

# Towards Efficient Replay in Federated Incremental Learning

Yichen Li  
 Huazhong University of Science  
 and Technology, China  
 ycli0204@hust.edu.cn

Qunwei Li  
 Ant Group, China  
 qunwei.qw@antgroup.com

Haozhao Wang  
 Huazhong University of Science  
 and Technology, China  
 hz\_wang@hust.edu.cn

Ruixuan Li\*  
 Huazhong University of Science  
 and Technology, China  
 rxli@hust.edu.cn

Wenliang Zhong  
 Ant Group, China  
 yice.zwl@antgroup.com

Guannan Zhang  
 Ant Group, China  
 zgn138592@antgroup.com

## Abstract

*In Federated Learning (FL), the data in each client is typically assumed fixed or static. However, data often comes in an incremental manner in real-world applications, where the data domain may increase dynamically. In this work, we study catastrophic forgetting with data heterogeneity in Federated Incremental Learning (FIL) scenarios where edge clients may lack enough storage space to retain full data. We propose to employ a simple, generic framework for FIL named Re-Fed, which can coordinate each client to cache important samples for replay. More specifically, when a new task arrives, each client first caches selected previous samples based on their global and local importance. Then, the client trains the local model with both the cached samples and the samples from the new task. Theoretically, we analyze the ability of Re-Fed to discover important samples for replay thus alleviating the catastrophic forgetting problem. Moreover, we empirically show that Re-Fed achieves competitive performance compared to state-of-the-art methods.*

## 1. Introduction

Federated learning (FL) is a distributed framework that allows multiple edge clients to learn a unified deep learning model cooperatively while preserving the data privacy of the local clients [23, 31, 44]. Recently, FL has attracted growing attention and been applied to various fields such as recommendation systems [7, 24] and smart healthcare [35, 46].

Typically, FL has been actively studied in a static setting, where the number of training samples does not change over

time. However, in a realistic FL application, each client may continue collecting new data. It is difficult to learn new data while retaining previous information in machine learning due to the notorious phenomenon known as catastrophic forgetting [9], leading to performance degradation on previous tasks. This challenge is further compounded in FL settings, where the data in a client remains inaccessible to other ones and clients lack enough storage to retain full previous samples.

To address this issue, researchers have studied federated incremental learning (FIL), which enables each client to continuously learn from a local private and incremental task stream [4, 20]. The authors in [48] aim to personalize the models for each client by decomposing the model parameters into shared parameters and adaptive parameters, facilitating the transfer of common knowledge across similar tasks among clients. FCIL is proposed in [6] which specifically focuses on the federated class-incremental learning scenario and a global model is developed by incorporating additional class-imbalance losses. It is studied in [29] to utilize extra distilled data at both the server and client sides, where knowledge distillation is employed to mitigate catastrophic forgetting. FedCIL is proposed in [38] to learn a generative network and reconstruct past samples for replay, improving the retention of previous information.

While these approaches may be effective at learning new tasks, it is essential to also consider the constraints of privacy concerns and data heterogeneity. For example, data-reconstruction techniques with gradient or adversarial training are employed in FCIL and FedCIL to alleviate catastrophic forgetting, but it may cause privacy leakage of local data. Moreover, existing works simply assume that each client collects incremental data of different tasks in an independently and identically distributed (IID) manner, ignoring the issue of data heterogeneity in real-world scenarios.

\*Ruixuan Li is the corresponding author.

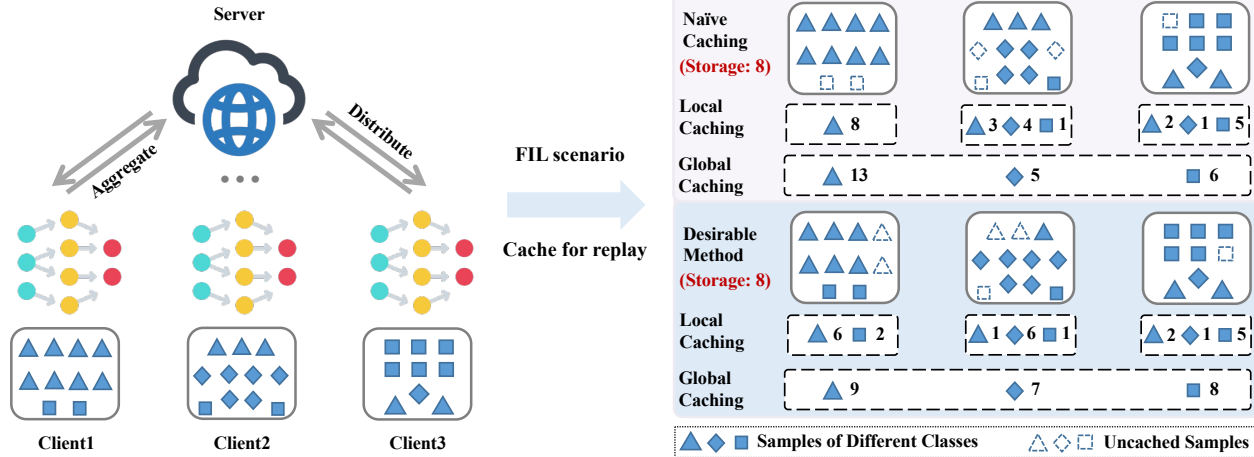


Figure 1. The motivation for our method: an example of 3-client in FIL scenario. When a new task arrives, each client needs to cache previous samples with limited storage for replay, alleviating catastrophic forgetting. Global caching represents all the samples cached by all clients collectively. With a naive caching method, the client may ignore the sample’s correlation across clients which increases the statistical data heterogeneity in global caching. With a desirable method, the client tends to cache samples which both considers the distribution of local samples and reduces the statistical data heterogeneity in the global caching.

In this paper, we investigate a simple and efficient method for catastrophic forgetting in FIL. We consider classification task in FIL and first identify two types of task for newly collected data: (1) Class-Incremental Task: the newly arrived data has different labels from previous data, and thus the label space of data is growing. (2) Domain-Incremental Task: the new data has a domain shift from previous data, and it does not change the label space of data. Then, we assume data heterogeneity as the data of each task in clients is Non-IID.

To tackle the catastrophic forgetting problem in the non-federated environment, data replay [32, 39] based methods have demonstrated great effectiveness by caching the important samples from previous tasks and replaying them when learning the new task. However, existing replay-based methods fail to consider the correlation between clients in FL, and the cached samples may not be globally optimal for tasks. Essentially, the cached samples should be strongly correlated with the statistical heterogeneity of all data across clients. The intuitive explanation can also be found in a simple 3-client example illustrated in Figure 1.

To explore this idea, we propose an efficient FIL framework named Re-Fed that can alleviate catastrophic forgetting by allowing all clients to synergistically cache data samples. More specifically, in Re-Fed, each client caches samples based not only on their importance in the local dataset but also on their correlation to the global dataset. Here we first employ an additional personalized informative model (PIM) for each client which can incorporate knowledge of data from global point of view in local caching such that the cached samples contribute to both local and global understanding of the data. Then, we quantify the sample

importance by calculating the gradient norm of previous local samples during the update of the PIM. Finally, the client caches the samples with higher importance scores, and trains the local model with both the cached samples from previous tasks and the samples from the new task.

Through extensive experiments on various datasets and two types of newly collected data (Class-Incremental Task and Domain-Incremental Task), we show that Re-Fed significantly improves the model accuracy compared to state-of-the-art approaches. The major contributions of this paper are summarized as follows:

- We are the first to study the problem of catastrophic forgetting with data heterogeneity in FIL. To address this problem, we propose a novel framework named Re-Fed which can be seen as an off-the-shelf personalization add-on for standard FIL and it inherits privacy protection and efficiency properties as traditional FL applications in FIL scenarios.
- Next, we theoretically show that Re-Fed can efficiently discover the important samples for data replay, with guaranteed convergence.
- Finally, we carry out extensive experiments on various datasets and different FIL task scenarios. Experimental results illustrate that our proposed model outperforms the state-of-the-art methods by up to 19.73% in terms of final accuracy on different tasks.

## 2. Related Work

**Federated Learning** FL is a technique to train a shared global model by aggregating models from multiple clients that are trained on their own local private datasets [23, 31,

44]. One effective architecture for FL is FedAvg [31], which optimizes the global model by aggregating the parameters of local models trained on private local data. However, traditional FL algorithms like FedAvg face challenges due to data heterogeneity, where the datasets in clients are Non-IID, resulting in degradation in model performance [15, 26]. To tackle the Non-IID issue in FL, a proximal term is introduced in optimization in [22] to mitigate the effects of heterogeneous and Non-IID data distribution across participating devices. Another approach of federated distillation [14], aims to distill the knowledge from multiple local models into the global model by aggregating only the soft predictions generated by each model. The authors in [25] proposed a knowledge distillation method that utilizes unlabeled training samples as a proxy dataset. Recently, there has been growing interest in data-free knowledge distillation methods that leverage adversarial approaches in generating data [45, 50, 51]. However, the aforementioned methods follow a framework designed to handle Non-IID static data with spatial heterogeneity, overlooking the potential challenges posed by incremental tasks with temporal heterogeneity in FL scenarios.

**Incremental Learning** Incremental Learning (IL) is a machine learning technique that allows a model to learn continuously from an incremental sequence of tasks while retaining knowledge gained from previous tasks [12, 43], including task-incremental learning [5, 30], class-incremental learning [39, 49], and domain-incremental learning [2, 33]. Existing approaches in IL can be classified into three main categories: replay-based methods [27, 39], regularization-based methods [16, 47], and parameter isolation methods [8, 28]. Replay-based methods select representative old samples to retain previously learned knowledge when training on a new task. Regularization-based methods protect existing knowledge from being overwritten with new knowledge by imposing constraints on the loss function of new tasks. Parameter isolation methods typically introduce additional parameters and computations to learn new tasks. Here we focus on the federated incremental learning scenario, which can be viewed as a combination of federal learning and incremental learning.

### 3. Methodology

We first formulate two FIL scenarios and propose a simple and scalable framework Re-Fed. Then, we present a scalable algorithm and provide rigorous analytical results to show the efficiency of the proposed method.

#### 3.1. Problem Formulation

In the standard IL (non-federated environment), a model learns from a sequence of streaming tasks  $\{\mathcal{T}^1, \mathcal{T}^2, \dots, \mathcal{T}^n\}$  where  $\mathcal{T}^t$  denotes the  $t$ -th task of the

dataset. Here  $\mathcal{T}^t = \sum_{i=1}^{N^t} (x_t^{(i)}, y_t^{(i)})$ , which has  $N^t$  pairs of sample data  $x_t^{(i)} \in \mathcal{X}^t$  and corresponding label  $y_t^{(i)} \in \mathcal{Y}^t$ . We use  $\mathcal{X}^t$  and  $\mathcal{Y}^t$  to represent the domain space and label space for the  $t$ -th task, which has  $|\mathcal{Y}^t|$  classes and  $\mathcal{Y} = \bigcup_{t=1}^n \mathcal{Y}^t$  where  $\mathcal{Y}$  denotes the total classes of all time. Similarly, we use  $\mathcal{X} = \bigcup_{t=1}^n \mathcal{X}^t$  to denote the total domain space for tasks of all time. In this paper, we focus on two types of IL scenarios: (1) Class-Incremental Task: all tasks share the same domain space, i.e.,  $\mathcal{X}^1 = \mathcal{X}^t, \forall t \in [n]$ . As the sequence of learning tasks arrives, the number of the classes may change, i.e.,  $\mathcal{Y}^1 \neq \mathcal{Y}^t, \forall t \in [n]$ . (2) Domain-Incremental Task: all tasks share the same number of classes i.e.,  $\mathcal{Y}^1 = \mathcal{Y}^t, \forall t \in [n]$ . As the sequence of tasks arrives, the client needs to learn the new task while their domain and data distribution changes, i.e.,  $\mathcal{X}^1 \neq \mathcal{X}^t, \forall t \in [n]$ .

We further consider IL in a federated setting. We aim to train a global model for  $K$  total clients and assume that client  $k$  can only access the local private streaming tasks  $\{\mathcal{T}_k^1, \mathcal{T}_k^2, \dots, \mathcal{T}_k^n\}$ . When the  $t$ -th task comes, while clients can cache all samples from previous tasks without forgetting, the goal is to train a global model  $w^t$  over all  $t$  tasks  $\mathcal{T}^t = \{\sum_{n=1}^t \sum_{k=1}^K \mathcal{T}_k^n\}$ , which can be formulated as :

$$w^t = \arg \min_w \sum_{n=1}^t \sum_{k=1}^K \sum_{i=1}^{N_k^n} \frac{1}{|\mathcal{T}^t|} l \left( f_{w_k}(x_{k,n}^{(i)}), y_{k,n}^{(i)} \right). \quad (1)$$

where  $f_{w_k}(\cdot)$  is the output of the model  $w_k$  in client  $k$  and  $l(\cdot)$  is the cross-entropy loss. Then, due to poor storage in common edge devices, each client caches partial samples for replay. Here we assume each client can only store total  $M$  samples and has to cache  $M - N_k^t$  samples from  $(t - 1)$  previous tasks, which is denoted as  $\mathcal{T}_{k,cached}^{t-1} = \sum_{i=1}^{M-N_k^t} (\tilde{x}_{k,t-1}^{(i)}, \tilde{y}_{k,t-1}^{(i)})$ . The goal is to train a global model  $w^t$  over both cached samples and the  $t$ -th new task, which can be formulated as:

$$w^t = \arg \min_w \sum_{k=1}^K \sum_{i=1}^M \frac{1}{|\mathcal{T}_{k,local}^t|} l \left( f_{w_k}(\tilde{x}_{k,t}^{(i)}), \tilde{y}_{k,t}^{(i)} \right). \quad (2)$$

where  $\mathcal{T}_{k,local}^t = \mathcal{T}_{k,cached}^{t-1} + \mathcal{T}_k^t = \sum_{i=1}^M (\tilde{x}_{k,t}^{(i)}, \tilde{y}_{k,t}^{(i)})$ .

#### 3.2. Re-Fed: Framework for FIL

The key idea of Re-Fed is to identify the sample importance and coordinate clients to cache important previous samples with limited local storage when the new task arrives. Specifically, in each communication round, the clients train the local models with the private sequence of tasks and the server aggregates local models from all participated clients. Then, when new tasks come, each client first trains an extra personalized informative model on previous local samples with the regulation of both the global model and local model.

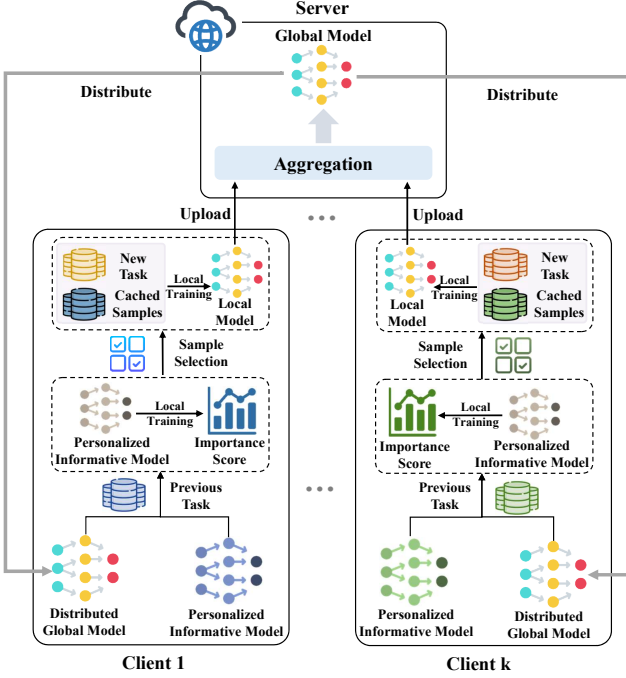


Figure 2. Illustration of the Re-Fed framework. When a new task arrives, each client first updates the personalized informative model on previous local samples with the distributed global model. Then, samples are selected to be cached by the sample importance scores that are calculated with the updated personalized informative model. Finally, each client trains the local model with both the new task and cached previous samples.

During the update of such model, gradient norms of individual samples are recorded to calculate sample importance scores. Finally, each client caches samples with higher importance scores based on local storage and continues training the local model on the cached samples and the new task. The workflow of the proposed framework is shown in Algorithm 1 and Figure 2 illustrates the Re-Fed framework.

**Personalized Informative Model.** In FIL with Non-IID client samples, sample importance should be defined based not only on their importance in the local dataset but also on their correlation to the global dataset across clients such that the model can be better trained with the cached samples. In a standard FL scenario, each client can only access its own local model and global model, which respectively contains local and global information. A straightforward idea here is to calculate two sample importance scores using local and global models and cache samples based on such two scores for data reply. Then, one can build upon such an idea by adding following capabilities: (1) The global model is aggregated by local models from participated clients and the gradient norm of the global model can be calculated locally without training the global model. (2) A control mechanism should be available to adjust the proportion of local

---

### Algorithm 1: Re-Fed

---

**Input:**  $T$ : communication round;  $K$ : client number;  $\eta$ : learning rate;  $\{\mathcal{T}^t\}_{t=1}^n$ : distributed dataset with  $n$  tasks;  $w$ : parameter of the model;  $v_k$ : personalized informative model in client  $k$ .

- 1 Initialize the parameter  $w$ ;
- 2 **for**  $c = 1$  **to**  $T$  **do** // the  $t$ -th new task
- 3     Server randomly selects a subset of devices  $S_t$  and send  $w^{t-1}$
- 4     **for each selected client**  $k \in S_t$  **in parallel do**
- 5         Update  $v_k^{t-1}$  in  $s$  local iterations with (3).
- 6         **for** *During the update of*  $v_k^{t-1}$  **do**
- 7             Calculate the importance score after total  $s$  iterations for the sample  $(\tilde{x}_{k,t-1}^{(i)}, \tilde{y}_{k,t-1}^{(i)})$  with (5).
- 8         **end**
- 9         Cache previous samples with higher importance scores;
- 10         Training the local model with cached samples and the new task with (6);
- 11         Send the model  $w_k^t$  back to the server.
- 12     **end**
- 13      $w^t \leftarrow \text{ServerAggregation}(\{w_k^t\}_{k \in S_t})$
- 14 **end**

---

and global information in the sample importance.

Toward the above goals, here we introduce an additional personalized informative model (PIM) for each client, which digests both the information from the local model and the global model. Then, we propose to adopt a ratio factor to adjust the information proportion from local and global sides. Finally, we can record the gradient norms of the samples to calculate the importance scores during the update of PIM, resulting in sample importance scores with both local and global information. Suppose that the client  $k$  receives the global model  $w^{t-1}$  and then the  $t$ -th new task arrives, and the clients update PIM  $v_k^{t-1}$  with previous local samples  $\mathcal{T}_{k,local}^{t-1}$  in  $s$  iterations as follows:

$$v_{k,s}^{t-1} = v_{k,s-1}^{t-1} - \eta \left( \sum_{i=1}^M \nabla l \left( f_{v_{k,s-1}^{t-1}}(\tilde{x}_{k,t-1}^{(i)}, \tilde{y}_{k,t-1}^{(i)}) + q(\lambda)(v_{k,s-1}^{t-1} - w^{t-1}) \right) \right). \quad (3)$$

where  $q(\lambda) = \frac{1-\lambda}{2\lambda}$ ,  $\lambda \in (0, 1)$ , and  $\eta$  is the rate to control the step size of the update. The hyper-parameter  $\lambda$  adjusts the balance between the local and global information incorporated in the update.

To better understand the update, we can draw an analogy to momentum methods in optimization. Momentum-based

Table 1. Performance comparison of various methods in two incremental learning scenarios.

Scenario	Dataset	FedAvg	FedProx	Fixed	DANN+FL	Shared	FCIL	FedCIL	Re-Fed
Class-Incremental	CIFAR10 ( $\alpha = 1.0$ )	26.73 $\pm$ 1.12	25.87 $\pm$ 0.68	19.21 $\pm$ 0.06	24.86 $\pm$ 2.31	23.91 $\pm$ 1.70	25.04 $\pm$ 0.11	27.35 $\pm$ 1.24	<b>29.22<math>\pm</math>0.49</b>
	CIFAR100 ( $\alpha = 5.0$ )	17.21 $\pm$ 1.35	18.03 $\pm$ 0.91	9.27 $\pm$ 0.22	19.73 $\pm$ 2.17	18.30 $\pm$ 1.53	23.02 $\pm$ 0.66	17.98 $\pm$ 1.46	<b>25.61<math>\pm</math>0.88</b>
	Tiny-ImageNet ( $\alpha = 10$ )	27.58 $\pm$ 0.74	21.82 $\pm$ 0.90	12.34 $\pm$ 0.23	20.77 $\pm$ 1.31	22.19 $\pm$ 0.54	29.58 $\pm$ 0.15	24.41 $\pm$ 0.95	<b>32.07<math>\pm</math>0.27</b>
Domain-Incremental	Digit10 ( $\alpha = 0.1$ )	77.59 $\pm$ 0.39	79.09 $\pm$ 0.58	71.26 $\pm$ 0.04	76.44 $\pm$ 1.05	74.77 $\pm$ 0.23	77.59 $\pm$ 0.39	83.85 $\pm$ 0.80	<b>85.96<math>\pm</math>0.14</b>
	Office31 ( $\alpha = 1$ )	39.25 $\pm$ 1.61	43.01 $\pm$ 1.59	37.44 $\pm$ 0.72	45.21 $\pm$ 2.10	37.55 $\pm$ 0.69	39.25 $\pm$ 1.61	46.26 $\pm$ 2.24	<b>50.80<math>\pm</math>0.77</b>
	DomainNet ( $\alpha = 10$ )	51.73 $\pm$ 2.32	49.12 $\pm$ 2.71	46.30 $\pm$ 1.42	50.01 $\pm$ 3.31	41.76 $\pm$ 1.26	51.73 $\pm$ 2.32	47.28 $\pm$ 3.01	<b>56.66<math>\pm</math>0.50</b>

methods leverage the past updates to guide the current update direction [1]. Similarly, the term  $q(\lambda)(v_{k,s-1}^{t-1} - w^{t-1})$  acts as a momentum component. It incorporates information from the global model  $w^{t-1}$  to influence the update of PIM  $v_k^{t-1}$ . The hyper-parameter  $\lambda$  controls the weight of this momentum component, and it lies within the range of (0,1). When  $\lambda$  is close to 0, PIM primarily focuses on recovering the global model  $w^{t-1}$ . In other words, it will align itself with the global data. On the other hand, as  $\lambda$  becomes larger, it leads to a stronger emphasis on local training.

**Theorem 3.1** (Convergence of PIM). *Assuming that the global model  $w^t$  converges to the optimal model  $\hat{w}$  at communication round  $t$  by  $g(t)$  as:  $\mathbb{E}[\|w^t - \hat{w}\|^2] \leq g(t)$ ,  $\lim_{t \rightarrow \infty} g(t) = 0$  and  $g(t+1) \leq g(t)$ , there exists a constant  $C < \infty$  such that for any client  $k \in [K]$  the personalized informative model  $v_k^t$  can converge to the optimal model  $\hat{v}_k$  by  $Cg(t)$ .*

With Theorem 3.1, we ensure the convergence of PIM and thus can calculate the gradient norms of samples during the training stage. We provide the proof in Appendix F.1.

**Sample Importance.** To quantify the sample importance, we investigate the impact of samples on the generalization ability of the model, and allocate the samples that can enhance the generalization with higher importance. We calculate the gradient norm with respect to model parameters of PIM as the importance scores to samples. The gradient norm of samples during training is recorded, which can be regarded as the contribution of the sample to the model update. A similar flavor of such a method can be found in [36, 42]. Suppose that the client  $k$  has converged on  $(t-1)$ -th tasks with local samples  $\mathcal{T}_{k,local}^{t-1}$ , when it comes to the incremental  $t$ -th task, the client should calculate the importance scores for all local samples. Here we denote  $G^p(\tilde{x}_{k,t-1}^{(i)})$  is the gradient norm of the sample  $(\tilde{x}_{k,t-1}^{(i)}, \tilde{y}_{k,t-1}^{(i)})$  in  $p$ -th iteration during the update of PIM  $v_{k,p}^{t-1}$ , which is:

$$G^p(\tilde{x}_{k,t-1}^{(i)}) = \left\| \nabla l \left( f_{v_{k,p}^{t-1}}(\tilde{x}_{k,t-1}^{(i)}, \tilde{y}_{k,t-1}^{(i)}) \right) \right\|^2. \quad (4)$$

According to [42], the difference in the gradient of the loss function with and without a sample  $(\tilde{x}_{k,t-1}^{(j)}, \tilde{y}_{k,t-1}^{(j)})$  is up-

per bounded and the bound is linearly dependent on the sample gradient norm defined in Eq. 4. Thus, caching samples based on sample gradient norms can least affect the gradient and best preserve the dynamics of training.

As PIM integrates both local and global models, a greater gradient norm of a sample with PIM indicates that such a sample drives PIM more to fitting the task with local and global knowledge. Such an effect could be more prominent at early training during  $s$  iterations for PIM, where fluctuation around optima is rare than that later in the training. Thus, we accumulate the gradient norm during the training of PIM and emphasis on the early training stage to calculate the sample importance as:

$$I(\tilde{x}_{k,t-1}^{(i)}) = \sum_{p=1}^s \frac{1}{p} G^p(\tilde{x}_{k,t-1}^{(i)}). \quad (5)$$

We also provide illustrative experimental results with  $I(\tilde{x}_{k,t-1}^{(i)}) = \sum_{p=1}^s G^p(\tilde{x}_{k,t-1}^{(i)})$  in the Appendix E and show the performance improvement with the sample importance adopted in Eq. 5.

**Local Training.** After caching important samples with higher importance scores, each client continues to train the local model  $w_k^t$  with local samples  $\mathcal{T}_{k,local}^t$  in iteration  $p \in [1, s]$  as follows:

$$w_{k,p}^t = w_{k,p-1}^t - \eta \sum_{i=1}^M \nabla l \left( f_{w_{k,p-1}^t}(\tilde{x}_{k,t}^{(i)}, \tilde{y}_{k,t}^{(i)}) \right). \quad (6)$$

**Modularity of Re-Fed.** From the Re-Fed framework and Algorithm 1, we can see that a key feature of Re-Fed is its unique modularity. One can readily use prior art developed for FIL algorithm, and employ Re-Fed as a useful off-the-shelf add-on. Our method has several advantages:

- **Optimization:** It is possible to plug in other aggregation methods beyond FedAvg [31] in Algorithm 1 to update the global model, and inherit the convergence benefits. In the subsequent experimental design, we investigate the performance of our Re-Fed framework with FedAvg algorithm and provide a detailed algorithm definition using FedAvg in Appendix D.
- **Privacy:** Re-Fed transmits no more extra information over the network than typical FL algorithms. This is dif-

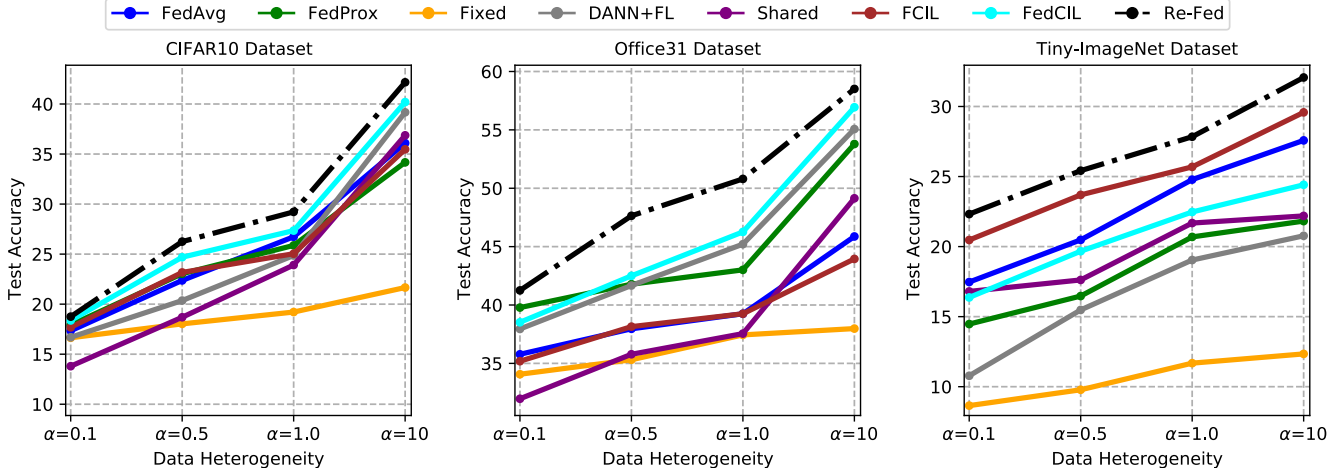


Figure 3. Performance w.r.t data heterogeneity  $\alpha$  for three datasets.

Table 2. Test accuracy for Re-Fed w.r.t data heterogeneity  $\alpha$  and hyper-parameter  $\lambda$  on CIFAR10, Office31 and Tiny-ImageNet.

Dataset	$\alpha = 0.1$			$\alpha = 0.5$			$\alpha = 1.0$			$\alpha = 10$		
	$\lambda = 0.2$	$\lambda = 0.5$	$\lambda = 0.8$	$\lambda = 0.2$	$\lambda = 0.5$	$\lambda = 0.8$	$\lambda = 0.2$	$\lambda = 0.5$	$\lambda = 0.8$	$\lambda = 0.2$	$\lambda = 0.5$	$\lambda = 0.8$
CIFAR10	<b>18.75</b> $\pm$ 1.30	18.61 $\pm$ 1.09	17.91 $\pm$ 0.81	<b>26.25</b> $\pm$ 1.64	26.00 $\pm$ 0.97	25.62 $\pm$ 0.52	27.05 $\pm$ 0.88	27.80 $\pm$ 0.21	<b>29.22</b> $\pm$ 0.49	38.43 $\pm$ 0.43	40.04 $\pm$ 0.19	<b>42.17</b> $\pm$ 0.25
Office31	<b>41.25</b> $\pm$ 1.01	39.29 $\pm$ 1.34	38.18 $\pm$ 0.68	46.86 $\pm$ 0.91	<b>47.64</b> $\pm$ 0.53	47.13 $\pm$ 1.16	43.81 $\pm$ 0.73	48.67 $\pm$ 0.99	<b>50.08</b> $\pm$ 0.77	52.79 $\pm$ 1.28	55.92 $\pm$ 0.38	<b>58.51</b> $\pm$ 0.46
Tiny-ImageNet	<b>22.32</b> $\pm$ 0.12	20.51 $\pm$ 0.98	18.00 $\pm$ 1.30	24.60 $\pm$ 0.48	<b>25.42</b> $\pm$ 0.59	24.39 $\pm$ 0.66	24.88 $\pm$ 0.87	27.15 $\pm$ 0.78	<b>27.84</b> $\pm$ 0.73	29.03 $\pm$ 0.30	30.26 $\pm$ 0.24	<b>32.07</b> $\pm$ 0.27

ferent from most other FIL methods where sample reconstruction methods are applied for data replay, which may raise privacy concerns.

- **Resource:** Re-Fed allows each client to train a backbone model using only its local training data, without employing additional distillation data or generated data, leading to extra either computation cost or storage overhead.

## 4. Experiments

In this section, we evaluate our proposed framework concerning two incremental learning scenarios. We investigate the relationship between the data heterogeneity and the balance between the local and global information incorporated in PIM. Additionally, we conduct parameter sensitivity analysis to verify the effectiveness of our method.

### 4.1. Experiment Setup

**Dataset:** We conduct our experiments with heterogeneously partitioned datasets over two federated incremental scenarios on six datasets. (1) **Class-Incremental Learning:** CIFAR10[17], CIFAR100[17] and Tiny-ImageNet[18]. (2) **Domain-Incremental Learning:** Digit10, Office31[41] and DomainNet[37]. Among them, the Digit10 dataset contains 10 digit categories in four domains: MNIST[19], EMNIST[3], USPS[13] and SVHN[34]. Details of datasets and data processing can be found in Appendix A.

**Baseline:** For a fair comparison with other key works, we follow the same protocols proposed by [31, 39] to set up FIL tasks. We evaluate all the methods with two representative FL models **FedAvg** [31] and **Fedprox** [22], two models designed for federated class-incremental learning: **FCIL** [6] and **FedCIL** [38], and three customized methods of **Fixed**: we train the model only from the first task and evaluate it for all the coming sequence of tasks; **DANN+FL**: here we adopt the robust adversarial-based method DANN[9] in local training for domain adaption; **Shared**: we adopt all front layers before the last fully connected layer as shared layers, and use relevant different fully-connected layers to obtain outputs for different tasks. Details of datasets and data processing can be found in Appendix B.

**Configurations:** Unless otherwise mentioned, we set the number of local training epoch  $E = 20$  and communication round  $T = 150$  for each task, which ensures the convergence of previous tasks before the arrival of new task. We use the Dirichlet distribution  $\text{Dir}(\alpha)$  to distribute local samples to yield data heterogeneity for all tasks where a smaller  $\alpha$  indicates higher data heterogeneity. We employ ResNet18 [11] as the basic backbone model in all methods. We calculate the Top-10 accuracy for the Tiny-ImageNet and DomainNet datasets and Top-1 accuracy for others. Each experiment setting is run three times and we record accuracy in the final 10 rounds and report the average value and standard deviation. The total clients number is 20/10 with an active ratio  $k = 0.4$  for {CIFAR10, CIFAR100,

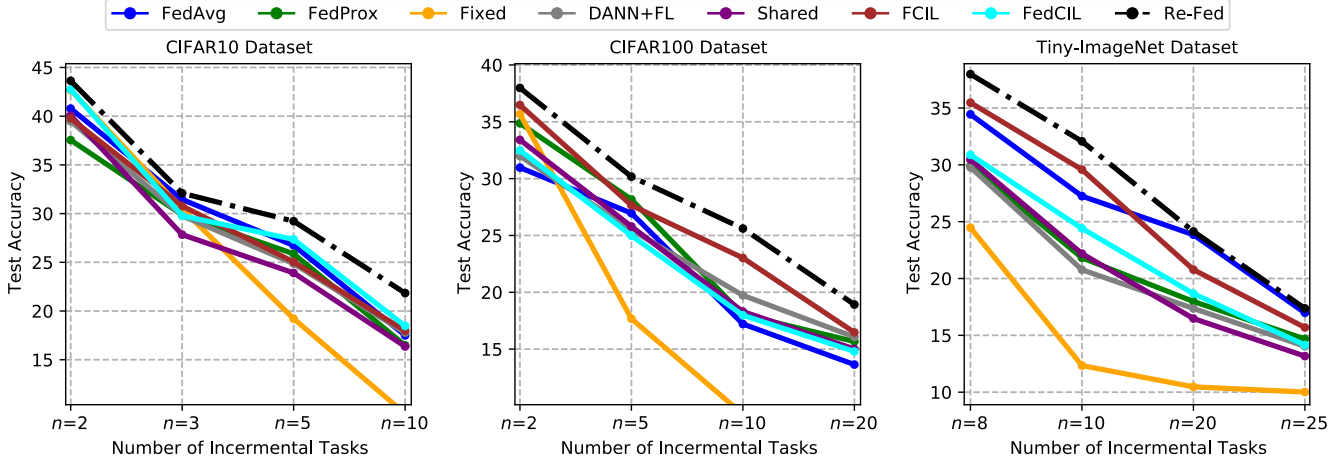


Figure 4. Performance w.r.t number of incremental tasks  $n$  for three class-incremental datasets.

Table 3. Evaluation of various methods, in terms of the communication rounds to reach the best test accuracy (150 communication rounds each task). We report the sum of communication rounds required to achieve the best performance on each task.

Scenario	Dataset	FedAvg	FedProx	Fixed	DANN+FL	Shared	FCIL	FedCIL	Re-Fed
Class-Incremental	CIFAR10 (Task:5)	613 $\pm$ 2.67	685 $\pm$ 3.00	142 $\pm$ 0.67	712 $\pm$ 3.67	574 $\pm$ 1.33	590 $\pm$ 2.67	738 $\pm$ 4.00	<b>562<math>\pm</math>1.67</b>
	CIFAR100 (Task:10)	1103 $\pm$ 2.33	1246 $\pm$ 3.00	137 $\pm$ 2.00	1258 $\pm$ 4.67	1154 $\pm$ 3.33	1095 $\pm$ 2.67	1311 $\pm$ 5.67	<b>1039<math>\pm</math>4.33</b>
	Tiny-ImageNet (Task:10)	1197 $\pm$ 2.67	1234 $\pm$ 2.67	132 $\pm$ 3.00	1305 $\pm$ 3.67	1278 $\pm$ 4.33	1185 $\pm$ 2.33	1317 $\pm$ 3.33	<b>1128<math>\pm</math>3.67</b>
Domain-Incremental	Digit10 (Task:4)	410 $\pm$ 1.67	412 $\pm$ 0.67	112 $\pm$ 0.33	483 $\pm$ 1.33	372 $\pm$ 2.00	410 $\pm$ 1.67	419 $\pm$ 2.67	<b>325<math>\pm</math>1.33</b>
	Office31 (Task:3)	413 $\pm$ 2.67	429 $\pm$ 2.00	144 $\pm$ 0.67	436 $\pm$ 3.67	391 $\pm$ 1.12	413 $\pm$ 2.67	431 $\pm$ 3.33	<b>388<math>\pm</math>1.67</b>
	DomainNet (Task:6)	726 $\pm$ 3.33	767 $\pm$ 2.67	141 $\pm$ 1.67	752 $\pm$ 4.00	694 $\pm$ 2.67	726 $\pm$ 3.33	791 $\pm$ 3.67	<b>661<math>\pm</math>2.33</b>

Tiny-ImageNet, Digit10, DomainNet}/\{Office31\}. The maximum number of cached samples  $M$  is 2000/1000/300 for \{Digit10, DomainNet, Tiny-ImageNet\}/\{CIFAR10, CIFAR100\}/\{Office31\}. We experiment with different numbers of incremental tasks and each task arrives with new classes: 10/10/5 tasks with 20/10/2 new classes in each task for \{Tiny-ImageNet\}/\{CIFAR100\}/\{CIFAR10\}. Details of configurations can be found in Appendix C.

## 4.2. Performance Overview

**Test Accuracy.** Table 1 shows the test accuracy of various methods with data heterogeneity across six datasets. We report the final accuracy of the global model when all clients finish their training on all tasks. FedCIL outperforms FedAvg and FedProx on CIFAR10, Digit10 and Office31 as the generator in FedCIL exhibits effective training on simple datasets, enabling the generation of high-quality samples for data replay. However, its performance experiences a significant decline with larger-scale datasets such as CIFAR100 and DomainNet, where the classes and domains become more complex. In the federated domain-incremental learning scenario, FCIL reverts to the FedAvg algorithm when there are no new incremental sample classes. Re-Fed achieves the best performance in all

cases by a margin of 1.87%~19.73% in terms of final accuracy. More discussions and results on model performance are available in Appendix E.

**Data Heterogeneity.** Figure 3 displays the test accuracy with different levels of data heterogeneity on three datasets. As shown in this figure, all methods achieve an improvement in test accuracy with the decline in data heterogeneity, and Re-Fed consistently achieves a leading improvement in performance with different levels of data heterogeneity.

Then, we conduct more research on the setting of hyperparameter  $\lambda$ . In our framework, we modify the  $\lambda$  value to adjust the global and local information proportion in PIM. As shown in Table 2, we select three different  $\lambda$  values with four different data heterogeneity settings and evaluate the final test accuracy on three datasets. Experimental results show that the value of  $\lambda$  should be chosen accordingly under different data heterogeneity. Nevertheless, results exhibit the same trend: as the degree of data heterogeneity increases, our Re-Fed framework performs better while  $\lambda$  decreases as PIM contains more global information. Empirically, striking a balance between global and local information is the key to addressing the data heterogeneity in FIL. Vice versa, as  $\alpha$  increases, the data distribution on the client side becomes more IID. At this point, the clients

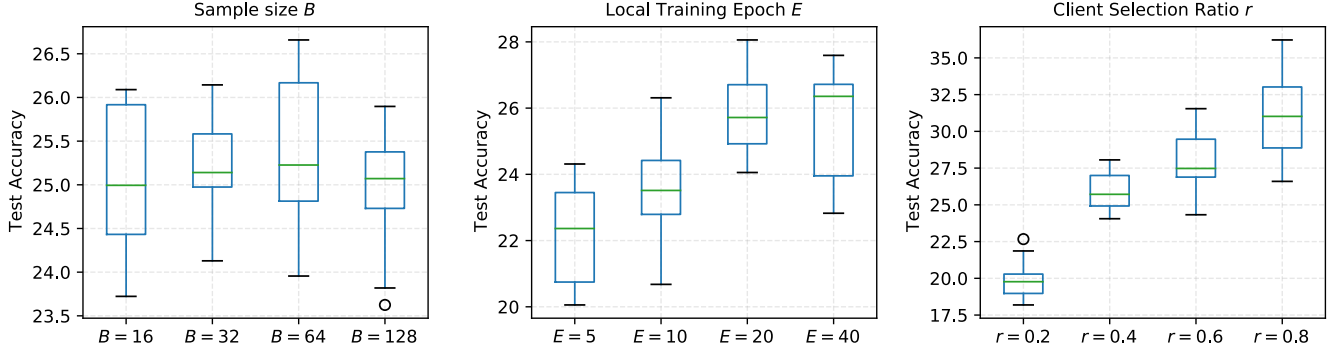


Figure 5. Performance of Re-Fed under different configurations (a) local training epoch  $E$ , (b) sample size  $B$  in the classifier, (c) client selection ratio  $r$  of all clients on CIFAR100 with  $\alpha = 5.0$ .

Table 4. Test accuracy for Re-Fed and other baselines w.r.t memory size  $M$  on Digit10, CIFAR10, CIFAR100 and Tiny-ImageNet.

Dataset	FedAvg			DANN+FL			FCIL			Re-Fed		
	$M = 1000$	$M = 2000$	$M = 3000$	$M = 1000$	$M = 2000$	$M = 3000$	$M = 1000$	$M = 2000$	$M = 3000$	$M = 1000$	$M = 2000$	$M = 3000$
Digit10	77.05 $\pm$ 0.24	77.59 $\pm$ 0.39	81.64 $\pm$ 0.31	75.16 $\pm$ 0.86	76.44 $\pm$ 1.05	83.34 $\pm$ 0.97	77.50 $\pm$ 0.24	77.59 $\pm$ 0.39	81.64 $\pm$ 0.31	<b>82.30<math>\pm</math>0.16</b>	<b>85.96<math>\pm</math>0.14</b>	<b>86.32<math>\pm</math>0.04</b>
CIFAR10	26.73 $\pm$ 1.12	30.19 $\pm$ 1.63	33.41 $\pm$ 0.92	24.86 $\pm$ 2.31	27.65 $\pm$ 1.34	30.09 $\pm$ 1.06	25.04 $\pm$ 0.11	29.14 $\pm$ 0.19	31.77 $\pm$ 0.28	<b>29.22<math>\pm</math>0.49</b>	<b>32.83<math>\pm</math>0.20</b>	<b>33.41<math>\pm</math>0.74</b>
CIFAR100	17.21 $\pm$ 1.35	23.58 $\pm$ 1.82	29.90 $\pm$ 2.11	19.73 $\pm$ 2.17	25.56 $\pm$ 2.68	28.49 $\pm$ 1.98	23.02 $\pm$ 0.66	26.12 $\pm$ 0.54	29.06 $\pm$ 0.89	<b>25.61<math>\pm</math>0.88</b>	<b>28.41<math>\pm</math>0.24</b>	<b>29.90<math>\pm</math>0.71</b>
Tiny-ImageNet	24.42 $\pm$ 0.59	27.58 $\pm$ 0.74	31.50 $\pm$ 0.63	18.06 $\pm$ 1.08	20.77 $\pm$ 1.31	26.93 $\pm$ 0.77	25.20 $\pm$ 1.12	29.58 $\pm$ 0.15	33.44 $\pm$ 0.38	<b>28.31<math>\pm</math>0.19</b>	<b>32.07<math>\pm</math>0.27</b>	<b>35.66<math>\pm</math>0.42</b>

require less global information and can rely more on their local information for caching important samples.

**Quantitative Analysis.** Figure 4 shows the qualitative analysis of the number of incremental tasks  $n$  on three class-incremental datasets. According to these curves, we can easily observe that our model performs better than other baselines across all tasks, with varying numbers of incremental tasks. It demonstrates that Re-Fed enables clients to learn new incremental classes better than other methods.

**Communication Efficiency.** Table 3 shows the evaluation of various methods in terms of the communication rounds to reach the test accuracy reported in Table 1. Here we show the sum of communication rounds required to achieve the performance on each task. Re-Fed requires the least communication rounds to achieve the reported test accuracy on all datasets with its caching method. Thus, the local model is easier to converge. More details of the communication round on each task are available in Appendix E.

**Parameter Sensitivity Analysis.** Figure 5 shows the performance of Re-Fed under different configurations with standard boxplots from ten trails with different random seeds. Re-Fed achieves similar performance with different sample sizes  $B$ , and it achieves a better result when we increase the local training epochs. However, Re-Fed has a comparable performance when the  $E$  is set to 20 and 40. In addition, a larger client selection ratio  $r$  contributes to higher test accuracy.

For the FIL scenario, we conduct additional research on the client storage  $M$ . When  $M$  is large enough for clients to cache samples from all tasks, our framework degrades to a normal FedAvg algorithm with a naive caching method.

As shown in Table 4, we select three different  $M$  values to conduct experiments. Experimental results demonstrate that larger memory size  $M$  contributes to the training and Re-Fed outperforms other baselines in all cases. To conclude, Re-Fed is only sensitive to few parameters and still robust to most parameters in a large range.

## 5. Conclusion and Future Work

We proposed Re-Fed, a simple framework, to address the catastrophic forgetting with data heterogeneity in federated incremental learning. Re-Fed can be thought of as a lightweight personalization add-on for any federated learning algorithms with global aggregation, which maintains privacy and communication efficiency. Extensive experiments conducted on various settings and baselines show that Re-Fed achieves significant improvement in test accuracy.

Although existing works and our method have demonstrated great effectiveness over the FIL scenarios, none have studied the dynamic requirements of the edge clients. To deploy the FL system in practical settings, it is necessary to consider personalized local factors such as storage, computation and even the different arrival time of the new task. In the future, we seek to work a step forward in this field.

## Acknowledgements

This work is supported by National Natural Science Foundation of China under grants 62376103, 62302184, 62206102, Science and Technology Support Program of Hubei Province under grant 2022BAA046, and CCF-AFSG Research Fund.



## References

- [1] Louis KC Chan, Narasimhan Jegadeesh, and Josef Lakonishok. Momentum strategies. *The journal of Finance*, 51(5): 1681–1713, 1996. 5
- [2] Nikhil Churamani, Ozgur Kara, and Hatice Gunes. Domain-incremental continual learning for mitigating bias in facial expression and action unit recognition. *ArXiv*, abs/2103.08637, 2021. 3
- [3] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*, pages 2921–2926. IEEE, 2017. 6, 1
- [4] Marcos F Criado, Fernando E Casado, Roberto Iglesias, Carlos V Regueiro, and Senén Barro. Non-iid data and continual learning processes in federated learning: A long road ahead. *Information Fusion*, 88:263–280, 2022. 1
- [5] Neil T. Dantam, Zachary K. Kingston, Swarat Chaudhuri, and Lydia E. Kavraki. Incremental task and motion planning: A constraint-based approach. In *Robotics: Science and Systems*, 2016. 3
- [6] Jiahua Dong, Lixu Wang, Zhen Fang, Gan Sun, Shichao Xu, Xiao Wang, and Qi Zhu. Federated class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10164–10173, 2022. 1, 6
- [7] Zhaoyang Du, Celimuge Wu, Tsutomu Yoshinaga, Kok-Lim Alvin Yau, Yusheng Ji, and Jie Li. Federated learning for vehicular internet of things: Recent advances and open issues. *IEEE Open Journal of the Computer Society*, 1:45–61, 2020. 1
- [8] Chrisantha Fernando, Dylan S. Banarse, Charles Blundell, Yori Zwols, David R Ha, Andrei A. Rusu, Alexander Pritzel, and Daan Wierstra. Pathnet: Evolution channels gradient descent in super neural networks. *ArXiv*, abs/1701.08734, 2017. 3
- [9] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016. 1, 6
- [10] Filip Hanzely and Peter Richtárik. Federated learning of a mixture of global and local models. *CoRR*, abs/2002.05516, 2020. 6
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 6
- [12] Yen-Chang Hsu, Yen-Cheng Liu, and Zsolt Kira. Re-evaluating continual learning scenarios: A categorization and case for strong baselines. *ArXiv*, abs/1810.12488, 2018. 3
- [13] Jonathan J. Hull. A database for handwritten text recognition research. *IEEE Transactions on pattern analysis and machine intelligence*, 16(5):550–554, 1994. 6, 1
- [14] Eunjeong Jeong, Seungeun Oh, Hyesung Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun Kim. Communication-efficient on-device machine learning: Federated distillation and augmentation under non-iid private data. *CoRR*, abs/1811.11479, 2018. 3
- [15] Eunjeong Jeong, Seungeun Oh, Hyesung Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun Kim. Communication-efficient on-device machine learning: Federated distillation and augmentation under non-iid private data. *ArXiv*, abs/1811.11479, 2018. 3
- [16] Sangwon Jung, Hongjoon Ahn, Sungmin Cha, and Taesup Moon. Continual learning with node-importance based adaptive group sparse regularization. *arXiv: Learning*, 2020. 3
- [17] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 6
- [18] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015. 6
- [19] Yann LeCun, Corinna Cortes, and Chris Burges. Mnist handwritten digit database, 2010. 6, 1
- [20] Yuan Lei, Shir Li Wang, Minghui Zhong, Meixia Wang, and Theam Foo Ng. A federated learning framework based on incremental weighting and diversity selection for internet of vehicles. *Electronics*, 11(22):3668, 2022. 1
- [21] Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Federated multi-task learning for competing constraints. *CoRR*, abs/2012.04221, 2020. 6
- [22] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems*, 2:429–450, 2020. 3, 6
- [23] Xin-Chun Li, Yi-Chu Xu, Shaoming Song, Bingshuai Li, Yinchuan Li, Yunfeng Shao, and De-Chuan Zhan. Federated learning with position-aware neurons. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 10072–10081. 1, 2
- [24] Yijing Li, Xiaofeng Tao, Xuefei Zhang, Junjie Liu, and Jin Xu. Privacy-preserved federated learning for autonomous driving. *IEEE Transactions on Intelligent Transportation Systems*, 23(7):8423–8434, 2021. 1
- [25] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. *Advances in Neural Information Processing Systems*, 33:2351–2363, 2020. 3
- [26] Lumin Liu, Jun Zhang, S. H. Song, and Khaled Ben Letaief. Edge-assisted hierarchical federated learning with non-iid data. *ArXiv*, abs/1905.06641, 2019. 3
- [27] Xialei Liu, Chenshen Wu, Mikel Menta, Luis Herranz, Bogdan Raducanu, Andrew D Bagdanov, Shangling Jui, and Joost van de Weijer. Generative feature replay for class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 226–227, 2020. 3
- [28] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I. Jordan. Learning transferable features with deep adaptation networks. *ArXiv*, abs/1502.02791, 2015. 3
- [29] Yuhang Ma, Zhongle Xie, Jue Wang, Ke Chen, and Lidan Shou. Continual federated learning based on knowledge distillation. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages

- 2182–2188. International Joint Conferences on Artificial Intelligence Organization, 2022. Main Track. [1](#)
- [30] Davide Maltoni and Vincenzo Lomonaco. Continuous learning in single-incremental-task scenarios. *Neural networks : the official journal of the International Neural Network Society*, 116:56–73, 2018. [3](#)
- [31] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017. [1](#), [2](#), [3](#), [5](#), [6](#)
- [32] Thomas Mensink, Jakob Verbeek, Florent Perronnin, and Gabriela Csurka. Distance-based image classification: Generalizing to new classes at near-zero cost. *IEEE transactions on pattern analysis and machine intelligence*, 35(11):2624–2637, 2013. [2](#)
- [33] M Jehanzeb Mirza, Marc Masana, Horst Possegger, and Horst Bischof. An efficient domain-incremental learning approach to drive in all weather conditions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3001–3011, 2022. [3](#)
- [34] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011. [6](#), [1](#)
- [35] Dinh C Nguyen, Quoc-Viet Pham, Pubudu N Pathirana, Ming Ding, Aruna Seneviratne, Zihuai Lin, Octavia Dobre, and Won-Joo Hwang. Federated learning for smart healthcare: A survey. *ACM Computing Surveys (CSUR)*, 55(3): 1–37, 2022. [1](#)
- [36] Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding important examples early in training. *CoRR*, abs/2107.07075, 2021. [5](#)
- [37] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1406–1415, 2019. [6](#)
- [38] Daiqing Qi, Handong Zhao, and Sheng Li. Better generative replay for continual federated learning. *arXiv preprint arXiv:2302.13001*, 2023. [1](#), [6](#)
- [39] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010, 2017. [2](#), [3](#), [6](#)
- [40] Sebastian Ruder. An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098*, 2017. [1](#)
- [41] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *Computer Vision—ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5–11, 2010, Proceedings, Part IV 11*, pages 213–226. Springer, 2010. [6](#)
- [42] Mariya Toneva, Alessandro Sordani, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J. Gordon. An empirical study of example forgetting during deep neural network learning. *CoRR*, abs/1812.05159, 2018. [5](#)
- [43] Guido M. van de Ven and Andreas Savas Tolias. Three scenarios for continual learning. *ArXiv*, abs/1904.07734, 2019. [3](#)
- [44] Chunnan Wang, Xiang Chen, Junzhe Wang, and Hongzhi Wang. ATPFL: automatic trajectory prediction model design under federated learning framework. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18–24, 2022*, pages 6553–6562. [1](#), [3](#)
- [45] Haozhao Wang, Yichen Li, Wenchao Xu, Ruixuan Li, Yufeng Zhan, and Zhigang Zeng. Dafkd: Domain-aware federated knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20412–20421, 2023. [3](#)
- [46] Jie Xu, Benjamin S Glicksberg, Chang Su, Peter Walker, Jiang Bian, and Fei Wang. Federated learning for healthcare informatics. *Journal of Healthcare Informatics Research*, 5: 1–19, 2021. [1](#)
- [47] Dong Yin, Mehrdad Farajtabar, and Ang Li. Sola: Continual learning with second-order loss approximation. *ArXiv*, abs/2006.10974, 2020. [3](#)
- [48] Jaehong Yoon, Wonyong Jeong, Giwoong Lee, Eunho Yang, and Sung Ju Hwang. Federated continual learning with weighted inter-client transfer. In *International Conference on Machine Learning*, pages 12073–12086. PMLR, 2021. [1](#)
- [49] Lu Yu, Bartłomiej Twardowski, Xialei Liu, Luis Herranz, Kai Wang, Yongmei Cheng, Shangling Jui, and Joost van de Weijer. Semantic drift compensation for class-incremental learning. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6980–6989, 2020. [3](#)
- [50] Lin Zhang, Li Shen, Liang Ding, Dacheng Tao, and Ling-Yu Duan. Fine-tuning global model via data-free knowledge distillation for non-iid federated learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18–24, 2022*, pages 10164–10173. [3](#)
- [51] Zhuangdi Zhu, Junyuan Hong, and Jiayu Zhou. Data-free knowledge distillation for heterogeneous federated learning. In *International Conference on Machine Learning*, pages 12878–12889. PMLR, 2021. [3](#)