingly, the samples from the public data that join the distillation will be different among different users and thus the size of the distillation logits at the server size will be larger than the upload communication cost. Simulation results will demonstrate that the proposed algorithm achieves a better performance under a limited communication cost and non-iid data distribution compared with baseline FD methods.

II. PRELIMINARY

In this section, we briefly introduce the notations for FL and KD, respectively.

A. Federated Learning

In FL, there are n clients with local private datasets $\{\mathcal{D}_i\}_{i=1}^n$ following the probability distribution $\{\mathcal{P}_i\}_{i=1}^n$. The dataset size of the i -th client is D_i . Based on the local dataset $\{\mathcal{D}_i\}$, the empirical local loss function for the i -th client is expressed as

$$L_i(\theta) = \frac{1}{D_i} \sum_{\{\mathbf{x}_j, y_j\} \in \mathcal{D}_i} \mathcal{L}(\theta, \mathbf{x}_j, y_j), \quad (1)$$

where $\mathcal{L}(\theta, \mathbf{x}_j, y_j)$ is the loss function of the training data sample \mathbf{x}_j and its label y_j , and θ denotes the model parameters. The target in FL is to learn a global model that performs well on the average of the local data distributions. Denote the joint dataset as $\mathcal{D} = \bigcup_{i=1}^n \mathcal{D}_i$ then the target training loss function in FL is given by

$$L(\theta) = \frac{1}{\sum_{i=1}^n D_i} \sum_{\xi_j \in \mathcal{D}} \mathcal{L}(\theta, \xi_j) = \frac{1}{\sum_{j=1}^n D_j} \sum_{i=1}^n D_i L_i(\theta). \quad (2)$$

The most commonly adopted training algorithm in FL is FedAvg, where each client periodically updates its model locally and averages the local model parameters through communications with a central server (e.g., at the cloud or edge). The parameters of the local model on the i -th client after t steps of stochastic gradient descent (SGD) iterations are denoted as θ_t^i . In this case, θ_t^i evolves as follows

$$\theta_t^i = \begin{cases} \theta_{t-1}^i - \eta \tilde{\nabla} L_i(\theta_{t-1}^i) & t \mid \tau \neq 0 \\ \frac{1}{n} \sum_{i=1}^n [\theta_{t-1}^i - \eta \tilde{\nabla} L_i(\theta_{t-1}^i)] & t \mid \tau = 0 \end{cases} \quad (3)$$

B. Knowledge Distillation

Knowledge Distillation (KD) is the process of distilling knowledge from a large and well-trained teacher model to

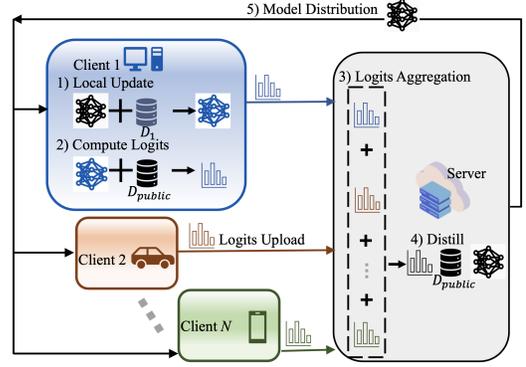


Figure 1: Illustration of FD meta-algorithm.

a small student model. For a classification problem with N_c classes, the logit of data sample \mathbf{x}_i is denoted as $t(\mathbf{x}_i)$ and it is the vector of the class probabilities which is obtained by using a softmax function on the neural network output. That is,

$$t(\mathbf{x}_i) = \text{softmax}(\theta(\mathbf{x}_i)), \quad (4)$$

where $\theta(\mathbf{x}_i) \in \mathcal{R}^{N_c}$ denotes the model output of input data sample \mathbf{x}_i , and $\theta(\cdot)$ is the neural network function parameterized by model parameters θ . Specifically, for the n -th element of logit $t(\mathbf{x}_i)$ of data sample \mathbf{x}_i ,

$$t^n(\mathbf{x}_i) = \frac{\exp(\theta(\mathbf{x}_i)^n/T)}{\sum_{m=1}^{N_c} \exp(\theta(\mathbf{x}_i)^m/T)}, \quad (5)$$

where T is the distillation temperature with a higher temperature producing a smoother probability distribution over classes.

The distillation loss of the trainset \mathcal{D} is the cross-entropy loss for the teacher logit t_t and the student logit t_s , which is

$$L_{\text{distill}} = - \sum_{\mathbf{x} \in \mathcal{D}} \sum_{n=1}^{N_c} t_t^n(\mathbf{x}) \log(t_s^n(\mathbf{x})). \quad (6)$$

In the distillation process, the student's objective function is an average of the distillation loss L_{distill} and the cross entropy loss with the groundtruth labels.

III. FEDERATED DISTILLATION META-ALGORITHM

In this section, we will first introduce the FD system and present a meta-algorithm, which is constituted of several key components. Then, we will investigate the impacts of these key components both empirically and theoretically.

A. FD Meta-Algorithm

For a FD system with n clients, the local private labeled dataset of the i -th client is denoted as $\mathcal{D}_i = \{\mathbf{x}_i^j, y_i^j\}_{j=1}^{D_i}$. A shared public unlabeled dataset $\mathcal{D}_{\text{pub}} = \{\mathbf{x}^j\}_{j=1}^{D_{\text{pub}}}$ is assumed accessible for each client and the server, where each data sample is identified by a unique and universal index. The local loss function of client i with local model parameters θ_i is denoted as $L_i(\theta_i)$. In the k -th communication round, the selected clients perform local updates on their local private datasets \mathcal{D}_i 's and get locally trained models θ_k^i 's.

Table I: Comparison of different algorithms.

	Upload	Aggregation	Auxiliary Dataset	Model Heterogeneity	Communication Cost (Uplink)	Model Performance
FedAvg [2]	Weights	Average	×	×	$\mathcal{O}(\theta)$	Baseline
FedDF [11]	Weights	Average & Distill	✓	✓	$\mathcal{O}(\theta)$	✓
FDA [7]	Label-logits	Average	×	✓	N_c^2	×
DSFL [9]	Logits	Entropy Reduction Average	✓	✓	$ \mathcal{D}_{logit} N_c$	✓
CEFD [10]	Delta-coded logits	Average	✓	✓	\setminus^1	✓
FD meta-algorithm	Logits	Average&Distill	✓	✓	$ \mathcal{D}_{logit} N_c$	✓

Algorithm 1: FD Meta-Algorithm

Initialize local model $\{\theta^i\}$ and server model θ
for $k = 0, 1, \dots, K-1$ **do**
 Download the server model $\theta_{k-1}^i = \theta_{k-1}$,
 Select clients \mathcal{C} from the n clients,
 Select a subset \mathcal{D}_{logit} of the public dataset \mathcal{D}_{pub} ,
 for *client* $i \in \mathcal{C}$ **do**
 Local update: $\theta_k^i = \theta_{k-1}^i - \eta \tilde{\nabla} L_i(\theta_{k-1}^i)$,
 Compute the logits:
 $t_i(\xi) = \text{softmax}(\theta_k^i(\xi))$ for $\xi \in \mathcal{D}_{logit}$
 Upload the logits and indexes:
 $\{t_i(\mathbf{x})\}_{\mathbf{x} \in \mathcal{D}_{logit}}, \mathcal{I}_{logit}$
 end
 Aggregate the logits: $t(\mathbf{x}) = \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} t_i(\mathbf{x})$
 Model distillation: $\theta_t = \theta_{k-1} - \eta \tilde{\nabla} L_{distill}(\theta_{k-1})$
end

The weight-averaging-based method will directly upload and average the model weights of different clients, and then the training proceeds to the next communication round. However, since the local models $\{\theta_i\}_{i=1}^n$ may have different neural network architectures, e.g., simple fully-connected neural networks and ResNets, it is infeasible to directly average the model weights of these heterogeneous clients.

To enable information sharing of the clients with heterogeneous neural architectures, in FD, the selected clients will compute the logits on a subset \mathcal{D}_{logit} of the public unlabeled dataset \mathcal{D}_{pub} , and the indexes of the data sample in \mathcal{D}_{logit} are denoted as \mathcal{I}_{logit} . The computed logits of the selected subset $\{t_i(\mathbf{x})\}_{\mathbf{x} \in \mathcal{D}_{logit}}$ and the index \mathcal{I}_{logit} are uploaded to the server for logits averaging. The averaged logits then serve as the teacher logits in the distillation loss in (6). The distilled model is then distributed back to the selected clients in the next communication round. The uploading communication cost is $N_c * |\mathcal{D}_{logit}|$ and the downloading communication cost is proportional to the local model size, i.e., $\mathcal{O}(|\theta_i|)$. The FD system and the detailed procedure of the algorithm are illustrated in Fig. 1 and Algorithm 1, respectively.

It is worth noting that in some existing works (e.g., [9]), the averaged logits are distributed to the clients and the distillation happens at the client side. Local distillation reduces

the downloading communication cost to $N_c * |\mathcal{D}_{logit}|$ and is completely free of the worry of the model heterogeneity. However, it also induces more local computation. In addition, partial client participation is not allowed if the averaged logits are sent back to clients. To allow a heterogeneous model for the local update, the server can distill the averaged weights into different models and then send back the weights to its corresponding client. We compared these two methods empirically and found that the weights downloading method exhibits a faster convergence. Thus, we will adopt the model weights downloading method. Finally, the differences of the typical algorithms mentioned in this paper are summarized in Table I. Given the enormous size of deep learning models, distillation-based methods achieve orders of magnitude reduction in the communication overhead and allow heterogeneous models for the local update. Comparable model performance can be achieved with an auxiliary public dataset. The FD-meta algorithm concluded the key components for the FD-based methods and can be extended to the existing work [9], [10] with slight modifications, e.g. in [9] the entropy of the averaged logits was reduced. With this FD meta-algorithm, we can better understand the design principles in a FD system.

B. Empirical Observations

There are some key components in the FD system which influence the communication cost and the final learning performance, i.e., the data distribution of public dataset \mathcal{D}_{pub} , the logits aggregation method, and the upload logits size. In this section, we investigate these key components with the FD-meta algorithm. With extensive simulations on the CIFAR-10 dataset, we will show their impacts in the following.

1) *Distillation dataset distribution:* A vital assumption in FD is the availability of a public unlabeled dataset which enables the distillation process. In practice, it is not difficult to collect or generate many unlabeled samples. However, it is difficult to collect or generate a public dataset which has the same data distribution as the private labeled dataset. In the empirical simulations of the FD works, the public dataset distribution problem is often ignored.

To investigate the impact of the distillation dataset distribution, we performed experiments with two distillation datasets, i.e., CIFAR-10, the dataset with exactly the same distribution, and STL-10, the dataset with a similar but broader distribution. The result is demonstrated in Fig. 2a, which shows the test accuracy of the server model after the clients upload their logits or weights to the server, which is one communication

¹Since in CEFD, delta coding is applied to the logits of the whole distillation datasets \mathcal{D}_{pub} , the communication cost is smaller than $|\mathcal{D}_{pub}|N_c$. But it varies in the training process.

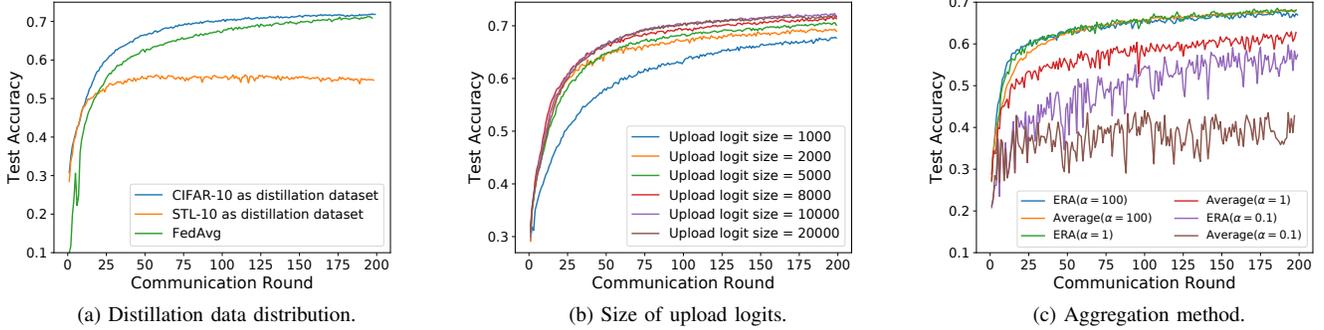


Figure 2: Empirical observations of the algorithm key components' impacts on the FD training performance. The figure lists the test accuracy versus the number of communication rounds between the clients and server.

round. It can be clearly seen that there exists a noticeable performance gap between the public dataset with similar distribution and the same distribution. And when distilling with CIFAR-10 dataset, the model reaches a comparable test accuracy with FedAvg.

2) *Upload Logits Size*: In the FD meta-algorithm, the uplink communication cost is determined by the number of selected data samples, i.e., the upload logits size. A straightforward way to further reduce the communication cost is to reduce the size of the selected public dataset subset, \mathcal{D}_{logit} . However, this will cause insufficient data for the distillation step at the server side. Hence, there exists a trade-off between the communication cost and accuracy.

To empirically investigate this trade-off, we perform experiments where the size of the uploaded logits ranges from 1,000 to 20,000. The empirical results are demonstrated in Fig. 2b. It is seen that increasing the upload logits size from a relatively small number improves the training performance. However, as the logits size increases to a very large number, the performance gain of more distillation data samples becomes marginal. For example, by uploading 20,000 logits, we barely see any performance gain compared with the one with 10,000.

3) *Logits Aggregation Method*: In the meta-algorithm, a simple average is adopted for the logits aggregation at the server side. However, the simple average method shows a bad performance when the local private data distribution is non-i.i.d.. Entropy reduction aggregation (ERA) is an aggregation method, which was proposed in [9]. There it was shown that it can achieve much better performance compared with the simple average method.

The main idea in ERA is to increase the confidence of the aggregated teacher logits during the server distillation step. ERA first averages the logits uploaded by the selected clients

$$t(\mathbf{x}) = \frac{1}{|\mathcal{C}|} \sum_{i \in \mathcal{C}} t_i(\mathbf{x}) \quad \text{for } \mathbf{x} \in \mathcal{D}_{logit}, \quad (7)$$

The entropy of the averaged logit $t(\mathbf{x})$ is then reduced by:

$$\hat{t}(\mathbf{x}) = \frac{\exp(t(\mathbf{x})/T)}{\sum_{m=1}^{N_c} \exp(t(\mathbf{x})^m/T)} \quad (8)$$

where T here should be set between 0 and 1 so as to sharpen the output and reduce the entropy of \hat{t} .

We adopt the Dirichlet distribution $Dir(\alpha)$ to simulate the non-i.i.d. data distribution in FL and perform experiments with three levels of non-i.i.d. data distribution, i.e., $\alpha = 100, 1$, and 0.1 . It is noted that the data heterogeneity increases as α decreases. The result is demonstrated in Fig. 2c. When $\alpha = 0.1$, i.e., the local data distribution is very non-i.i.d., reducing the entropy of the logits greatly improves the FD training performance.

C. Theoretical Verification

From the empirical observations, we have seen that for the FD meta-algorithm,

- 1) A public unlabeled dataset with the same input distribution is necessary to guarantee a good training performance;
- 2) The size of upload logits influences the convergence speed. More logits lead to a better performance, but the performance gain becomes marginal when there is a sufficient amount of uploaded logits;
- 3) ERA improves the model performance of non-i.i.d. data distribution.

In this subsection, we verify the latter two observations theoretically through a binary classification problem with Gaussian mixture models. Particularly, we show that for this setting, the FD meta-algorithm is equivalent to semi-supervised learning (SSL) with self-training [14].

We first give a definition of the binary classification problem, the Gaussian mixture models, and self-training. For the binary classification problem, suppose there is a labeled dataset $\mathcal{S} = (\mathbf{x}_i, y_i) \in \mathbb{R}^p \times \{-1, +1\}$ and $f: \mathbb{R}^p \rightarrow \mathbb{R}$ is the prediction function. The prediction rule is then defined as:

$$\hat{y}_f(\mathbf{x}) = \begin{cases} 1 & \text{if } f(\mathbf{x}) \geq 0 \\ -1 & \text{otherwise} \end{cases} \quad (9)$$

Definition 1 (Binary Gaussian Mixture Model (GMM)) The distribution $(\mathbf{x}, y) \sim \mathcal{D}$ is given as follows. Fix a unit vector $\mathbf{u} \in \mathbb{R}^p$ and a scalar $\sigma > 0$, and let y be a Rademacher random variable ($\mathbb{P}(y = 1) = 1 - \mathbb{P}(y = -1) = \frac{1}{2}$ and $\mathbf{x} \sim \mathcal{N}(y\mathbf{u}, \sigma \mathbf{I}_p)$).

The component mean \mathbf{u} is the optimal linear classifier for this binary classification problem, where the prediction func-

tion is $f(\mathbf{x}) = \mathbf{u}^T \mathbf{x}$. With a labeled dataset $\mathcal{S} = (\mathbf{x}_i, y_i)_{i=1}^n$, \mathbf{u} can be estimated by the following averaging estimator

$$\beta_{init} = \frac{1}{n} \sum_{i=1}^n y_i \mathbf{x}_i, \quad (10)$$

The self-training approach uses the predicted labels $\hat{y}_f(\mathbf{x})$ for an unlabeled dataset $\mathcal{U} = \{\mathbf{x}_i\}_{i=n+1}^{n+u}$ (a.k.a, pseudo labels) to self-train. Given the initial averaging estimator β_{init} of the labeled dataset in (10) and an acceptance threshold $\beta_{init}^T \mathbf{x} > \Gamma$, the updated estimator after self-training with the unlabeled dataset \mathcal{U} is then

$$\hat{\beta} = \frac{\sum_{i=n+1}^u \mathbb{1}(|\beta_{init}^T \mathbf{x}_i| > \Gamma) \text{sgn}(\beta_{init}^T \mathbf{x}_i) \mathbf{x}_i}{\sum_{i=n+1}^u \mathbb{1}(|\beta_{init}^T \mathbf{x}_i| > \Gamma)}. \quad (11)$$

In the following, we will show the training process of the FD meta-algorithm with the binary classification problem of GMM as the learning objective. In FD, there are N locally stored private datasets, $\mathcal{S}^i = (\mathbf{x}_i^k, y_i^k)_{i=1}^{n_k}$, and the unlabeled auxiliary distillation dataset is denoted as $\mathcal{U} = \{\mathbf{x}_i\}_{i=n}^{n+u}$. Following the steps in the FD meta-algorithm (Algorithm 1), the training proceeds as follows

- 1) **Local Update:** After the local updates, each user k will have a local averaging estimator as

$$\beta_{init}^k = \frac{1}{n_k} \sum_{i=1}^{n_k} y_i \mathbf{x}_i^k \quad (12)$$

- 2) **Logits Computation:** Each user will compute the local model output (logits) of the unlabeled dataset, i.e., $\{(\beta_{init}^k)^T \mathbf{x}_i\}_{i=n}^{n+u}$ and upload the logits to the server.
- 3) **Logits Aggregation:** The server averages the logits and we have the averaged logits of the distillation dataset $\mathcal{U} = \{\mathbf{x}_i\}_{i=n}^{n+u}$ as

$$\sum_{k=1}^n \frac{n_k}{n} \{(\beta_{init}^k)^T \mathbf{x}_i\} = \beta_s^T \mathbf{x}_i \quad (13)$$

- 4) **Model Distillation:** The server creates pseudo labels by choosing data samples in \mathcal{U} whose logits satisfy $|\beta_s^T \mathbf{x}_i| > \Gamma$ and the pseudo labels are generated by $\tilde{y} = \hat{y}_{\beta_s^T \mathbf{x}}(\mathbf{x})$. After distillation with the averaged logits, the estimator at the server side with the averaged logits is then

$$\hat{\beta} = \frac{\sum_{i=n+1}^u \mathbb{1}(|\beta_s^T \mathbf{x}_i| > \Gamma) \text{sgn}(\beta_s^T \mathbf{x}_i) \mathbf{x}_i}{\sum_{i=n+1}^u \mathbb{1}(|\beta_s^T \mathbf{x}_i| > \Gamma)} \quad (14)$$

where $\Gamma \geq 0$ is the acceptance threshold that eliminates low-confidence predictions. It is noted that this is similar to the ERA method, which also eliminates high-entropy, i.e., low confidence predictions in the distillation process.

We measure the estimator performance with the cotangent of the angle of the estimator β and the optimal classifier \mathbf{u} :

$$\cot(\beta, \mathbf{u}) = \frac{\rho(\beta, \mathbf{u})}{\sqrt{1 - \rho^2(\beta, \mathbf{u})}}. \quad (15)$$

With $\hat{\beta}$, we have the following theorem.

Theorem 1. ([14]) Let $\mathbf{u} \in \mathbb{R}^p$ be a uniform vector from Definition 1 and suppose $\beta_s \in \mathbb{R}^p$ as defined in (13) has correlation $\rho(\beta_s, \mathbf{u}) = \alpha > 0$. Set $\beta = \sqrt{1 - \alpha^2}$ and draw

i.i.d. unlabeled samples $\{\mathbf{x}_i\}_{i=n+1}^{n+u}$ from GMM. Let $\hat{\beta}$ be defined in (14). Define the normalized thresholds $\bar{\Gamma}_- = \frac{\alpha + \Gamma}{\sigma}$ and $\bar{\Gamma}_+ = \frac{\Gamma - \alpha}{\sigma}$ and the quantities

$$\Lambda = \frac{1}{2\pi\rho} (\exp(-\bar{\Gamma}_+^2/2) + \exp(-\bar{\Gamma}_-^2/2))$$

$$\rho = Q(\bar{\Gamma}_+) + Q(\bar{\Gamma}_-) \quad (16)$$

$$\nu = Q(\bar{\Gamma}_-)/\rho$$

where $Q(\cdot)$ is the tail of standard normal variable. Then, by fixing $\bar{u} = u/p$ and letting $p \rightarrow \infty$, we have

$$\cot(\hat{\beta}, \boldsymbol{\mu}) \xrightarrow{\mathbb{P}} \frac{1 + \sigma\alpha\Lambda - 2\nu}{\sigma\sqrt{(1 - \alpha^2)\Lambda^2 + 1/\bar{u}\rho}}. \quad (17)$$

Proof. From [14], it is proved that for the self-training algorithm with the initial estimator $\beta_{init} = \frac{1}{n} \sum_{i=1}^n y_i \mathbf{x}_i$ in (10) and $\hat{\beta}$ in (11) Theorem 1 holds.

From (12) and (13),

$$\beta_s = \sum_{k=1}^n \frac{n_k}{n} \beta_{init}^k = \frac{1}{n} \sum_{i=1}^n y_i \mathbf{x}_i. \quad (18)$$

The estimator by distributed training of the labeled data samples is the same as the self-training algorithms. Thus, the result still holds for the FD meta-algorithm. \square

Remark 1 Theorem 1 shows that for the GMM binary classification problem, the FD algorithm can obtain a higher correlation for the estimator than the initial estimator β_{init} , i.e., a better model is obtained after the model distillation step. The distilled model $\hat{\beta}$ benefits from a larger unlabeled dataset and a higher accepting threshold, which is consistent with the empirical observations in Section III-B.

IV. PROPOSED ALGORITHM VIA ADAPTIVE DATA SAMPLING

The theoretical and empirical results suggest two approaches to improve the training performance: 1) increase the size of the logits; or 2) choose the public data with low-entropy logits. Thus, we propose a communication-efficient FD algorithm with Active data Sampling (FAS). To increase the size of the distillation logits while maintaining the communication cost, each user will actively select the low-entropy logits to be uploaded with its locally trained model.

The main difference between the FD meta-algorithm and the proposed FAS algorithm is the active data sampling step. To generate better teacher logits, the entropy of the selected logits should be low, which means that the local model is confident. However, using the low entropy as the only criterion may lead to the result that every client is very confident about its uploaded logits, but the selected data samples for distillation are very easy to classify, which may degrade the final performance. This is similar to the process of human learning. If we always learn tasks that we are already very confident about, then we cannot learn new things. We need to learn something basic but we also need to explore new and challenging things. Thus, we propose the following mixed active data selection strategy. For a selected client i ,

assuming the communication budget is N_{logit} logit samples, then the active data sampling step proceeds as follows:

- 1) Generate pseudo labels of D_{pub} with the locally trained model θ_i ;
- 2) Select $N_{logit}/2$ logits from D_{pub} with an ascending order in entropy to generate half of D_{logit}^i , and the pseudo label distribution in this half D_{logit}^i needs be close to the local data distribution;
- 3) Randomly select $N_{logit}/2$ from D_{pub} to generate the other half of D_{logit}^i .

We next provide experimental results to demonstrate the effectiveness of the proposed FAS algorithm. In the simulated FD system, there are 20 clients. The local private training data are a subset with 20,000 data samples of the CIFAR-10 dataset, which means each user has only 1,000 local private data samples. The distillation dataset is the other 20,000 data samples of the CIFAR-10 dataset. The neural network model is ResNet-8. In each communication round, 8 clients are selected randomly to participate in the learning process. For the local update and distillation, we adopt Adam with a batch size of 8 as the optimizer. The local update steps and the distillation steps are set as 20 epochs in each communication round. The step sizes of the local update and distillation are set to 0.02 and 0.001, respectively. The step size decays at the 300-th and 600-th epochs by a rate of 0.1.

We compare the following 4 data sampling methods for FAS under different data distributions:

- 1) No data sampling (D_{logit} is the same);
- 2) Random data sampling (D_{logit}^i is randomly sampled from D_{pub});
- 3) Low-entropy data sampling (D_{logit}^i is sampled from D_{pub} assuming an entropy ascending order);
- 4) Mixed-random-low-entropy sampling.

The results are listed in Table II and Table III. To ensure that there is an overlap of the selected logits, we select 400 logits from 2000 unlabeled public data samples for the simulations in Table II and 2000 logits from 8000 unlabeled public data for the simulations in Table III in each communication round. We perform experiments with different degrees of non-i.i.d. data distribution controlled by α in the Dirichlet distribution. A smaller α leads to a more non-i.i.d. data distribution.

All the methods with data sampling exhibit better performance than the No data sampling method due to a larger distillation dataset. The performance of the random sampling method degrades evidently when the data becomes non-i.i.d., i.e., α decreases. And the performance of low-entropy sampling increases with more non-i.i.d. data. The mixed sampling strategy provides consistently better or comparable performance in terms of test accuracy compared with the other sampling methods, under different degrees of non-i.i.d. local data distributions.

V. CONCLUSIONS

In this paper, we presented an FD meta-algorithm that incorporates existing FD methods and investigated the effects of key parameters both experimentally and theoretically to

Table II: Test accuracy of ResNet-8 on CIFAR-10, $N_{logit} = 500$.

	NoSample	Random	Low-Entropy	Mixed
$\alpha = 100$	0.6835	0.6956	0.6694	0.7058
$\alpha = 1$	0.6376	0.6516	0.6468	0.6506
$\alpha = 0.1$	0.4644	0.468	0.5219	0.5498

Table III: Test accuracy of ResNet-8 on CIFAR-10, $N_{logit} = 2000$.

	NoSample	Random	Low-Entropy	Mixed
$\alpha = 100$	0.7314	0.7381	0.7031	0.7400
$\alpha = 1$	0.6947	0.6827	0.7014	0.6947
$\alpha = 0.1$	0.5511	0.5308	0.5573	0.5788

provide several design guidelines. Inspired by the design guidelines, a simple but effective FD algorithm with active data sampling was proposed. Experiments showed that the proposed algorithm performs consistently well under different distributions of heterogeneous data. Analyzing the FD meta-algorithm for neural networks and adapting it to the auxiliary dataset with similar distribution are left for future work.

REFERENCES

- [1] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.
- [2] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. Y. Arcas, "Communication-efficient learning of deep networks from decentralized data," *Artificial Intelligence and Statistics*, pp. 1273–1282, 2017.
- [3] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, "Federated learning for mobile keyboard prediction," *arXiv preprint arXiv:1811.03604*, 2018.
- [4] L. Liu, J. Zhang, S. Song, and K. B. Letaief, "Client-edge-cloud hierarchical federated learning," *IEEE International Conference on Communications (ICC)*, pp. 1–6, 2020.
- [5] A. Reiszadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani, "Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization," *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, vol. 108, pp. 2021–2031, 2020.
- [6] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [7] E. Jeong, S. Oh, H. Kim, J. Park, M. Bennis, and S.-L. Kim, "Communication-efficient on-device machine learning: Federated distillation and augmentation under non-iid private data," *arXiv preprint arXiv:1811.11479*, 2018.
- [8] S. Oh, J. Park, E. Jeong, H. Kim, M. Bennis, and S.-L. Kim, "Mix2fld: Downlink federated learning after uplink federated distillation with two-way mixup," *IEEE Communications Letters*, vol. 24, no. 10, pp. 2211–2215, 2020.
- [9] 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*, 2021.
- [10] F. Sattler, A. Marban, R. Rischke, and W. Samek, "Communication-efficient federated distillation," *arXiv preprint arXiv:2012.00632*, 2020.
- [11] T. Lin, L. Kong, S. U. Stich, and M. Jaggi, "Ensemble distillation for robust model fusion in federated learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 2351–2363, 2020.
- [12] I. Bistriz, A. Mann, and N. Bambos, "Distributed distillation for on-device learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 22 593–22 604, 2020.
- [13] S. Wan, J. Lu, P. Fan, Y. Shao, C. Peng, and K. B. Letaief, "Convergence analysis and system design for federated learning over wireless networks," *IEEE J. Select. Areas Commun. Early Access*, 2021.
- [14] S. Oymak and T. Cihad Gulcu, "A theoretical characterization of semi-supervised learning with self-training for gaussian mixture models," *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, vol. 130, pp. 3601–3609, 13–15 Apr 2021.