

Self-supervised On-device Federated Learning from Unlabeled Streams

Jiahe Shi, Yawen Wu, Dewen Zeng, Jun Tao, *senior member, IEEE*, Jingtong Hu, *senior member, IEEE*,
Yiyu Shi, *senior member, IEEE*,

Abstract—The ubiquity of edge devices has led to a growing amount of unlabeled data produced at the edge. Deep learning models deployed on edge devices are required to learn from these unlabeled data to continuously improve accuracy. Self-supervised representation learning has achieved promising performances using centralized unlabeled data. However, the increasing awareness of privacy protection limits centralizing the distributed unlabeled image data on edge devices. While federated learning has been widely adopted to enable distributed machine learning with privacy preservation, without a data selection method to efficiently select streaming data, the traditional federated learning framework fails to handle these huge amounts of decentralized unlabeled data with limited storage resources on edge. To address these challenges, we propose a Self-supervised On-device Federated learning framework with coreset selection, which we call SOFed, to automatically select a *coreset* that consists of the most representative samples into the replay buffer on each device. It preserves data privacy as each client does not share raw data while learning good visual representations. Experiments demonstrate the effectiveness and significance of the proposed method in visual representation learning.

Index Terms—Learning Representations, On-Device Learning, Federated Contrastive Learning, Self-Supervised Learning

I. INTRODUCTION

Deep learning models deployed on edge devices have made significant improvements in the last few years to extract useful visual representations in real-world tasks, such as sensors for fault diagnosis in smart manufacturing [1], and agricultural robots for fruit anomalies monitoring [2]. The significance of the deep neural networks highly relies on the centralized supervised training paradigm, where the models are trained on labeled data in the server. However, a tremendous amount of unlabeled image data that cannot be centralized is produced by edge devices every day. On the one hand, the decentralized data is usually of private nature, e.g., photos taken from medical monitoring may contain patient faces. Sending the data to the server may violate data privacy regulations. On the other hand, the cost of labeling is prohibitive. Hence, it is desirable for the deep learning models deployed on edge devices to directly learn from this unlabeled decentralized data stream. Besides, compared with the traditional training method, deep

learning models learn better representations for the deployed environment by leveraging local data.

To learn effective visual representations without supervision, researchers have proposed many self-supervised representation learning methods [3]–[9]. Among them, contrastive learning [5]–[8] is the state-of-the-art method. Contrastive learning methods instruct neural networks to minimize dissimilarity between representations of two similar data points or to maximize dissimilarity between representations of two different data points. Although contrastive learning has demonstrated promising performance in learning from unlabeled data, it is conducted on centralized data in servers or distributed devices without storage restrictions.

However, it is impractical to collect a large centralized dataset from edge devices due to privacy concerns and to store all the data locally on each device due to storage overhead. First, due to privacy regulations, it is unfeasible to directly send the data samples in the local data buffer to a high-performance server. Existing methods learning from unlabeled data [5]–[8] cannot collaboratively leverage decentralized unlabeled image data to learn a visual representation under the restriction of limited storage resources while keeping data private. Second, the images generated by the edge devices are usually captured by the integrated cameras in sequence and normally are non-independent-and-identically distributed (non-iid). Storing the constantly generated streaming data on distributed edge devices is prohibitive because of the storage and energy overhead. Therefore, it is necessary to select only a “coreset” into the small data buffer on edge device where only the most representative samples are stored.

Federated learning [10] is a distributed training technique with privacy guaranteed. Existing works that combine contrastive learning with federated learning [11]–[14] either have privacy leakage risks by sharing features or do not consider the restriction of limited storage volumes. To learn visual representations from unlabeled data generated by networked devices with limited on-device storage, a small data buffer can be used for each edge device to maintain the most effective data for learning. However, there are several challenges we need to address to select high-quality data into the data buffer in the federated contrastive learning (FCL) scenario. First, privacy cannot be compromised when measuring the sample’s quality and impact on the global model. The features of clients’ local data should not be exposed to other parties, including the server. Otherwise, the secrets of the participating client in the network may be revealed by others. For example, in a network of distributed industrial robots in car plants, if a

Jiahe Shi and Jun Tao are with the State Key Laboratory of ASIC and System, School of Microelectronics, Fudan University, Shanghai 200433, China (e-mail: jhshi21@m.fudan.edu.cn, taojun@fudan.edu.cn).

Yawen Wu, Dewen Zeng and Yiyu Shi are with the Department of Computer Science and Engineering, University of Notre Dame, Notre Dame, IN 46556 USA (e-mail: ywu37@nd.edu, dzeng2@nd.edu, yshi4@nd.edu).

Jingtong Hu is with the Department of Electrical and Computer Engineering, University of Pittsburgh, Pittsburgh, PA 15261 USA (e-mail: jthu@pitt.edu).

robot produces a feature distinguished from others, it may reveal the manufacturing of a new vehicle model. Hence, a sample selection policy must be carefully designed to protect the privacy of clients. Second, the selected collection of samples should be effective for the global federated learning model. The distributed unlabeled data in FCL complicate the online data selection. As a result of the lack of knowledge of other devices' data distribution, the samples selected by different devices may overlap, hence, deviate from the distribution of training data collected by all devices. Moreover, due to the restricted computation, storage, and communication resources of local edge devices, the selection process should introduce marginal overheads. Without an efficient data selection policy, the performance of mobile systems may degrade.

To address these challenges, we propose a novel Self-supervised On-device Federated learning framework with coreset selection, SOFed, to learn visual representations from unlabeled data streams collaboratively on edge devices. A server coordinates distributed edge devices to conduct federated contrastive learning. We use a coreset selection method based on importance scoring to maintain the most representative samples in the local buffer. A larger importance score indicates the sample has a larger impact on the model. Selecting samples with higher importance scores into the replay buffer enables the replay buffer to maintain more valuable data. Each client conducts unsupervised representation learning on the data in the local buffer and collaboratively optimizes the global model. After federated contrastive learning, we utilize the global model as the backbone. Then we train an additional new classifier on top of the backbone model with a limited number of labels.

In summary, the main contributions of the paper include:

- **Federated on-device contrastive learning framework.** We propose an on-device FCL framework, SOFed, to form a small data buffer on each edge device for collaborative self-supervised learning from real-time streaming data. There is only a small data buffer on each edge device to eliminate unaffordable storage overhead. *This is the first work on self-supervised on-device federated learning from unlabeled streams.*
- **Coreset selection method in a federated contrastive learning scenario.** We propose a coreset selection method based on importance scoring to maintain the most effective data in the buffer for local training. Labels are not needed in the coreset selection and privacy is guaranteed by keeping features and raw data local. The computation of importance scores is efficient and induces low overhead.
- **Better classification accuracy and label efficiency.** Experiments on various public datasets show that SOFed outperforms state-of-the-art baselines with different sampling methods and self-supervised learning methods in classification accuracy and label efficiency.

II. BACKGROUND AND RELATED WORK

A. Contrastive Learning

Contrastive learning is a self-supervised approach to extracting visual representations from unlabeled data. It trains an encoder to provide generalized information for downstream

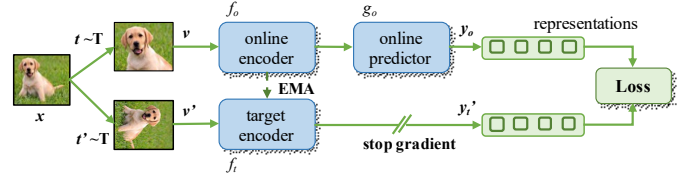


Fig. 1: Illustration of contrastive learning. It comprises an online network and a target network. An input image x is transformed by two augmentation methods $t \sim T$ and $t' \sim T$ to produce two views. These two views are then sent into the network to generate representations y_o and y_t . Contrastive learning minimizes a dissimilarity loss between y_o and y_t .

tasks by minimizing a contrastive loss. In this work, we employ the contrastive learning approach from [8], since it performs best when learning representations from decentralized unlabeled data collaboratively [15]. In the training process, the contrastive loss is evaluated on a pair of features extracted from augmented views of one image. By optimizing the contrastive loss, contrastive learning maximizes the agreement of representations between these two views. As shown in Fig. 1, contrastive learning utilizes an asymmetric network containing an online network with an online encoder f_o and a predictor g_o and a target network with a target encoder f_t . Optimizing the contrastive loss only contributes to the online network update, while stop gradient prevents the gradient optimization of the target network. The target network is updated by the exponential moving average (EMA) of the online network. For an unlabeled input image x , two data augmentations methods $t \sim T$ and $t' \sim T$ are applied to x to produce two augmented views $v = t(x)$, $v' = t'(x)$, where T is a distribution of image augmentations. They are fed into the online network and target network respectively and output representations $y_o = g_o(f_o(v))$, $y'_t = f_t(v')$. The contrastive loss is the dissimilarity between the representations:

$$\mathcal{L}(x) = \| \bar{y}_o - \bar{y}'_t \|_2^2 = 2 - 2 \cdot \frac{\langle y_o, y'_t \rangle}{\| y_o \|_2 \cdot \| y'_t \|_2}, \quad (1)$$

where $\bar{y}_o = y_o / \| y_o \|_2$ and $\bar{y}'_t = y'_t / \| y'_t \|_2$ are the ℓ_2 -normalized representations. This contrastive learning approach is compatible with federated learning since it doesn't require negative pairs formed by samples in the whole training batch that is distributed on devices.

B. Federated Learning

Federated learning is a distributed training framework to train a global model without centralizing local raw data [10]. FedAvg is a typical algorithm of federated learning methods. In each round of FedAvg, a subset of clients C are activated to perform training on local datasets $\{D_i, i \in C\}$ and transmit the updated model parameters $\{W_i, i \in C\}$ to server. Then the server aggregates the updated models and updates the global model W_g by averaging the local model parameters $W_g \leftarrow \sum_{i \in C} \frac{|D_i|}{\sum_{i \in C} |D_i|} W_i$. This communication round repeats to update the global model until convergence. While the traditional federated learning methods [10], [16]–[18] have

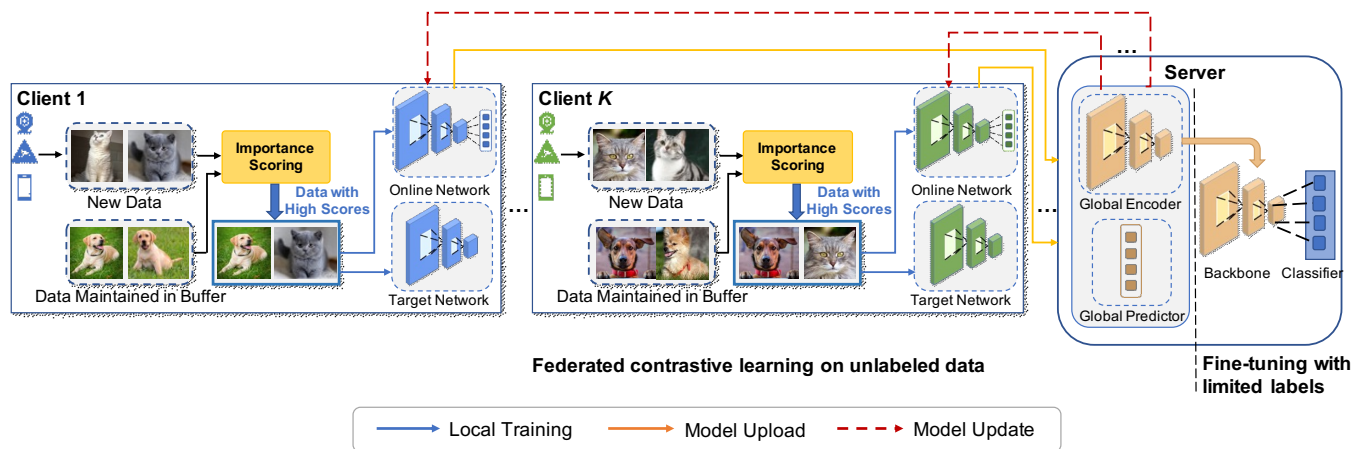


Fig. 2: Overview of the proposed federated contrastive learning framework SOFed. The network contains a server and K devices (clients). The online encoder is first collaboratively trained on clients with data selected from streaming unlabeled data to generate good representations. Then a classifier is trained with few labeled data on top of the online encoder.

privacy guarantee, they assume the devices have enough storage resources and the input data stream is fully labeled.

C. Related Work

Self-supervised Visual Representation Learning. Self-supervised learning tries to learn effective visual representations from unlabeled data [3]–[6], [8], [9]. Among those methods, contrastive learning [5]–[8] demonstrates state-of-the-art performance. Some contrastive learning methods [5], [6] rely on negative pairs to prevent model collapse. For example, SimCLR produces the negative pairs from other samples in a large training batch. MoCo generates the negative pairs from a memory bank. Recently, some contrastive learning methods are proposed to eliminate the negative pairs by adopting stop-gradient operation to prevent the collapsing of Siamese networks [7], [8]. All these methods require training on a large amount of accessible data. During training, minibatches are generated by randomly splitting the dataset and sent to the server. By reshuffling the training batches over epochs, the training data can maintain independent and identically distributed. However, in distributed on-device learning, the input streaming data distributed across devices are usually non-iid and it is not feasible to maintain a large dataset on edge devices. Simply applying contrastive learning on each device with a limited number of training data will degrade the performance. Therefore, it is necessary to conduct contrastive learning collaboratively to obtain good visual representation from online unlabeled data in a network of edge devices.

Federated Contrastive Learning. There are several federated contrastive learning methods [11]–[14] to collaboratively extract representations from unlabeled distributed data while preserving privacy. [11], [12], [14] address the non-iid challenge by sharing features, hence, having the risk of privacy leakage. [13] is based on BYOL and achieves state-of-the-art accuracy in downstream tasks. However, these methods neglect the restriction of storage resources. Therefore, existing methods cannot be applied in real-life scenarios where storage resources on edge devices are limited.

Coreset Selection in Streaming and Continual Learning. Training on streaming data with non-IID distribution results in catastrophic forgetting. To overcome catastrophic forgetting, replay buffer technique is widely adopted in continual learning. Existing methods [19]–[22] select coresets from streaming data into the replay buffer to be included in the training. To overcome the catastrophic forgetting in networked scenario, some previous works [23]–[25] propose online data selection for federated learning. [23], [24] select high-quality samples based on metrics such as gradient upper bound of each sample and the projection of the local gradient onto the global gradient. [25] trains a relevant data selector (RDS) to select data and needs the server to intervene in the training of RDS based on the distribution of labels. All the existing approaches require data labels to select coreset. However, the requirement of labeling all the streaming data on edge devices is stringent. Hence an efficient approach to select coreset into the replay buffer in FCL is needed.

III. FEDERATED SELF-SUPERVISED LEARNING WITH SELECTIVE DATA CONTRAST

This paper proposes a framework SOFed aiming at training a global model to extract good visual representations from unlabeled streaming data in a network of edge devices with limited storage. This is the first work on self-supervised on-device federated learning from unlabeled streams. To reduce the storage overhead on edge devices while learning representations efficiently, we propose a data selection method based on importance scoring. The data selection method facilitates the edge device to select the most representative samples into the local buffer and benefits representation learning.

In this section, we will first introduce the overview of the proposed federated contrastive learning with coreset selection framework in Section III-A. Then we will present the local contrastive learning process in Section III-B. After that, we will introduce how to perform coreset selection on each device in Section III-C.

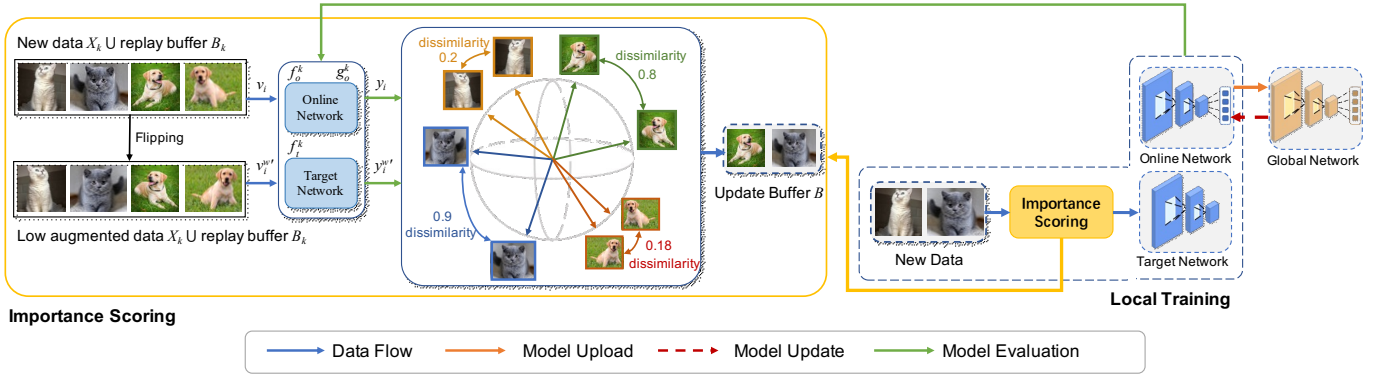


Fig. 3: Local contrastive learning process with replay buffer on the k -th device. The replay buffer employs importance scoring for online data selection.

A. Framework Overview

As shown in Fig. 2, there are two stages in the proposed FCL framework SOFed. In the first stage, the model is collaboratively trained on streaming unlabeled data by distributed devices to generate good visual representations. An additional classifier is then trained on top of the online encoder from the learned model with few (e.g. 1%) labeled data, in the second stage. In the rest of this paper, we focus on the FCL framework in the first stage to obtain better visual representations.

In the first stage, there is a network containing a server and many distributed devices (i.e., clients). The server is used to coordinate the process of learning, aggregate the model updates and distribute the model to clients. SOFed is performed round by round. There are four steps in one round. (1) First, the server sends the global model to devices to initialize the online network. (2) Then, each client conducts local training. As data keep streaming in, each device independently adopts a coreset selection policy based on importance score to maintain a replay buffer. More specifically, in round r , once a device k receives a new batch of samples X_k , it evaluates the importance score of each sample in both X_k and the local buffer B_k . For simplicity, we assume the size of X_k and the size of B_k are the same for each k . The device k selects $|B_k|$ data with the highest scores from $B_k \cup X_k$ and updates buffer B_k with these data. With this importance-based data selection process, each device can maintain the most valuable samples from streaming data for model training, reducing storage costs and enhancing the performance of the global model simultaneously. Local models are trained on the local replay buffers using contrastive learning to extract better visual representations. (3) After training, the online neural networks in local models are uploaded to the server. (4) Finally, the server aggregates the clients' online neural networks to update the global model using as in [10]. The details of local coreset selection and the model update will be discussed in the next subsections.

B. Local Model Architecture

The local contrastive learning process is shown in Fig. 3. BYOL architecture is applied as the local model since its loss function of one sample only uses the representations of this sample and does not depend on other samples. This decoupled

formulation of loss function makes it easier to measure and interpret the importance of each sample. Different from BYOL, other popular contrastive learning models such as SimCLR, MoCo [5], [6] use other samples in the same mini-batch when computing the loss of one sample. This correlation makes it challenging to measure the importance of each sample. To be more specific, SimCLR and MoCo use a contrastive loss which minimizes the similarity of representations between different samples:

$$\mathcal{L}(x) = -\log \frac{\exp(\text{sim}(y_o, y_t')/\tau)}{\sum_{i \in Q} \exp(\text{sim}(y_o, y_t^{i'})/\tau)}, \quad (2)$$

where $y_t^{i'}$ is the output of target network of x_i in a batch or a memory bank Q . In contrast to (1), the denominator in (2) involves other samples. Therefore, this contrastive loss is prone to be affected by other samples and is not preferable to evaluate the performance of an individual sample.

For each device k , the local model contains an online model with an online encoder f_o^k and an online predictor g_o^k and a target model with a target encoder f_t^k . The model parameters of the online network and the target network can be denoted as $W_o^k = (W_o^k, W_p^k)$ and W_t^k , respectively. We utilize the online model as the main model to be aggregated at each round and as the backbone to generate representations for the second stage. To train the online model from the input data stream, we adopt importance scoring to select the most representative data into the replay buffer and train the online model on the replay buffer. The target model is updated by EMA with the online model: $W_t^k = \tau W_o^k + (1 - \tau)W_t^k$, where the EMA weight τ is empirically set to 0.99.

C. Local Coreset Selection Policy

Importance Scoring. We calculate the importance scores of each sample to indicate its effectiveness during training. For each input x_i on the k -th device, the importance scoring function $S(x_i)$ aims to measure the quality of the generated representation $y_i = g_o^k(f_o^k(x_i))$. Intuitively, the input x_i with bad representation is more valuable to be kept for further training since it will provide new knowledge for the model to learn. To achieve this, we generate a weakly augmented view $x_i^{w'}$ of each sample x_i from buffer B_k and data stream

X_k . We feed x_i and $x_i^{w'}$ to the online network and the target network respectively. The output projected representations can be denoted as y_i and $y_i^{w'}$, where $y_i^{w'} = f_t^k(x_i^{w'})$. Ideally, if the network has learned to generate effective representations of x_i , the representations will be similar. Hence, we use BYOL loss as the importance score of this positive pair $\{x_i, x_i^{w'}\}$. The importance scoring function $S(\cdot)$ can be denoted as:

$$\begin{aligned} S(x_i) &= \text{dissim}(x_i, x_i^{w'}) \\ &= 1 - \text{similarity}(\bar{y}_i, \bar{y}_i^{w'}) \\ &= 1 - \frac{\langle y_i, y_i^{w'} \rangle}{\|y_i\|_2 \cdot \|y_i^{w'}\|_2}, x_i \in \{B_k \cup X_k\} \end{aligned} \quad (3)$$

where \bar{y}_i and $\bar{y}_i^{w'}$ are ℓ_2 -normalized vectors to enforce $\|y_i\|_2 = \|y_i^{w'}\|_2 = 1$. Therefore, the importance score $S(x_i)$ is positive and in the range $[0, 2]$. The BYOL loss can be interpreted as the dissimilarity between the projected representation vectors $y_i, y_i^{w'}$ of x_i and its weakly augmented view $x_i^{w'}$. If the dissimilarity between y_i and $y_i^{w'}$ is significant, the importance score will be large. Otherwise, the importance score will be approximately zero. From the perspective of model effectiveness, the representations generated by views of a positive pair need to be invariant to the transformation. Hence, the two representations generated by a more efficient model are more similar. Since a higher importance score represents a larger dissimilarity, input x_i with a higher score indicates that the model has not effectively learned it. Therefore, learning from inputs with higher scores can provide more valuable information to the model. Minimizing the contrastive loss in (1) to update the online network with input x_i aims at minimizing the dissimilarity between two strongly augmented views. Thus the importance score of x_i in (3) will also decrease, resulting in a lower possibility to select x_i . Samples that will generate good representations tend to be discarded and the more valuable samples will be kept for further learning. In this way, the dynamics of the importance scoring also enable the online model to update with more valuable data consistently.

Note that different augmentation methods can induce variations in generating a positive pair $\{x_i, x_i^{w'}\}$ for importance scoring. However, importance scoring aims at evaluating the model effectiveness of each sample independently from transformations. To eliminate randomness, we use a single augmentation method to generate weakly augmented views. Horizontal flipping is utilized as the default augmentation method. Hence, the importance scores $S(\cdot)$ of different images can be compared fairly, since $S(\cdot)$ is only deterministic to the input image.

Compatibility with Other Contrastive Learning Architectures. Although we use BYOL loss to introduce the importance score, (3) can also be implemented with other contrastive learning architectures. Since contrastive learning architectures can be unified as a Siamese network containing an online network and a target network [15], the importance scoring in Fig. 3 can also be generalized to other contrastive learning frameworks. To be more specific, contrastive learning methods train a Siamese network constituted by an online network and a target network. For example, SimCLR [5] method shares the same weights between these two networks, and MoCo

[6], as well as BYOL, updates the target model with the moving average of the online encoder. With different contrastive learning methods, we can feed the x_i and its flipped view $x_i^{w'}$ to online and target networks respectively and compute the importance score based on the projected representations as (3).

Importance Scoring Based Coreset Selection. The client selects the most effective samples into its buffer according to the importance scores. Once receiving a new segment of input X_k on device k , the client k has to make an online decision on how to update the buffer B_k with the new samples. The goal of this process is to select and store the valuable data, i.e., a coreset, into the buffer for model training in the coming rounds. To achieve this, we apply the importance scoring function $S(\cdot)$ on each sample from $B_k \cup X_k$. The buffer B_k is then updated to B'_k by selecting data with the highest importance scores in $B_k \cup X_k$:

$$B'_k = \{x_i | x_i \in B_k \cup X_k, i \in \text{topN}(\{S(x_i)\}_{i=1}^{2N})\}, \quad (4)$$

where N is the size of B and X , and $\text{topN}(\cdot)$ returns the indices of x_i with N largest importance scores.

The main steps of the federated contrastive learning with importance scoring based coreset selection are summarized in Algorithm 1. At each round, the clients conduct local training in parallel. For a client C_k , the local model is first initialized by the global model (W^g, W_p^g) . Suppose that on each client, local data have the same velocity v . At the r -th round, with a sequence of online data $X_{k,r}^i$, where $i = 1, 2, \dots, v$, the client C_k calculates the importance scores using (3) of data from both $X_{k,r}^i$ and B_k . Then it updates the buffer B_k using (4) to maintain the most important data in the buffer. The local online model $W_o^k = (W^k, W_p^k)$ is trained using the local buffer B_k . The local target model W_t^k is updated by EMA with the local online model W^k . After local training, each client C_k uploads the online model W^k, W_p^k to the server. The server aggregates the local models to update the global model W^g, W_p^g as $W^g \leftarrow \sum_{k=1}^K \frac{1}{K} W^k$ and $W_p^g \leftarrow \sum_{k=1}^K \frac{1}{K} W_p^k$. The global model can gain better representations after R communication rounds. After that, the global encoder W^g can be applied to downstream tasks.

D. Lazy Scoring

Due to the limited computation resources on edge devices, we utilize the lazy scoring technique to minimize the computation overhead induced by the importance scoring [26]. By adopting lazy scoring, the client reuses the scores of samples in the buffer. It is based on two observations. First, the score $S(x_i)$ of the sample x_i changes slightly in successive iterations. Since the importance score is only dependent on the sample itself and the Siamese model, the score changes slowly with a slowly updated model. Therefore, the client can utilize the scores computed several iterations before. Second, when a new data batch arrives, the importance scoring-based coreset selection method maintains most of the data in the buffer and discards most of the new data. As will be demonstrated in Section IV, over 80% new data are directly dropped. Since we assume the size of the new data batch is the same as the buffer, over

Algorithm 1 SOFed FCL Framework

Input: clients $\mathbb{C} = \{C_1, \dots, C_K\}$ with local buffers $\{B_1, \dots, B_K\}$, learning rate η , weight decay m , local data velocity v , the number of communication rounds R and EMA weight τ .

Output: global encoder W^g .

- 1: Initialize W^g, W_p^g
- 2: **for** each round $r = 1, 2, \dots, R$ **do**
- 3: **for** client $C_k \in \mathbb{C}$ in parallel **do**
- 4: $W^k, W_p^k \leftarrow \text{LocalTrain}(W^g, W_p^g)$
- 5: **end for**
- 6: $W^g \leftarrow \sum_{k=1}^K \frac{1}{K} W^k$
- 7: $W_p^g \leftarrow \sum_{k=1}^K \frac{1}{K} W_p^k$
- 8: **end for**
- 9: **return** W^g
- 10: **LocalTrain**(W^g, W_p^g)
- 11: $W^k \leftarrow W^g, W_p^k \leftarrow W_p^g$
- 12: **for** each $X_{k,r}^i, i = 1, 2, \dots, v$ **do**
- 13: **for** each x_j in $X_{k,r}^i \cup B_k$ **do**
- 14: Calculate $S(x_j)$ using (3)
- 15: **end for**
- 16: Update B_k using (4)
- 17: $(W^k, W_p^k) \leftarrow W_o^k - \eta \nabla \mathcal{L}(W_o^k; B_k) - \eta m W_o^k$
- 18: $W_t^k \leftarrow \tau W_t^k + (1 - \tau) W^k$
- 19: **end for**

80% data in the buffer are preserved at each round. Therefore, reusing the scores can reduce lots of computation.

To be more specific, on the k -th client, we recompute the score of the sample in the buffer B_k every T iterations, where T is the interval of lazy scoring. For the sample x_i on the k -th client, we record the number of iterations that it has stayed in the buffer as $\text{age}(x_i)$. When a new data batch X_k arrives, for the samples in the buffer, we only update the ones that haven't been recomputed in T iterations. For buffer B_k , the scores can be updated as:

$$S_k^t = \begin{cases} \text{dissim}(x_i, x_i^{w'}), & \text{if } \text{age}(x_i) \bmod T = 0 \\ S_k^{t-1}(x_i), & \text{otherwise} \end{cases}, \quad (5)$$

where $S_k^{t-1}(x_i)$ is the score of x_i from the last iteration. If the score of x_i needs to be updated, the score is computed by (3). Otherwise, the score remains the same as the last iteration. For example, suppose T is 5 and sample x_i is in the new data batch X_k . When x_i first gets to the k -th client at the t -th iteration, its score $S_k^t(x_i)$ is computed and restored. The age of x_i is initialized as 0. In the next 4 iterations, the score of x_i will not be updated, while the age of x_i increases. The buffer is updated according to the old score $S_k^t(x_i)$ of x_i as (4). If x_i maintains in the buffer for 4 iterations, at iteration $t+5$, since $\text{age}(x_i) \bmod T = 0$, the score of x_i is updated. The buffer is updated according to the new score $S_k^{t+5}(x_i)$ at iteration $t+5$. The recomputing ratio of the importance scores can be reduced to about $1/T$ by utilizing lazy scoring. Therefore, the computation cost can be reduced and the scoring efficiency can be improved.

IV. EXPERIMENTS

In this section, we evaluate the performance of the visual representation learned by SOFed on CIFAR-10, CIFAR-100, and SVHN. We first explain the experiment setup. Then, we evaluate the enhanced performance of the proposed FCL framework. After that, we evaluate the improved accuracy and the higher learning speed of the proposed coreset selection policy in the proposed framework. Finally, we evaluate the impacts of lazy scoring.

A. Experimental Setup

Datasets and Evaluation Protocols. We use CIFAR-10, CIFAR-100 [27] and SVHN [28] datasets to evaluate the proposed approach. CIFAR-10 and CIFAR-100 contain 10 classes and 100 classes of an equal number of images per class, respectively. SVHN consists of 73,257 training and 26,032 testing images in 10 classes. We use the proposed SOFed method to train the model by distributed devices without labels. To evaluate the performance of the representations generated by the encoder in the trained model, we use linear evaluation as described in [29]. We freeze the parameters in the encoder and train a classifier on top of it with 1%, 10%, or 100% labeled data. We report the classification accuracy on different datasets.

Federated and Continual Learning Setting. To simulate the temporally correlated data stream, the data stream is temporally ordered by class. We use Strength of Temporal Correlation (STC) to measure the level of temporal correlation of the input stream [19]. STC represents the number of consecutive data in the input stream that are from the same class. A higher STC implies a higher level of temporal correlation. We use 5 devices for federated learning and divide the training dataset into 5 partitions. The temporally reordered dataset is randomly split and distributed to different devices. At each communication round, we set the local data velocity to be $v = \frac{\text{\#training samples}}{5}$. To simulate a more non-IID situation, we increase the STC and distribute 2 classes and 20 classes to each client in CIFAR-10 and CIFAR-100 respectively.

Training Setting. We use ResNet-18 as the architecture for encoders and use a multi-layer perceptron (MLP) as the predictor. We use rean SGD optimizer to train the network with contrastive loss by default. For a fair comparison with other methods, we train the network for 300 rounds from scratch. Unless otherwise specified, the learning rate η is set to 0.06, and the default batch size $|B|$ is 128 with a weight decay m of 0.0001.

Baselines. The proposed coreset selection policy is referred to as *Importance Scoring (IS)* method. We implement three more unlabeled data selection methods to select coreset into the replay buffer from the input data stream on the fly. *Random Replacement (RR)* is a continual learning variant of reservoir sampling [19]. It randomly selects data from the input batch and the buffer to update the buffer. *First-In, First-Out (FIFO) Replacement* is recently used for incremental batch learning [19]. FIFO replacement replaces the oldest samples in the buffer with new samples. Although these two methods are straightforward, they yield superior results to maintain

TABLE I: Accuracy on CIFAR-10, CIFAR-100, SVHN datasets under different non-iid settings with different FCL frameworks. The strength of Temporal Correlation (STC) is the number of consecutive data in the input stream that are from the same class. A higher STC represents a more non-iid situation.

Method	CIFAR-10		CIFAR-100		SVHN	
	STC 500	STC 2500	STC 500	STC 2500	STC 500	STC 2500
FedU-FIFO	76.61	73.74	45.99	46.20	91.72	88.86
FedSimCLR-FIFO	67.21	65.45	38.37	38.40	82.73	81.14
FedU-RR	76.88	73.63	46.19	46.88	91.51	89.32
FedSimCLR-RR	68.20	66.45	40.21	40.01	82.30	81.06
FedU-IS	77.71	74.71	49.24	48.93	92.10	87.41
FedSimCLR-IS	73.39	70.50	44.70	45.26	88.94	84.51
SOFed (ours)	78.66	74.98	49.58	49.95	92.20	89.46
BYOL (Centralized upper-bound)	80.50	76.80	50.45	50.78	92.61	89.80

data without labels in continual learning [30]. *K-center* is a SOTA active learning method to select coreset without labels. It solves a *k*-center problem in the feature space to select the data with the largest impacts [31]. We compare the IS coreset selection policy with different data selection policies in the proposed FCL framework. Since SOFed is the first federated contrastive learning framework with a small data buffer on each client, to demonstrate the efficacy of the proposed FCL framework, we compare the performance of SOFed with baseline methods that combine FCL frameworks with different data selection methods. We implement two more recently proposed FCL methods. *FedSimCLR* is implemented by combining FedAvg with SimCLR [5]. *FedU* [13] is a SOTA FCL framework that addresses the non-iid data problem. Among the data selection methods, IS is the best policy and RR, FIFO also have promising performance. Therefore, we combine the aforementioned FCL methods with IS, RR and FIFO for a fair comparison. We use *FedU-IS*, *FedSimCLR-IS*, *FedU-RR*, *FedSimCLR-RR*, *FedU-FIFO*, and *FedSimCLR-FIFO* as baselines to compare with SOFed. We also compare with a centralized contrastive learning method *BYOL* [8], since the local model in SOFed is based on it. The centralized BYOL is implemented with the proposed importance scoring method to enable an online coreset selection with temporally correlated data stream as in [26]. Kindly note that this is the first work on self-supervised on-device federated learning from unlabeled streams, which is extended from our previous work in [26] for a single edge device. Different from conventional self-supervised learning settings, the streaming data are temporally correlated and the on-device buffer size is limited. To have a fair comparison with existing works, we reproduce their results using the same setup as our methods.

B. Effectiveness of SOFed

We evaluate the proposed FCL framework SOFed and other FCL frameworks on CIFAR-10, SVHN and CIFAR-100 by training an additional classifier with 100% labeled data. We report the accuracy of the proposed method and other baselines in Table I. When STC is set to 500, the data are distributed to clients in a weakly non-iid manner. When STC is set to 2500, a client is assigned to data from four classes. Therefore, the data distribution is more non-iid. On CIFAR-10 in a weakly non-iid setting, the proposed approach achieves 78.66%

top-1 accuracy, only 1.84% below the upper bound method centralized BYOL. Although the proposed approach does not significantly outperform the combination of FedU and the proposed coreset selection policy IS (i.e., FedU-IS), SOFed can achieve competitive performance without the communication module and the empirical hyperparameter selection process in FedU. The proposed framework outperforms FedSimCLR-IS, FedSimCLR-RR, FedU-RR, FedSimCLR-FIFO, FedU-FIFO by {5.27%, 10.46%, 1.78%, 11.45%, 2.05%} respectively. The proposed framework also achieves the best accuracy in non-iid setting on CIFAR-10, where it outperforms FedSimCLR-IS, FedSimCLR-RR, FedU-RR, FedSimCLR-FIFO, FedU-FIFO by {4.48%, 8.53%, 1.35%, 9.43%, 1.24%} respectively. The results demonstrate that the proposed approach achieves improved accuracy on various datasets and collaborative learning settings. Since the experiments are conducted in an online learning scenario where temporally correlated data stream in, the vanilla method dealing with the data stream is FIFO. FIFO facilitates the models to be trained directly on the new data without selection. So the FCL scheme without core dataset selection in an online learning scenario is the conventional FCL schemes [13], [15] combined with FIFO. Compared with FedSimCLR-FIFO and FedU-FIFO, the proposed method SOFed exhibits better accuracy than FCL schemes without core dataset selection.

TABLE II: Accuracy on CIFAR-10 with different FCL frameworks and various numbers of devices.

Method	Number of devices		
	10	50	100
FedU-FIFO	78.75	70.47	68.12
FedSimCLR-FIFO	67.43	61.58	57.39
FedU-RR	77.25	75.09	74.53
FedSimCLR-RR	69.72	62.95	58.90
FedU-IS	80.26	78.02	75.76
FedSimCLR-IS	73.30	62.09	57.49
SOFed (ours)	80.79	78.49	76.46

To demonstrate the scalability of the proposed method, we evaluate the proposed FCL framework with 10, 50, and 100 clients on CIFAR-10 when STC is set to 200. We report the accuracy of the proposed method and other baselines in Table II. The performance of the proposed method is robust and outperforms the baselines with various numbers of

devices. When the number of devices is 10, SOFed outperforms FedSimCLR-IS, FedU-IS, FedSimCLR-RR, FedU-RR, FedSimCLR-FIFO, FedU-FIFO by {7.49%, 0.53%, 11.07%, 3.54%, 13.36%, 2.04%} respectively. As discussed before, although SOFed does not outperform FedU-IS with a clear margin, our proposed framework can achieve competitive performance without the empirical hyperparameter selection process in FedU. The decent performance of FedU-IS can also demonstrate the compatibility of our coreset selection policy with other FCL paradigms. As shown in Table II, all the approaches degrade when the number of devices becomes larger. However, the degradation of the proposed approach is substantially smaller than other approaches. The accuracy of SOFed drops 4.33% when the count of devices increases from 10 to 100, while the accuracy of FedSimCLR-IS drops 15.81%. The proposed framework is robust in different scales of edge devices network.

C. Effectiveness of Coreset Selection Policy

Improved Accuracy with Different Labeling Ratios. We compare the proposed coreset selection approaches by substituting the importance scoring for *Random Replacement*, *FIFO Replacement*, *K-Center* in the proposed FCL framework. We first perform federated contrastive learning with unlabeled data under different iid settings (i.e., different STC values), then train a classifier on top of the fixed encoder with 1% or 10% labeled data.

TABLE III: Accuracy on CIFAR-10 dataset under different non-iid settings with different coreset selection methods. We train a classifier with 1% and 10% labeled data on top of the encoder trained with FCL.

Method	CIFAR 10			
	STC 500		STC 2500	
	1%	10%	1%	10%
Random Replacement	59.25(-2.74)	70.53(-2.73)	52.71(-2.65)	66.76(-1.77)
FIFO Replacement	57.78(-5.21)	70.6(-2.60)	53.66(-2.70)	66.53(-2.00)
K-Center	53.87(-8.12)	63.44(-9.76)	50.63(-5.73)	62.56(-5.97)
Importance Scoring (ours)	61.99	73.2	56.36	68.53

The proposed coreset selection approach by importance scoring outperforms the SOTA baselines. The accuracy with different labeling ratios (i.e., 1%, 10%) on CIFAR-10 is shown in Table. III. With 1% and 10% labeled data for learning the classifier, the proposed importance scoring achieves the accuracy of 61.99%, 73.2% with STC 500, 56.36%, and 68.53% with STC 2500. The proposed method demonstrates better accuracy in different iid settings. When STC is set to 500, importance scoring outperforms other three approaches by {2.74%, 5.21%, 8.12%}, {2.73%, 2.6%, 9.76%} training the classifier with 1% or 10% labeled data respectively. When STC is set to 2500, the proposed method outperforms other approaches by {2.65%, 2.7%, 5.73%}, {1.77%, 2.0%, 5.97%} with 1% or 10% labeled data in the linear evaluation respectively. The proposed method learns better visual representations with different data distributions and demonstrates enhanced accuracy after training a classifier with different number of labeled data.

The learning curves of the proposed approach and the two best baselines trained on CIFAR-10 with different iid settings

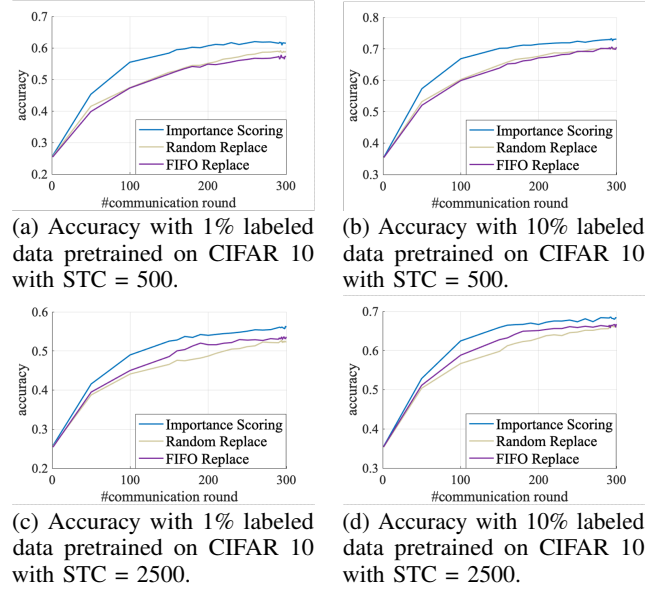


Fig. 4: Accuracy on CIFAR-10 with 1% and 10% labeled data and different STC values.

TABLE IV: Accuracy on CIFAR-10 dataset with different buffer sizes.

Buffer Size	Method	Accuracy
32	Random Replacement	67.07(-2.18)
	FIFO Replacement	67.72(-1.53)
	K-Center	62.73(-6.52)
	Importance Scoring (ours)	69.25
64	Random Replacement	67.43(-2.35)
	FIFO Replacement	67.63(-2.15)
	K-Center	66.65(-3.23)
	Importance Scoring (ours)	69.78
128	Random Replacement	66.70(-3.47)
	FIFO Replacement	64.89(-5.28)
	K-Center	68.56(-1.61)
	Importance Scoring (ours)	70.17
256	Random Replacement	62.51(-5.69)
	FIFO Replacement	61.73(-6.47)
	K-Center	65.80(-2.40)
	Importance Scoring (ours)	68.20

are shown in Fig. 4. The accuracy is obtained by training a classifier on top of the learned visual representations with 1% or 10% labeled data. The learning curves represent the speed of the model to learn high-quality visual representations. We compare the proposed coreset selection method with the two most competitive baselines. The proposed coreset selection policy quickly learns data representations and achieves a significantly faster learning speed and a higher accuracy than the baselines. When the federated contrastive learning is trained with STC equal to 500 and the linear evaluation is performed with 10% labeled data, the accuracy of the proposed approach quickly increases to higher than 70% at the 150-th round, which is 1.86 \times faster than the random replacement policy that achieves 70.12% at the 280-th round. The proposed approach achieves a faster learning speed than the baselines.

Improved Accuracy with Different Buffer Sizes. We

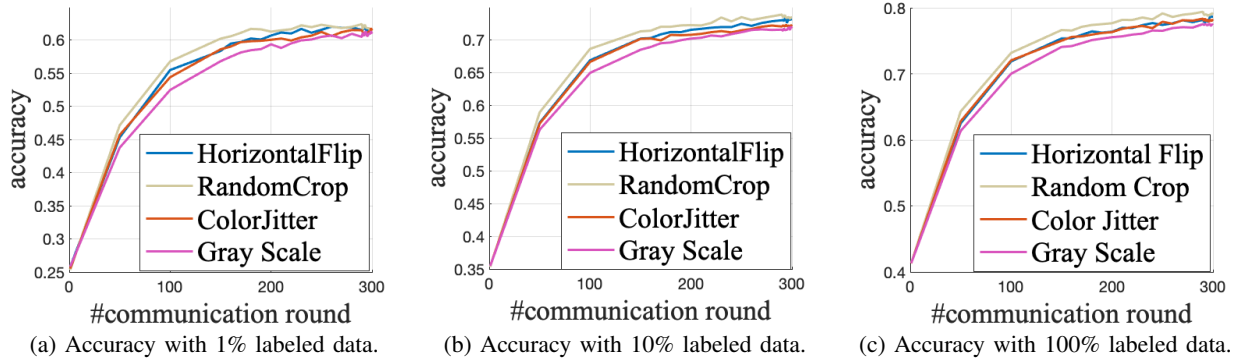


Fig. 5: Learning curves on CIFAR-10 datasets with 1%, 10%, and 100% labeled data using different augmentation methods.

evaluate the impact of buffer size on the performance of the proposed approach. The model is trained on the CIFAR-10 dataset with a buffer size in $\{32, 64, 128, 256\}$ for 100 rounds. The corresponding learning rate is scaled to $\{0.04, 0.05, 0.06, 0.08\}$.

As shown in Table IV, the proposed approach outperforms the baselines under different buffer sizes. The results show that with different buffer sizes, the proposed approach maintains significantly better than the baselines. When the buffer size is larger, on one hand, the framework can train the model on more informative data. On the other hand, with the same number of input data, the number of training times decreases. Hence, the accuracy doesn't change monotonically with the buffer size. The proposed method achieves the highest accuracy with a buffer size of 128.

TABLE V: Accuracy on CIFAR-10 dataset with different weak augmentation methods in importance scoring. We train a classifier with 1%, 10%, and 100% labeled data on top of the encoder trained with FCL.

Augmentation Method	Accuracy		
	1%	10%	100%
Horizontal Flipping	61.99	73.21	78.66
Random Cropping	62.38	73.81	79.39
Gray Scale	61.31	71.97	77.58
Color Jittering	61.79	72.29	78.39

Robustness with Different Augmentation Methods. Since strong augmentation can introduce randomness into importance scoring, we use weak augmentation to generate views for the evaluation of importance scores. We evaluate the impact of different augmentation methods on the performance of the proposed approach. The model is trained on the CIFAR-10 dataset for 300 rounds with different augmentation methods in importance scoring and various percentages of labeled data in classifier training.

The selection of augmentation methods has a marginal impact on the performance of the proposed approach as shown in Table V. The learning curves of the proposed framework with different augmentation methods are displayed in Fig. 5. When the default augmentation method, horizontal flipping, is applied as the augmentation method in importance scoring, the proposed framework achieves the accuracy of 61.99%, 73.21%, and 78.66% with 1%, 10%, and 100% labeled data respectively. When random cropping is selected as the augmentation method in importance scoring, the proposed framework achieves the

accuracy of 62.38%, 73.81%, and 79.39% with 1%, 10%, and 100% labeled data respectively. The differences between the accuracy of these two augmentation methods are 0.39%, 0.6%, and 0.63% with different percentages of labeled data. The differences in terms of accuracy between the framework using horizontal flipping and gray scaling or color jittering are also around 1%. Therefore, the selection of the weak augmentation approach is insignificant in the proposed framework.

D. Validation of Coreset Selection Efficacy

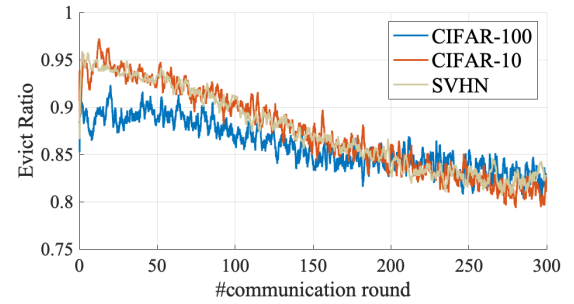


Fig. 6: The ratio of new data batch directly discarded in importance scoring on CIFAR-10, CIFAR-100, and SVHN after smoothing.

We evaluate the ratio of data that is directly discarded during training by using the proposed importance scoring-based coreset selection policy. The model is trained on CIFAR-10, CIFAR-100, and SVHN datasets for 300 rounds with buffer size 128 and STC value 500 in a centralized manner. Fig. 6 demonstrates the ratio of data that is directly dropped in each iteration on CIFAR-10, CIFAR-100, and SVHN datasets using importance scoring. Note that the curve of ratio with respect to iterations is smoothed using a moving average filter for better visibility.

We can gain two insights into the coreset selection policy from these results. First, over 80% of the data in the new data batch are dropped, hence, most data in the buffer are preserved. By incorporating more information from previous iterations, the proposed coreset selection policy can preserve knowledge learned before. Therefore, the model trained on continual non-iid data stream can learn more effective visual representations. Second, the ratio of new data directly discarded decreases with respect to iteration. For Random Replacement, the ratio is around 50% at each iteration. If the ratio is approaching 50%,

the coreset selection method could not discriminate data with various levels of importance, which implies the model has been fully trained. So the evict ratio of new data decreases during training as a result of a more accurate model.

E. Impacts of Lazy Scoring

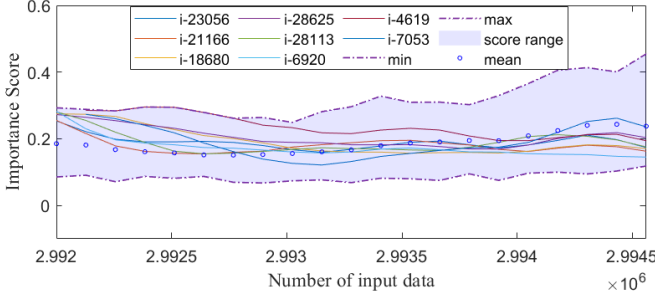


Fig. 7: Importance scores of data that stay in the replay buffer longer than the lazy scoring interval. The continuous line “i-index” is the trend of importance score of the *index*-th sample. The blue shadow represents the range of scores of samples in the buffer.

We evaluate the importance scores of the same samples in the replay buffer during training on CIFAR-10 when the lazy scoring interval T is 20. The score of each sample changes slowly in successive iterations. The trends of importance scores of samples in the buffer of the first device in 20 iterations are shown in Fig. 7. The x -axis is the number of seen inputs on one device and the y -axis is the importance score. Since some samples are discarded during these 20 iterations, for clarity, only the scores of the samples that stay in the replay buffer for the whole 20 iterations are depicted in the figure. The continuous line “i-index” represents the importance score of the *index*-th sample. The two dot lines are the minimal and maximal scores of the samples in both the buffer and the new data batch. Hence, the blue shadow represents the range of scores of samples in the buffer. The blue circle represents the mean score of the samples in each iteration.

In the first 10 iterations, the importance scores change slowly and remain higher than the mean score. Then, the scores of the less representative samples slowly decrease. Thus, these less representative samples are recomputed and will be discarded, while some of the old and representative samples are kept in the buffer. The results verify the motivation of the lazy scoring and rationalize its effectiveness.

We also evaluate the accuracy and the reduced computation overhead by utilizing lazy scoring. The model is trained on the CIFAR-10 dataset with a buffer size of 128.

Table VI demonstrates the accuracy, the rescoring ratio of buffer, and relative batch time with different lazy scoring intervals T . When the interval increases, the average ratio of data in the buffer that needs to be recomputed decreases. Hence, the time to process a new data batch is reduced. We aim to achieve the same level of accuracy when lazy scoring is enabled as when it is disabled. With the re-scoring interval T set to smaller values 4 or 10, the accuracy is very similar to

TABLE VI: Accuracy, rescoring ratio and batch time on CIFAR-10 dataset with different lazy scoring intervals.

Interval	Accuracy(%)	Rescoring Ratio(%)	Relative Batch Time(%)
disabled	78.66	100	100
4	78.48 (-0.18)	19.53	80.86
10	78.45 (-0.21)	7.42	74.33
20	79.18 (+0.52)	3.03	73.26
50	79.79 (+1.13)	1.54	72.20
100	76.12 (-2.54)	0.73	70.92

that with lazy scoring disabled, as shown in Table VI. When T is further increased to 20 and 50, lazy scoring can slightly increase accuracy. The increase in accuracy may be attributed to the incorporation of previous knowledge which could be beneficial to the data selection [6]. With a lazy scoring interval of 50, the accuracy increases by 1.13%, and the rescoring ratio is reduced to 1.54%. Therefore, with a proper value of T , the accuracy can be improved by the stability of samples in the buffer. However, if the interval keeps increasing, the score that the sample possesses may be stale and cannot represent the sample’s quality well. The accuracy decreases when the interval T is set to 100, as shown in Table VI. So the hyper-parameter T should be empirically set to 50 to boost the accuracy and reduce the computation overhead at the same time.

V. CONCLUSION

This work aims to enable contrastive learning from input streaming data in a network of edge devices. We propose a framework to maintain a small data buffer filled with the most representative data for training for collaborative visual representation contrastive learning. To achieve the online selection of streaming unlabeled data, we propose a coreset selection to maintain the data buffer based on importance scores. Experiments show superior accuracy of the proposed framework compared with state-of-the-art methods.

VI. ACKNOWLEDGEMENT

This work is supported in part by National Science Foundation under grants CNS-1822099, CNS-2122220, and CNS-2122320.

REFERENCES

- [1] R. A. Khalil, N. Saeed, M. Masood, Y. M. Fard, M.-S. Alouini, and T. Y. Al-Naffouri, “Deep learning in the industrial internet of things: Potentials, challenges, and emerging applications,” *IEEE Internet of Things Journal*, vol. 8, no. 14, pp. 11 016–11 040, 2021.
- [2] T. Choi, O. Would, A. Salazar-Gomez, and G. Cielniak, “Self-supervised representation learning for reliable robotic monitoring of fruit anomalies,” in *2022 International Conference on Robotics and Automation (ICRA)*, 2022, pp. 2266–2272.
- [3] C. Doersch, A. Gupta, and A. A. Efros, “Unsupervised visual representation learning by context prediction,” in *2015 IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 1422–1430.
- [4] M. Noroozi and P. Favaro, “Unsupervised learning of visual representations by solving jigsaw puzzles,” in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer, 2016, pp. 69–84.

- [5] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *Proceedings of the 37th International Conference on Machine Learning*. PMLR, 2020, pp. 1597–1607.
- [6] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9729–9738.
- [7] X. Chen and K. He, "Exploring simple siamese representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 15 750–15 758.
- [8] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar, B. Piot, k. kavukcuoglu, R. Munos, and M. Valko, "Bootstrap your own latent - a new approach to self-supervised learning," in *Advances in Neural Information Processing Systems*, 2020, pp. 21 271–21 284.
- [9] L. Jing and Y. Tian, "Self-supervised visual feature learning with deep neural networks: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 11, pp. 4037–4058, 2021.
- [10] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, "Communication-Efficient Learning of Deep Networks from Decentralized Data," in *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*. PMLR, 2017, pp. 1273–1282.
- [11] Y. Wu, D. Zeng, Z. Wang, Y. Sheng, L. Yang, A. J. James, Y. Shi, and J. Hu, "Federated contrastive learning for dermatological disease diagnosis via on-device learning (invited paper)," in *2021 IEEE/ACM International Conference On Computer Aided Design (ICCAD)*. IEEE, 2021, p. 1–7.
- [12] F. Zhang, K. Kuang, Z. You, T. Shen, J. Xiao, Y. Zhang, C. Wu, Y. Zhuang, and X. Li. (2020) *Federated Unsupervised Representation Learning*. [Online]. Available: <https://arxiv.org/abs/2010.08982>
- [13] W. Zhuang, X. Gan, Y. Wen, S. Zhang, and S. Yi, "Collaborative unsupervised visual representation learning from decentralized data," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 4892–4901.
- [14] Y. Wu, Z. Wang, D. Zeng, M. Li, Y. Shi, and J. Hu, "Decentralized unsupervised learning of visual representations," in *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence (IJCAI)*, 2022, pp. 2326–2333.
- [15] W. Zhuang, Y. Wen, and S. Zhang, "Divergence-aware federated self-supervised learning," in *International Conference on Learning Representations*, 2022.
- [16] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra. (2018) *Federated Learning with Non-IID Data*. [Online]. Available: <https://doi.org/10.48550/arXiv.1806.00582>
- [17] H. Wang, M. Yurochkin, Y. Sun, D. Papailiopoulos, and Y. Khazaeni, "Federated learning with matched averaging," *International Conference on Learning Representations (ICLR)*, 2020.
- [18] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "SCAFFOLD: Stochastic controlled averaging for federated learning," in *Proceedings of the 37th International Conference on Machine Learning*. PMLR, 2020, pp. 5132–5143.
- [19] T. L. Hayes, N. D. Cahill, and C. Kanan, "Memory efficient experience replay for streaming learning," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 9769–9776.
- [20] R. Aljundi, M. Lin, B. Goujaud, and Y. Bengio, "Gradient based sample selection for online continual learning," in *Advances in Neural Information Processing Systems*, 2019, pp. 11 816–11 825.
- [21] Z. Borsos, M. Mutn̄y, and A. Krause. (2020) *Coresets via Bilevel Optimization for Continual Learning and Streaming*. [Online]. Available: <https://doi.org/10.48550/arXiv.2006.03875>
- [22] Y. Wu, Z. Wang, Y. Shi, and J. Hu, "Enabling on-device cnn training by self-supervised instance filtering and error map pruning," *International Conference on Compilers, Architecture, and Synthesis for Embedded Systems (CASES) in conjunction with Embedded Systems Week (ESWEEK)*, 2020.
- [23] A. Li, L. Zhang, J. Tan, Y. Qin, J. Wang, and X.-Y. Li, "Sample-level data selection for federated learning," in *IEEE INFOCOM 2021 - IEEE Conference on Computer Communications*. IEEE, 2021, p. 1–10.
- [24] C. Gong, Z. Zheng, F. Wu, B. Li, Y. Shao, and G. Chen. (2022) Online data selection for federated learning with limited storage. [Online]. Available: <https://doi.org/10.48550/arXiv.2209.00195>
- [25] L. Nagalapatti, R. S. Mittal, and R. Narayanam, "Is your data relevant?: Dynamic selection of relevant data for federated learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022, pp. 7859–7867.
- [26] Y. Wu, Z. Wang, D. Zeng, Y. Shi, and J. Hu, "Enabling on-device self-supervised contrastive learning with selective data contrast," in *2021 58th ACM/IEEE Design Automation Conference (DAC)*, 2021, pp. 655–660.
- [27] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [28] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," 2011.
- [29] A. Kolesnikov, X. Zhai, and L. Beyer, "Revisiting self-supervised visual representation learning," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1920–1929.
- [30] A. Chaudhry, M. Rohrbach, M. Elhoseiny, T. Ajanthan, P. K. Dokania, P. H. Torr, and M. Ranzato. (2019) Continual learning with tiny episodic memories. [Online]. Available: <https://doi.org/10.48550/arXiv.1902.10486>
- [31] O. Sener and S. Savarese, "Active learning for convolutional neural networks: A core-set approach," in *International Conference on Learning Representations (ICLR)*, 2018.