

Federated Continual Learning with Weighted Inter-client Transfer

Jaehong Yoon^{1*} Wonyong Jeong^{1,2*} Giwoong Lee¹ Eunho Yang^{1,2} Sung Ju Hwang^{1,2}

Abstract

There has been a surge of interest in continual learning and federated learning, both of which are important in deep neural networks in real-world scenarios. Yet little research has been done regarding the scenario where each client learns on a sequence of tasks from a private local data stream. This problem of *federated continual learning* poses new challenges to continual learning, such as utilizing knowledge from other clients, while preventing interference from irrelevant knowledge. To resolve these issues, we propose a novel federated continual learning framework, *Federated Weighted Inter-client Transfer (FedWeIT)*, which decomposes the network weights into global federated parameters and sparse task-specific parameters, and each client receives selective knowledge from other clients by taking a weighted combination of their task-specific parameters. *FedWeIT* minimizes interference between incompatible tasks, and also allows positive knowledge transfer across clients during learning. We validate our *FedWeIT* against existing federated learning and continual learning methods under varying degrees of task similarity across clients, and our model significantly outperforms them with a large reduction in the communication cost. Code is available at <https://github.com/wyjeong/FedWeIT>.

1. Introduction

Continual learning (Thrun, 1995; Kumar & Daume III, 2012; Ruvolo & Eaton, 2013; Kirkpatrick et al., 2017; Schwarz et al., 2018) describes a learning scenario where a model continuously trains on a sequence of tasks; it is inspired by the human learning process, as a person learns to perform numerous tasks with large diversity over his/her lifespan,

*Equal contribution ¹Korea Advanced Institute of Science and Technology (KAIST), South Korea ²AITRICS, South Korea. Correspondence to: Jaehong Yoon <jaehong.yoon@kaist.ac.kr>, Wonyong Jeong <wyjeong@kaist.ac.kr>.

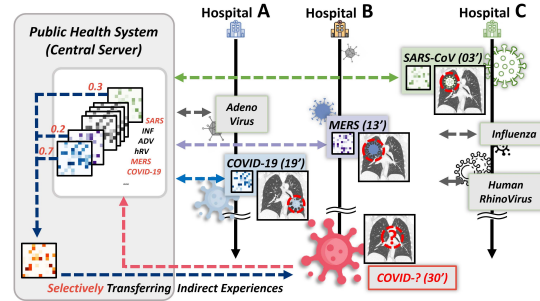


Figure 1. **Concept.** A continual learner at a hospital which learns on sequence of disease prediction tasks may want to utilize relevant task parameters from other hospitals. FCL allows such inter-client knowledge transfer via the communication of task-decomposed parameters.

making use of the past knowledge to learn about new tasks without forgetting previously learned ones. Continual learning is a long-studied topic since having such an ability leads to the potential of building a general artificial intelligence. However, there are crucial challenges in implementing it with conventional models such as deep neural networks (DNNs), such as *catastrophic forgetting*, which describes the problem where parameters or semantic representations learned for the past tasks drift to the direction of new tasks during training. The problem has been tackled by various prior work (Kirkpatrick et al., 2017; Shin et al., 2017; Riemer et al., 2019). More recent works tackle other issues, such as scalability or order-robustness (Schwarz et al., 2018; Hung et al., 2019; Yoon et al., 2020).

However, all of these models are fundamentally limited in that the models can only learn from its direct experience - they only learn from the sequence of the tasks they have trained on. Contrarily, humans can learn from *indirect experience* from others, through different means (e.g. verbal communications, books, or various media). Then wouldn't it be beneficial to implement such an ability to a continual learning framework, such that multiple models learning on different machines can learn from the knowledge of the tasks that have been already experienced by other clients? One problem that arises here, is that due to data privacy on individual clients and exorbitant communication cost, it may not be possible to communicate data directly between the clients or between the server and clients. Federated learning (McMahan et al., 2016; Li et al., 2018; Yurochkin

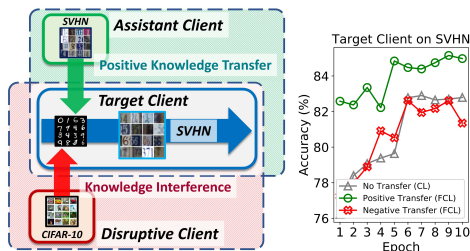


Figure 2. Challenge of Federated Continual Learning. Interference from other clients, resulting from sharing irrelevant knowledge, may hinder an optimal training of target clients (Red) while relevant knowledge from other clients will be beneficial for their learning (Green).

et al., 2019) is a learning paradigm that tackles this issue by communicating the parameters instead of the raw data itself. We may have a server that receives the parameters locally trained on multiple clients, aggregates it into a single model parameter, and sends it back to the clients. Motivated by our intuition on learning from indirect experience, we tackle the problem of *Federated Continual Learning (FCL)* where we perform continual learning with multiple clients trained on private task sequences, which communicate their task-specific parameters via a global server.

Figure 1 depicts an example scenario of FCL. Suppose that we are building a network of hospitals, each of which has a disease diagnosis model which continuously learns to perform diagnosis given CT scans, for new types of diseases. Then, under our framework, any diagnosis model which has learned about a new type of disease (e.g., COVID-19) will transmit the task-specific parameters to the global server, which will redistribute them to other hospitals for the local models to utilize. This allows all participants to benefit from the new task knowledge without compromising the data privacy.

Yet, the problem of federated continual learning also brings new challenges. First, there is not only the catastrophic forgetting from continual learning, but also the **threat of potential interference from other clients**. Figure 2 describes this challenge with the results of a simple experiment. Here, we train a model for MNIST digit recognition while communicating the parameters from another client trained on a different dataset. When the knowledge transferred from the other client is relevant to the target task (SVHN), the model starts with high accuracy, converge faster and reach higher accuracy (**green line**), whereas the model underperforms the base model if the transferred knowledge is from a task highly different from the target task (CIFAR-10, **red line**). Thus, we need to *selective utilize* knowledge from other clients to minimize the *inter-client interference* and maximize *inter-client knowledge transfer*. Another problem with the federated learning is efficient communication, as communication cost could become excessively large when

utilizing the knowledge of the other clients, since the communication cost could be the main bottleneck in practical scenarios when working with edge devices. Thus we want the knowledge to be represented as compactly as possible.

To tackle these challenges, we propose a novel framework for federated continual learning, *Federated Weighted Inter-client Transfer (FedWeIT)*, which decomposes the local model parameters into a dense base parameter and sparse task-adaptive parameters. FedWeIT reduces the interference between different tasks since the base parameters will encode task-generic knowledge, while the task-specific knowledge will be encoded into the task-adaptive parameters. When we utilize the generic knowledge, we also want the client to selectively utilize task-specific knowledge obtained at other clients. To this end, we allow each model to take a weighted combination of the task-adaptive parameters broadcast from the server, such that it can select task-specific knowledge helpful for the task at hand. FedWeIT is communication-efficient, since the task-adaptive parameters are *highly sparse* and only need to be communicated once when created. Moreover, when communication efficiency is not a critical issue as in cross-silo federated learning (Kairouz et al., 2019), we can use our framework to incentivize each client based on the attention weights on its task-adaptive parameters. We validate our method on multiple different scenarios with varying degree of task similarity across clients against various federated learning and local continual learning models. The results show that our model obtains significantly superior performance over all baselines, adapts faster to new tasks, with largely reduced communication cost. The main contributions of this paper are as follows:

- We introduce a **new problem of Federated Continual Learning (FCL)**, where multiple models continuously learn on distributed clients, which poses new challenges such as prevention of inter-client interference and inter-client knowledge transfer.
- We propose a **novel and communication-efficient framework for federated continual learning**, which allows each client to adaptively update the federated parameter and selectively utilize the past knowledge from other clients, by communicating sparse parameters.

2. Related Work

Continual learning While continual learning (Kumar & Daume III, 2012; Ruvolo & Eaton, 2013) is a long-studied topic with a vast literature, we only discuss recent relevant works. **Regularization-based:** EWC (Kirkpatrick et al., 2017) leverages Fisher Information Matrix to restrict the change of the model parameters such that the model finds solution that is good for both previous and the current task, and IMM (Lee et al., 2017) proposes to learn the posterior

distribution for multiple tasks as a mixture of Gaussians. Stable SGD (Mirzadeh et al., 2020) shows impressive performance gain through controlling essential hyperparameters and gradually decreasing learning rate each time a new task arrives. **Architecture-based:** DEN (Yoon et al., 2018) tackles this issue by expanding the networks size that are necessary via iterative neuron/filter pruning and splitting, and RCL (Xu & Zhu, 2018) tackles the same problem using reinforcement learning. APD (Yoon et al., 2020) additively decomposes the parameters into shared and task-specific parameters to minimize the increase in the network complexity. **Coreset-based:** GEM variants (Lopez-Paz & Ranzato, 2017; Chaudhry et al., 2019) minimize the loss on both of actual dataset and stored episodic memory. FRCL (Titsias et al., 2020) memorizes approximated posteriors of previous tasks with sophisticatedly constructed inducing points. To the best of our knowledge, none of the existing approaches considered the communicability for continual learning of deep neural networks, which we tackle. CoLLA (Rostami et al., 2018) aims at solving multi-agent lifelong learning with sparse dictionary learning, it does not have a central server to guide collaboration among clients and is formulated by a simple dictionary learning problem, thus not applicable to modern neural networks. Also, CoLLA is restricted to synchronous training with homogeneous clients.

Federated learning Federated learning is a distributed learning framework under differential privacy, which aims to learn a global model on a server while aggregating the parameters learned at the clients on their private data. FedAvg (McMahan et al., 2016) aggregates the model trained across multiple clients by computing a weighted average of them based on the number of data points trained. FedProx (Li et al., 2018) trains the local models with a proximal term which restricts their updates to be close to the global model. FedCurv (Shoham et al., 2019) aims to minimize the model disparity across clients during federated learning by adopting a modified version of EWC. Recent works (Yurochkin et al., 2019; Wang et al., 2020) introduce well-designed aggregation policies by leveraging Bayesian non-parametric methods. A crucial challenge of federated learning is **the reduction of communication cost**. TWAFL (Chen et al., 2019) tackles this problem by performing layer-wise parameter aggregation, where shallow layers are aggregated at every step, but deep layers are aggregated in the last few steps of a loop. (Karimireddy et al., 2020) suggests an algorithm for rapid convergence, which minimizes the interference among discrepant tasks at clients by sacrificing the local optimality. This is an opposite direction from personalized federated learning methods (Fallah et al., 2020; Lange et al., 2020; Deng et al., 2020) which put more emphasis on the performance of local models. FCL is a parallel research direction to both, and to the best of our knowledge, ours is the first work that considers task-

incremental learning of clients under federated learning framework.

3. Federated Continual Learning with FedWeIT

Motivated by the human learning process from indirect experiences, we introduce a novel continual learning under federated learning setting, which we refer to as *Federated Continual Learning (FCL)*. FCL assumes that multiple clients are trained on a sequence of tasks from private data stream, while communicating the learned parameters with a global server. We first formally define the problem in Section 3.1, and then propose naive solutions that straightforwardly combine the existing federated learning and continual learning methods in Section 3.2. Then, following Sections 3.3 and 3.4, we discuss about two novel challenges that are introduced by federated continual learning, and propose a novel framework, *Federated Weighted Inter-client Transfer (FedWeIT)* which can effectively handle the two problems while also reducing the client-to-server communication cost.

3.1. Problem Definition

In the standard continual learning (on a single machine), the model iteratively learns from a sequence of tasks $\{\mathcal{T}^{(1)}, \mathcal{T}^{(2)}, \dots, \mathcal{T}^{(T)}\}$ where $\mathcal{T}^{(t)}$ is a labeled dataset of t^{th} task, $\mathcal{T}^{(t)} = \{\mathbf{x}_i^{(t)}, \mathbf{y}_i^{(t)}\}_{i=1}^{N_t}$, which consists of N_t pairs of instances $\mathbf{x}_i^{(t)}$ and their corresponding labels $\mathbf{y}_i^{(t)}$. Assuming the most realistic situation, we consider the case where the task sequence is a task stream with an unknown arriving order, such that the model can access $\mathcal{T}^{(t)}$ only at the training period of task t which becomes inaccessible afterwards. Given $\mathcal{T}^{(t)}$ and the model learned so far, the learning objective at task t is as follows: minimize $\mathcal{L}(\boldsymbol{\theta}^{(t)}; \boldsymbol{\theta}^{(t-1)}, \mathcal{T}^{(t)})$,

where $\boldsymbol{\theta}^{(t)}$ is a set of the model parameters at task t .

We now extend the conventional continual learning to the federated learning setting with multiple clients and a global server. Let us assume that we have C clients, where at each client $c_c \in \{c_1, \dots, c_C\}$ trains a model on a *privately accessible* sequence of tasks $\{\mathcal{T}_c^{(1)}, \mathcal{T}_c^{(2)}, \dots, \mathcal{T}_c^{(t)}\} \subseteq \mathcal{T}$. Please note that there is no relation among the tasks $\mathcal{T}_{1:C}^{(t)}$ received at step t , across clients. Now the goal is to effectively train C continual learning models on their own private task streams, via communicating the model parameters with the global server, which aggregates the parameters sent from each client, and redistributes them to clients.

3.2. Communicable Continual Learning

In conventional federated learning settings, the learning is done with multiple rounds of local learning and parameter aggregation. At each round of communication r , each

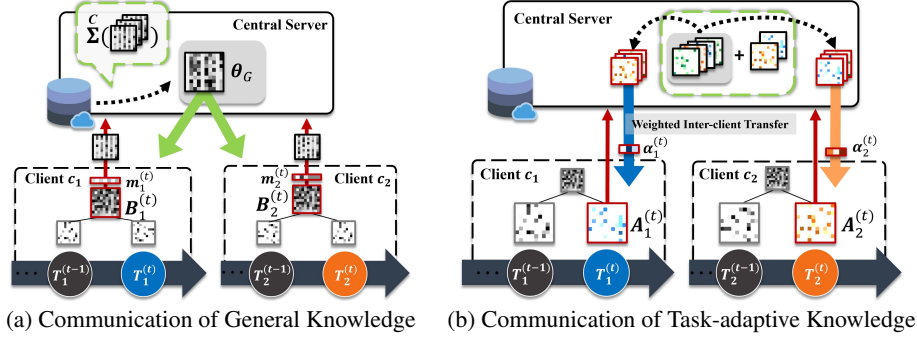


Figure 3. Updates of FedWeIT. (a) A client sends sparsified federated parameter $B_c \odot m_c^{(t)}$. After that, the server redistributes aggregated parameters to the clients. (b) The knowledge base stores previous tasks-adaptive parameters of clients, and each client selectively utilizes them with an attention mask.

client c_c and the server s perform the following two procedures: *local parameter transmission* and *parameter aggregation & broadcasting*. In the local parameter transmission step, for a randomly selected subset of clients at round r , $\mathcal{C}^{(r)} \subseteq \{c_1, c_2, \dots, c_C\}$, each client $c_c \in \mathcal{C}^{(r)}$ sends updated parameters $\theta_c^{(r)}$ to the server. The server-clients transmission is not done at every client because some of the clients may be temporarily disconnected. Then the server aggregates the parameters $\theta_c^{(r)}$ sent from the clients into a single parameter. The most popular frameworks for this aggregation are FedAvg (McMahan et al., 2016) and FedProx (Li et al., 2018). However, naive federated continual learning with these two algorithms on local sequences of tasks may result in catastrophic forgetting. One simple solution is to use a regularization-based, such as Elastic Weight Consolidation (EWC) (Kirkpatrick et al., 2017), which allows the model to obtain a solution that is optimal for both the previous and the current tasks. There exist other advanced solutions (Nguyen et al., 2018; Chaudhry et al., 2019) that successfully prevents catastrophic forgetting. However, the prevention of catastrophic forgetting at the client level is an orthogonal problem from federated learning.

Thus we focus on challenges that newly arise in this federated continual learning setting. In the federated continual learning framework, the aggregation of the parameters into a global parameter θ_G allows inter-client knowledge transfer across clients, since a task $\mathcal{T}_i^{(q)}$ learned at client c_i at round q may be similar or related to $\mathcal{T}_j^{(r)}$ learned at client c_j at round r . Yet, using a single aggregated parameter θ_G may be suboptimal in achieving this goal since knowledge from irrelevant tasks may not be useful or even hinder the training at each client by altering its parameters into incorrect directions, which we describe as *inter-client interference*. Another problem that is also practically important, is the *communication-efficiency*. Both the parameter transmission from the client to the server, and server to client will incur large communication cost, which will be problematic for the continual learning setting, since the clients may train on

possibly unlimited streams of tasks.

3.3. Federated Weighted Inter-client Transfer

How can we then maximize the *knowledge transfer* between clients while minimizing the *inter-client interference*, and communication cost? We now describe our model, *Federated Weighted Inter-client Transfer (FedWeIT)*, which can resolve these two problems that arise with a naive combination of continual learning approaches with federated learning framework.

The main cause of the problems, as briefly alluded to earlier, is that the knowledge of all tasks learned at multiple clients is stored into a single set of parameters θ_G . However, for the knowledge transfer to be effective, each client should *selectively* utilize only the knowledge of the *relevant* tasks that is trained at other clients. This **selective transfer** is also the key to minimize the inter-client interference as well as it will disregard the knowledge of irrelevant tasks that may interfere with learning.

We tackle this problem by decomposing the parameters, into three different types of the parameters with different roles: *global parameters* (θ_G) that capture the global and generic knowledge across all clients, *local base parameters* (B) which capture generic knowledge for each client, and *task-adaptive parameters* (A) for each specific task per client, motivated by Yoon et al. (2020). A set of the model parameters $\theta_c^{(t)}$ for task t at continual learning client c_c is then defined as follows:

$$\theta_c^{(t)} = B_c^{(t)} \odot m_c^{(t)} + A_c^{(t)} + \sum_{i \in \mathcal{C} \setminus c} \sum_{j < |t|} \alpha_{i,j}^{(t)} A_i^{(j)} \quad (1)$$

where $B_c^{(t)} \in \{\mathbb{R}^{I_l \times O_l}\}_{l=1}^L$ is the set of base parameters for c^{th} client shared across all tasks in the client, $m_c^{(t)} \in \{\mathbb{R}^{O_l}\}_{l=1}^L$ is the set of sparse vector masks which allows to adaptively transform $B_c^{(t)}$ for the task t , $A_c^{(t)} \in \{\mathbb{R}^{I_l \times O_l}\}_{l=1}^L$ is the set of a sparse task-adaptive parameters at client c_c . Here, L is the number of the layer in the neural

network, and I_l, O_l are input and output dimension of the weights at layer l , respectively.

The first term allows selective utilization of the global knowledge. We want the base parameter $\mathbf{B}_c^{(t)}$ at each client to capture generic knowledge across all tasks across all clients. In Figure 3 (a), we initialize it at each round t with the global parameter from the previous iteration, $\theta_G^{(t-1)}$ which aggregates the parameters sent from the client. This allows $\mathbf{B}_c^{(t)}$ to also benefit from the *global* knowledge about all the tasks. However, since $\theta_G^{(t-1)}$ also contains knowledge irrelevant to the current task, instead of using it as is, we learn the sparse mask $\mathbf{m}_c^{(t)}$ to select only the relevant parameters for the given task. This sparse parameter selection helps minimize inter-client interference, and also allows for efficient communication. The second term is the task-adaptive parameters $\mathbf{A}_c^{(t)}$. Since we additively decompose the parameters, this will learn to capture knowledge about the task that is not captured by the first term, and thus will capture specific knowledge about the task $\mathcal{T}_c^{(t)}$. The final term describes weighted inter-client knowledge transfer. We have a set of parameters that are *transmitted* from the server, which contain all task-adaptive parameters from all the clients. To selectively utilizes these indirect experiences from other clients, we further allocate attention $\alpha_c^{(t)}$ on these parameters, to take a weighted combination of them. By learning this attention, each client can select only the relevant task-adaptive parameters that help learn the given task. Although we design $\mathbf{A}_i^{(j)}$ to be highly sparse, using about 2 – 3% of memory of full parameter in practice, sending all task knowledge is not desirable. Thus we transmit the randomly sampled task-adaptive parameters across all time steps from knowledge base, which we empirically find to achieve good results in practice.

Training. We learn the decomposable parameter $\theta_c^{(t)}$ by optimizing for the following objective:

$$\begin{aligned} \underset{\mathbf{B}_c^{(t)}, \mathbf{m}_c^{(t)}, \mathbf{A}_c^{(1:t)}, \alpha_c^{(t)}}{\text{minimize}} \quad & \mathcal{L}(\theta_c^{(t)}; \mathcal{T}_c^{(t)}) + \lambda_1 \Omega(\{\mathbf{m}_c^{(t)}, \mathbf{A}_c^{(1:t)}\}) \\ & + \lambda_2 \sum_{i=1}^{t-1} \|\Delta \mathbf{B}_c^{(t)} \odot \mathbf{m}_c^{(i)} + \Delta \mathbf{A}_c^{(i)}\|_2^2, \end{aligned} \quad (2)$$

where \mathcal{L} is a loss function and $\Omega(\cdot)$ is a sparsity-inducing regularization term for all task-adaptive parameters and the masking variable (we use ℓ_1 -norm regularization), to make them sparse. The second regularization term is used for retroactive update of the past task-adaptive parameters, which helps the task-adaptive parameters to maintain the original solutions for the target tasks, by reflecting the change of the base parameter. Here, $\Delta \mathbf{B}_c^{(t)} = \mathbf{B}_c^{(t)} - \mathbf{B}_c^{(t-1)}$ is the difference between the base parameter at the current and previous timestep, and $\Delta \mathbf{A}_c^{(i)}$ is the difference between the task-adaptive parameter for task i at the current and pre-

vious timestep. This regularization is essential for preventing catastrophic forgetting. λ_1 and λ_2 are hyperparameters controlling the effect of the two regularizers.

Algorithm 1 Federated Weighted Inter-client Transfer

input Dataset $\{\mathcal{D}_c^{(1:t)}\}_{c=1}^C$, global parameter θ_G ,
 hyperparameters λ_1, λ_2 , knowledge base $kb \leftarrow \{\}$
output $\{\mathbf{B}_c, \mathbf{m}_c^{(1:t)}, \alpha_c^{(1:t)}, \mathbf{A}_c^{(1:t)}\}_{c=1}^C$
 1: Initialize \mathbf{B}_c to θ_G for all clients $\mathcal{C} \equiv \{c_1, \dots, c_C\}$
 2: **for** task $t = 1, 2, \dots$ **do**
 3: Randomly sample knowledge base $kb^{(t)} \sim kb$
 4: **for** round $r = 1, 2, \dots$ **do**
 5: Collect communicable clients $\mathcal{C}^{(r)} \sim \mathcal{C}$
 6: Distribute θ_G and $kb^{(t)}$ to client $c_c \in \mathcal{C}^{(r)}$ **if** c_c meets $kb^{(t)}$ first, **otherwise** distribute only θ_G
 7: Minimize Equation 2 for solving local CL problems
 8: $\mathbf{B}_c^{(t,r)} \odot \mathbf{m}_c^{(t,r)}$ are transmitted from $\mathcal{C}^{(r)}$ to the server
 9: Update $\theta_G \leftarrow \frac{1}{|\mathcal{C}^{(r)}|} \sum_c \mathbf{B}_c^{(t,r)} \odot \mathbf{m}_c^{(t,r)}$
 10: **end for**
 11: Update knowledge base $kb \leftarrow kb \cup \{\mathbf{A}_j^{(t)}\}_{j \in \mathcal{C}}$
 12: **end for**

3.4. Efficient Communication via Sparse Parameters

FedWeIT learns via server-to-client communication. As discussed earlier, a crucial challenge here is to reduce the communication cost. We describe what happens at the client and the server at each step.

Client: At each round r , each client c_c partially updates its base parameter with the nonzero components of the global parameter sent from the server; that is, $\mathbf{B}_c(n) = \theta_G(n)$ where n is a nonzero element of the global parameter. After training the model using Equation 2, it obtains a sparsified base parameter $\hat{\mathbf{B}}_c^{(t)} = \mathbf{B}_c^{(t)} \odot \mathbf{m}_c^{(t)}$ and task-adaptive parameter $\mathbf{A}_c^{(t)}$ for the new task, both of which are sent to the server, at smaller cost compared to naive FCL baselines. While naive FCL baselines require $|\mathcal{C}| \times R \times |\theta|$ for client-to-server communication, FedWeIT requires $|\mathcal{C}| \times (R \times |\hat{\mathbf{B}}| + |\mathbf{A}|)$ where R is the number of communication round per task and $|\cdot|$ is the number of parameters.

Server: The server first aggregates the base parameters sent from all the clients by taking an weighted average of them: $\theta_G = \frac{1}{c} \sum_c \hat{\mathbf{B}}_i^{(t)}$. Then, it broadcasts θ_G to all the clients. Task adaptive parameters of $t-1$, $\{\mathbf{A}_i^{(t-1)}\}_{i=1}^{C \setminus c}$ are broadcast at once per client during training task t . While naive FCL baselines requires $|\mathcal{C}| \times R \times |\theta|$ for server-to-client communication cost, FedWeIT requires $|\mathcal{C}| \times (R \times |\theta_G| + (|\mathcal{C}| - 1) \times |\mathbf{A}|)$ in which θ_G, \mathbf{A} are highly sparse. We describe the FedWeIT algorithm in Algorithm 1.

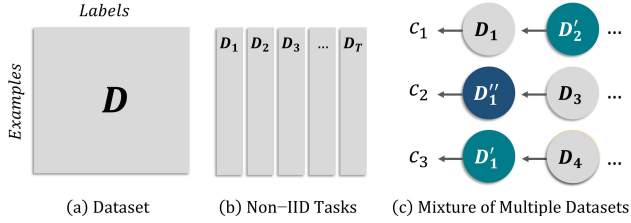


Figure 4. **Configuration of task sequences:** We first split a dataset D into multiple sub-tasks in non-IID manner ((a) and (b)). Then, we distribute them to multiple clients ($C_{\#}$). Mixed tasks from multiple datasets (colored circles) are distributed across all clients ((c)).

4. Experiments

We validate our **FedWeIT** under different configurations of task sequences against baselines which are namely Overlapped-CIFAR-100 and NonIID-50. **1) Overlapped-CIFAR-100:** We group 100 classes of CIFAR-100 dataset into 20 non-iid superclasses tasks. Then, we randomly sample 10 tasks out of 20 tasks and split instances to create a task sequence for each of the clients with overlapping tasks. **2) NonIID-50:** We use the following eight benchmark datasets: MNIST (LeCun et al., 1998), CIFAR-10/-100 (Krizhevsky & Hinton, 2009), SVHN (Netzer et al., 2011), Fashion-MNIST (Xiao et al., 2017), Not-MNIST (Bulatov, 2011), FaceScrub (Ng & Winkler, 2014), and TrafficSigns (Stal Kamp et al., 2011). We split the classes in the 8 datasets into 50 non-IID tasks, each of which is composed of 5 classes that are disjoint from the classes used for the other tasks. This is a large-scale experiment, containing 280,000 images of 293 classes from 8 heterogeneous datasets. After generating and processing tasks, we randomly distribute them to multiple clients as illustrated in Figure 4. We followed metrics for accuracy and forgetting from recent works (Chaudhry et al., 2020; Mirzadeh et al., 2020; 2021).

Experimental setup We use a modified version of LeNet (LeCun et al., 1998) for the experiments with both Overlapped-CIFAR-100 and NonIID-50 dataset. Further, we use ResNet-18 (He et al., 2016) with NonIID-50 dataset. We followed other experimental setups from (Serrà et al., 2018) and (Yoon et al., 2020). For detailed descriptions of the task configuration, metrics, hyperparameters, and more experimental results, please see **supplementary file**.

Baselines and our model **1) STL:** Single Task Learning at each arriving task. **2) EWC:** Individual continual learning with *EWC* (Kirkpatrick et al., 2017) per client. **3) Stable-SGD:** Individual continual learning with *Stable-SGD* (Mirzadeh et al., 2020) per client. **4) APD:** Individual continual learning with *APD* (Yoon et al., 2020) per client. **5) FedProx:** FCL using *FedProx* (Li et al., 2018) algorithm. **6) Scaffold:** FCL using *Scaffold* (Karimireddy et al., 2020)

Table 1. Average Per-task Performance on Overlapped-CIFAR-100 during FCL with 100 clients.

100 clients ($F=0.05, R=20, 1,000$ tasks in total)			
Methods	Accuracy	Forgetting	Model Size
FedProx	24.11 (± 0.44)	0.14 (± 0.01)	1.22 GB
FedCurv	29.11 (± 0.20)	0.09 (± 0.02)	1.22 GB
FedProx-SSGD	22.29 (± 0.51)	0.14 (± 0.01)	1.22 GB
FedProx-APD	32.55 (± 0.29)	0.02 (± 0.01)	6.97 GB
FedWeIT (Ours)	39.58 (± 0.27)	0.01 (± 0.00)	4.03 GB
STL	32.96 (± 0.23)	—	12.20 GB

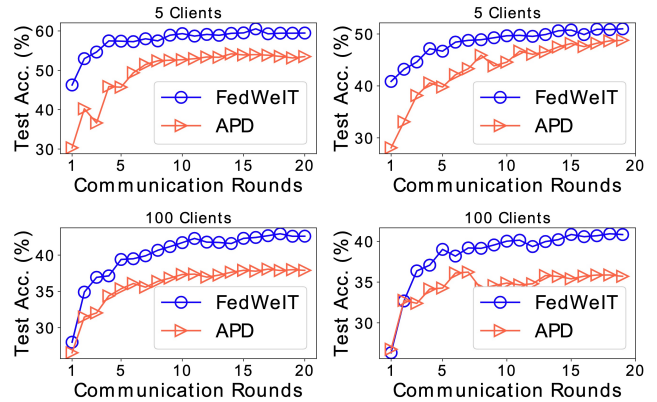


Figure 5. Averaged task adaptation during training last two (9^{th} and 10^{th}) tasks with 5 and 100 clients.

algorithm. **7) FedCurv:** FCL using *FedCurv* (Shoham et al., 2019) algorithm. **8) FedProx-[model]:** FCL, that is trained using *FedProx* algorithm with [model]. **9) FedWeIT:** Our FedWeIT algorithm.

4.1. Experimental Results

We first validate our model on both Overlapped-CIFAR-100 and NonIID-50 task sequences against single task learning (STL), continual learning (EWC, APD), federated learning (FedProx, Scaffold, FedCurv), and naive federated continual learning (FedProx-based) baselines. Table 2 shows the final average per-task performance after the completion of (federated) continual learning on both datasets. We observe that FedProx-based federated continual learning (FCL) approaches degenerate the performance of continual learning (CL) methods over the same methods without federated learning. This is because the aggregation of all client parameters that are learned on irrelevant tasks results in severe interference in the learning for each task, which leads to catastrophic forgetting and suboptimal task adaptation. Scaffold achieves poor performance on FCL, as its regularization on the local gradients is harmful for FCL, where all clients learn from a different task sequences. While FedCurv reduces inter-task disparity in parameters, it cannot minimize

Table 2. Averaged Per-task performance on both dataset during FCL with 5 clients (fraction=1.0). We measured task accuracy and model size after completing all learning phases over 3 individual trials. We also measured C2S/S2C communication cost for training each task.

NonIID-50 Dataset ($F=1.0, R=20$)						
Methods	FCL	Accuracy	Forgetting	Model Size	Client to Server Cost	Server to Client Cost
EWC (Kirkpatrick et al., 2017)	×	74.24 (± 0.11)	0.10 (± 0.01)	61 MB	N/A	N/A
Stable SGD (Mirzadeh et al., 2020)	×	76.22 (± 0.26)	0.14 (± 0.01)	61 MB	N/A	N/A
APD (Yoon et al., 2020)	×	81.42 (± 0.89)	0.02 (± 0.01)	90 MB	N/A	N/A
FedProx (Li et al., 2018)	✓	68.03 (± 2.14)	0.17 (± 0.01)	61 MB	1.22 GB	1.22 GB
Scaffold (Karimireddy et al., 2020)	✓	30.84 (± 1.41)	0.11 (± 0.02)	61 MB	2.44 GB	2.44 GB
FedCurv (Shoham et al., 2019)	✓	72.39 (± 0.32)	0.13 (± 0.02)	61 MB	1.22 GB	1.22 GB
FedProx-EWC	✓	68.27 (± 0.72)	0.12 (± 0.01)	61 MB	1.22 GB	1.22 GB
FedProx-Stable-SGD	✓	75.02 (± 1.44)	0.12 (± 0.01)	79 MB	1.22 GB	1.22 GB
FedProx-APD	✓	81.20 (± 1.52)	0.01 (± 0.01)	79 MB	1.22 GB	1.22 GB
FedWeIT (Ours)	✓	84.11 (± 0.27)	0.00 (± 0.00)	78 MB	0.37 GB	1.08 GB
Single Task Learning	×	85.78 (± 0.17)	—	610 MB	N/A	N/A

Overlapped CIFAR-100 Dataset ($F=1.0, R=20$)						
Methods	FCL	Accuracy	Forgetting	Model Size	Client to Server Cost	Server to Client Cost
EWC (Kirkpatrick et al., 2017)	×	44.26 (± 0.53)	0.13 (± 0.01)	61 MB	N/A	N/A
Stable SGD (Mirzadeh et al., 2020)	×	43.31 (± 0.44)	0.08 (± 0.01)	61 MB	N/A	N/A
APD (Yoon et al., 2020)	×	50.82 (± 0.41)	0.02 (± 0.01)	73 MB	N/A	N/A
FedProx (Li et al., 2018)	✓	38.96 (± 0.37)	0.13 (± 0.02)	61 MB	1.22 GB	1.22 GB
Scaffold (Karimireddy et al., 2020)	✓	22.80 (± 0.47)	0.09 (± 0.01)	61 MB	2.44 GB	2.44 GB
FedCurv (Shoham et al., 2019)	✓	40.36 (± 0.44)	0.15 (± 0.02)	61 MB	1.22 GB	1.22 GB
FedProx-EWC	✓	41.53 (± 0.39)	0.13 (± 0.01)	61 MB	1.22 GB	1.22 GB
FedProx-Stable-SGD	✓	43.29 (± 1.45)	0.07 (± 0.01)	61 MB	1.22 GB	1.22 GB
FedProx-APD	✓	52.20 (± 0.50)	0.02 (± 0.01)	75 MB	1.22 GB	1.22 GB
FedWeIT (Ours)	✓	55.16 (± 0.19)	0.01 (± 0.00)	75 MB	0.37 GB	1.07 GB
Single Task Learning	×	57.15 (± 0.07)	—	610 MB	N/A	N/A

inter-task interference, which results it to underperform single-machine CL methods. On the other hand, FedWeIT significantly outperforms both single-machine CL baselines and naive FCL baselines on both datasets. Even with larger number of clients ($C = 100$), FedWeIT consistently outperforms all baselines (Figure 5). This improvement largely owes to FedWeIT’s ability to selectively utilize the knowledge from other clients to rapidly adapt to the target task, and obtain better final performance.

The fast adaptation to new task is another clear advantage of inter-client knowledge transfer. To further demonstrate the practicality of our method with larger networks, we experiment on Non-IID dtaset with ResNet-18 (Table 3), on which FedWeIT still significantly outperforms the strongest baseline (FedProx-APD) while using fewer parameters.

Efficiency of FedWeIT We also report the accuracy as a function of network capacity in Tables 1 to 3, which we measure by the number of parameters used. We observe that FedWeIT obtains much higher accuracy while utilizing less number of parameters compared to FedProx-APD. This efficiency mainly comes from the reuse of task-adaptive parameters from other clients, which is not possible with single-machine CL methods or naive FCL methods.

We also examine the communication cost (the size of non-zero parameters transmitted) of each method. Table 2 re-

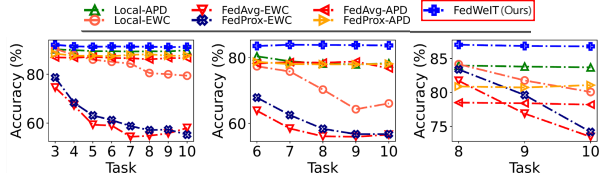


Figure 6. **Catastrophic forgetting.** Performance comparison about current task adaptation at 3rd, 6th and 8th tasks during federated continual learning on *NonIID-50*. We provide full version in our supplementary file.

ports both the *client-to-server* (C2S) / *server-to-client* (S2C) communication cost at training each task. FedWeIT, uses only 30% and 3% of parameters for \hat{B} and \hat{A} of the dense models respectively. We observe that FedWeIT is significantly more communication-efficient than FCL baselines although it broadcasts task-adaptive parameters, due to high sparsity of the parameters. Figure 7 shows the accuracy as a function of C2S cost according to a transmission of top- $\kappa\%$ informative parameters. Since FedWeIT selectively utilizes task-specific parameters learned from other clients, it results in superior performance over APD-baselines especially with sparse communication of model parameters.

Catastrophic forgetting Further, we examine how the performance of the past tasks change during continual learning, to see the severity of catastrophic forgetting with each method. Figure 6 shows the performance of FedWeIT and

Table 3. FCL results on NonIID-50 dataset with ResNet-18.

ResNet-18 ($F=1.0, R=20$)			
Methods	Accuracy	Forgetting	Model Size
APD	92.44 (± 0.17)	0.02 (± 0.00)	1.86 GB
FedProx-APD	92.89 (± 0.22)	0.02 (± 0.01)	2.05 GB
FedWeIT (Ours)	94.86 (± 0.13)	0.00 (± 0.00)	1.84 GB

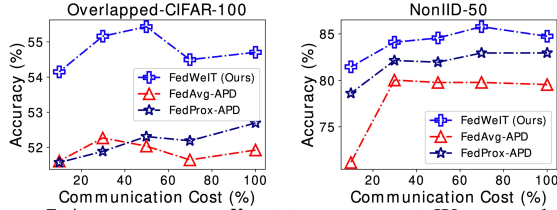


Figure 7. Accuracy over client-to-server cost. We report the relative communication cost to the original network. All results are averaged over the 5 clients.

FCL baselines on the 3rd, 6th and 8th tasks, at the end of training for later tasks. We observe that naive FCL baselines suffer from more severe catastrophic forgetting than local continual learning with EWC because of the *inter-client interference*, where the knowledge of irrelevant tasks from other clients overwrites the knowledge of the past tasks. Contrarily, our model shows no sign of catastrophic forgetting. This is mainly due to the selective utilization of the prior knowledge learned from other clients through the global/task-adaptive parameters, which allows it to effectively alleviate *inter-client interference*. FedProx-APD also does not suffer from catastrophic forgetting, but they yield inferior performance due to ineffective knowledge transfer.

Weighted inter-client knowledge transfer By analyzing the attention α in Equation 1, we examine which task parameters from other clients each client selected. Figure 8, shows example of the attention weights that are learned for the 0th split of *MNIST* and 10th split of *CIFAR-100*. We observe that large attentions are allocated to the task parameters from the same dataset (*CIFAR-100* utilizes parameters from *CIFAR-100* tasks with disjoint classes), or from a similar dataset (*MNIST* utilizes parameters from *Traffic Sign* and *SVHN*). This shows that FedWeIT effectively selects

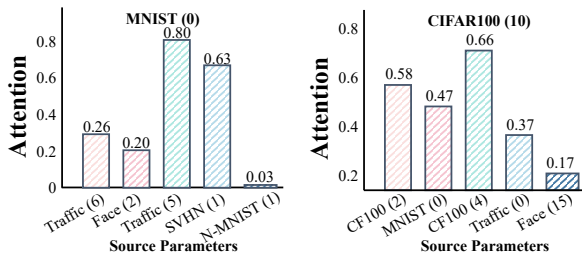
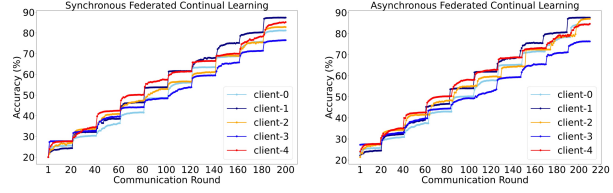


Figure 8. Inter-client transfer for NonIID-50. We compare the scale of the attentions at first FC layer which gives the weights on transferred task-adaptive parameters from other clients.


 Figure 9. FedWeIT with asynchronous federated continual learning on *Non-iid 50* dataset. We measure the test accuracy of all tasks per client.

beneficial parameters to maximize *inter-client* knowledge transfer. This is an impressive result since it does not know which datasets the parameters are trained on.

Asynchronous Federated Continual Learning We now consider FedWeIT under the asynchronous federated continual learning scenario, where there is no synchronization across clients for each task. This is a more realistic scenario since each task may require different training rounds to converge during federated continual learning. Here, *asynchronous* implies that each task requires different training costs (i.e., time, epochs, or rounds) for training. Under the asynchronous federated learning scenario, FedWeIT transfers any available task-adaptive parameters from the knowledge base to each client. In Figure 9, we plot the average test accuracy over all tasks during synchronous / asynchronous federated continual learning. In asynchronous FedWeIT, each task requires different training rounds and receives new tasks and task-adaptive parameters in an asynchronous manner and the performance of asynchronous FedWeIT is almost similar to that of the synchronous FedWeIT.

5. Conclusion

We tackled a novel problem of federated continual learning, which continuously learns local models at each client while allowing it to utilize indirect experience (task knowledge) from other clients. This poses new challenges such as *inter-client knowledge transfer* and prevention of *inter-client interference* between irrelevant tasks. To tackle these challenges, we additively decomposed the model parameters at each client into the global parameters that are shared across all clients, and sparse local task-adaptive parameters that are specific to each task. Further, we allowed each model to selectively update the global task-shared parameters and selectively utilize the task-adaptive parameters from other clients. The experimental validation of our model under various task similarity across clients, against existing federated learning and continual learning baselines shows that our model obtains significantly outperforms baselines with reduced communication cost. We believe that federated continual learning is a practically important topic of large interests to both research communities of continual learning and federated learning, that will lead to new research directions.

6. Acknowledgement

This work was supported by Samsung Research Funding Center of Samsung Electronics (No. SRFC-IT1502-51), Samsung Advanced Institute of Technology, Samsung Electronics Co., Ltd., Next-Generation Information Computing Development Program through the National Research Foundation of Korea(NRF) funded by the Ministry of Science, ICT & Future Plannig (No. 2016M3C4A7952634), the National Research Foundation of Korea(NRF) grant funded by the Korea government(MSIT) (2018R1A5A1059921), and Center for Applied Research in Artificial Intelligence (CARAI) grant funded by DAPA and ADD (UDI190031RD).

References

- Bulatov, Y. Not-mnist dataset. 2011.
- Chaudhry, A., Ranzato, M., Rohrbach, M., and Elhoseiny, M. Efficient lifelong learning with a-gem. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- Chaudhry, A., Khan, N., Dokania, P. K., and Torr, P. H. Continual learning in low-rank orthogonal subspaces. In *Advances in Neural Information Processing Systems (NIPS)*, 2020.
- Chen, Y., Sun, X., and Jin, Y. Communication-efficient federated deep learning with asynchronous model update and temporally weighted aggregation. *arXiv preprint arXiv:1903.07424*, 2019.
- Deng, Y., Kamani, M. M., and Mahdavi, M. Adaptive personalized federated learning. *arXiv preprint arXiv:2003.13461*, 2020.
- Fallah, A., Mokhtari, A., and Ozdaglar, A. Personalized federated learning: A meta-learning approach. *arXiv preprint arXiv:2002.07948*, 2020.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- Hung, C.-Y., Tu, C.-H., Wu, C.-E., Chen, C.-H., Chan, Y.-M., and Chen, C.-S. Compacting, picking and growing for unforgetting continual learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2019.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- Karimireddy, S. P., Kale, S., Mohri, M., Reddi, S. J., Stich, S. U., and Suresh, A. T. Scaffold: Stochastic controlled averaging for on-device federated learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2020.
- Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, pp. 201611835, 2017.
- Krizhevsky, A. and Hinton, G. E. Learning multiple layers of features from tiny images. Technical report, Computer Science Department, University of Toronto, 2009.
- Kumar, A. and Daume III, H. Learning task grouping and overlap in multi-task learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2012.
- Lange, M. D., Jia, X., Parisot, S., Leonardis, A., Slabaugh, G., and Tuytelaars, T. Unsupervised model personalization while preserving privacy and scalability: An open problem. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 1998.
- Lee, S.-W., Kim, J.-H., Jun, J., Ha, J.-W., and Zhang, B.-T. Overcoming catastrophic forgetting by incremental moment matching. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.
- Li, T., Sahu, A. K., Zaheer, M., Sanjabi, M., Talwalkar, A., and Smith, V. Federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*, 2018.
- Lopez-Paz, D. and Ranzato, M. Gradient episodic memory for continual learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.
- McMahan, H. B., Moore, E., Ramage, D., Hampson, S., et al. Communication-efficient learning of deep networks from decentralized data. *arXiv preprint arXiv:1602.05629*, 2016.
- Mirzadeh, S. I., Farajtabar, M., Pascanu, R., and Ghasemzadeh, H. Understanding the role of training regimes in continual learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2020.
- Mirzadeh, S. I., Farajtabar, M., Gorur, D., Pascanu, R., and Ghasemzadeh, H. Linear mode connectivity in multitask

- and continual learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images with unsupervised feature learning. 2011.
- Ng, H.-W. and Winkler, S. A data-driven approach to cleaning large face datasets. In *2014 IEEE international conference on image processing (ICIP)*, pp. 343–347. IEEE, 2014.
- Nguyen, C. V., Li, Y., Bui, T. D., and Turner, R. E. Variational continual learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- Riemer, M., Cases, I., Ajemian, R., Liu, M., Rish, I., Tu, Y., and Tesauro, G. Learning to learn without forgetting by maximizing transfer and minimizing interference. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- Rostami, M., Kolouri, S., Kim, K., and Eaton, E. Multi-agent distributed lifelong learning for collective knowledge acquisition. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents and Multiagent Systems, 2018.
- Ruvolo, P. and Eaton, E. Ella: An efficient lifelong learning algorithm. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2013.
- Schwarz, J., Luketina, J., Czarnecki, W. M., Grabska-Barwinska, A., Teh, Y. W., Pascanu, R., and Hadsell, R. Progress & compress: A scalable framework for continual learning. *arXiv preprint arXiv:1805.06370*, 2018.
- Serrà, J., Surís, D., Miron, M., and Karatzoglou, A. Overcoming catastrophic forgetting with hard attention to the task. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.
- Shin, H., Lee, J. K., Kim, J., and Kim, J. Continual learning with deep generative replay. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.
- Shoham, N., Avidor, T., Keren, A., Israel, N., Benditkis, D., Mor-Yosef, L., and Zeitak, I. Overcoming forgetting in federated learning on non-iid data. *arXiv preprint arXiv:1910.07796*, 2019.
- Stallkamp, J., Schlipsing, M., Salmen, J., and Igel, C. The german traffic sign recognition benchmark: a multi-class classification competition. In *The 2011 international joint conference on neural networks*, 2011.
- Thrun, S. *A Lifelong Learning Perspective for Mobile Robot Control*. Elsevier, 1995.
- Titsias, M. K., Schwarz, J., Matthews, A. G. d. G., Pascanu, R., and Teh, Y. W. Functional regularisation for continual learning with gaussian processes. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- Wang, H., Yurochkin, M., Sun, Y., Papailiopoulos, D., and Khazaeni, Y. Federated learning with matched averaging. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=BkluqlSFDS>.
- Xiao, H., Rasul, K., and Vollgraf, R. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Xu, J. and Zhu, Z. Reinforced continual learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2018.
- Yoon, J., Yang, E., Lee, J., and Hwang, S. J. Lifelong learning with dynamically expandable networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- Yoon, J., Kim, S., Yang, E., and Hwang, S. J. Scalable and order-robust continual learning with additive parameter decomposition. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- Yurochkin, M., Agarwal, M., Ghosh, S., Greenewald, K., Hoang, T. N., and Khazaeni, Y. Bayesian nonparametric federated learning of neural networks. *Proceedings of the International Conference on Machine Learning (ICML)*, 2019.

Federated Continual Learning with Weighted Inter-client Transfer: Supplementary File

Organization We provide in-depth descriptions and explanations that are not covered in the main document, and additionally report more experiments in this supplementary document, which organized as follows:

- **Section A** - We further describe the experimental details, including the network architecture, training configurations, forgetting measures, and datasets.
- **Section B** - We report additional experimental results, such as the effectiveness of the communication frequency (Section B.1) and ablation study on Overlapped-CIFAR-100 dataset (Section B.2).

A. Experimental Details

We further provide the experimental settings in detail, including the descriptions of the network architectures, hyperparameters, and dataset configuration.

Network Architecture We utilize a modified version of LeNet and a conventional ResNet-18 as the backbone network architectures for validation. In the LeNet, the first two layers are convolutional neural layers of 20 and 50 filters with the 5×5 convolutional kernels, which are followed by the two fully-connected layers of 800 and 500 units each. Rectified linear units activations and local response normalization are subsequently applied to each layers. We use 2×2 max-pooling after each convolutional layer. All layers are initialized based on the variance scaling method. Detailed description of the architecture for LeNet is given in Table A.4.

Configurations We use an Adam optimizer with adaptive learning rate decay, which decays the learning rate by a factor of 3 for every 5 epochs with no consecutive decrease in the validation loss. We stop training in advance and start learning the next task (if available) when the learning rate reaches ρ . The experiment for LeNet with 5 clients, we initialize by $1e^{-3} \times \frac{1}{3}$ at the beginning of each new task and $\rho = 1e^{-7}$. Mini-batch size is 100, the rounds per task is 20, an the epoch per round is 1. The setting for ResNet-18 is identical, excluding the initial learning rate, $1e^{-4}$. In the case of experiments with 20 and 100 clients, we set the same settings except reducing minibatch size from 100 to 10 with an initial learning rate $1e^{-4}$. We use client fraction 0.25 and 0.05, respectively, at each communication round.

Table A.4. Implementation Details of Base Network Architecture (LeNet). Note that T indicates the number of tasks that each client sequentially learns on.

Layer	Filter Shape	Stride	Output
Input	N/A	N/A	$32 \times 32 \times 3$
Conv 1	$5 \times 5 \times 20$	1	$32 \times 32 \times 20$
Max Pooling 1	3×3	2	$16 \times 16 \times 20$
Conv 2	$5 \times 5 \times 50$	1	$16 \times 16 \times 50$
Max Pooling 2	3×3	2	$8 \times 8 \times 50$
Flatten	3200	N/A	$1 \times 1 \times 3200$
FC 1	800	N/A	$1 \times 1 \times 800$
FC 2	500	N/A	$1 \times 1 \times 500$
Softmax	Classifier	N/A	$1 \times 1 \times 5 \times T$
Total Number of Parameters			3,012,920

we set $\lambda_1 = [1e^{-1}, 4e^{-1}]$ and $\lambda_2 = 100$ for all experiments. Further, we use $\mu = 5e^{-3}$ for FedProx, $\lambda = [1e^{-2}, 1.0]$ for EWC and FedCurv. We initialize the attention parameter $\alpha_c^{(t)}$ as sum to one, $\alpha_{c,j}^{(t)} \leftarrow 1/|\alpha_c^{(t)}|$.

Metrics. We evaluate all the methods on two metrics following the continual learning literature (Chaudhry et al., 2019; Mirzadeh et al., 2020).

1. **Averaged Accuracy:** We measure averaged test accuracy of all the tasks after the completion of a continual learning at task t by $A_t = \frac{1}{t} \sum_{i=1}^t a_{t,i}$, where $a_{t,i}$ is the test accuracy of task i after learning on task t .
2. **Averaged Forgetting:** We measure the forgetting as the averaged disparity between minimum task accuracy during continuous training. More formally, for T tasks, the forgetting can be defined as $F = \frac{1}{T-1} \sum_{i=1}^{T-1} \max_{t \in \{1, \dots, T-1\}} (a_{t,i} - a_{T,i})$.

Datasets We create both Overlapped-CIFAR-100 and NonIID-50 datasets. For Overlapped-CIFAR-100, we generate 20 non-iid tasks based on 20 superclasses, which hold 5 subclasses. We split instances of 20 tasks according to the number of clients (5, 20, and 100) and then distribute the tasks across all clients. For NonIID-50 dataset, we utilize 8 heterogenous datasets and create 50 non-iid tasks in total as shown in Table A.5. Then we arbitrarily select 10 tasks without duplication and distribute them to 5 clients. The average performance of single task learning on the dataset is $85.78 \pm 0.17(\%)$, measured by our base LeNet architecture.

Table A.5. Dataset Details of *NonIID-50* Task. We provide dataset details of *NonIID-50* dataset, including 8 heterogeneous datasets, number of sub-tasks, classes per sub-task, and instances of train, valid, and test sets.

NonIID-50						
Dataset	Num. Classes	Num. Tasks	Num. Classes (Task)	Num. Train	Num. Valid	Num. Test
CIFAR-100	100	15	5	36,750	10,500	5,250
Face Scrub	100	16	5	13,859	3,959	1,979
Traffic Signs	43	9	5 (3)	32,170	9,191	4,595
SVHN	10	2	5	61,810	17,660	8,830
MNIST	10	2	5	42,700	12,200	6,100
CIFAR-10	10	2	5	36,750	10,500	5,250
Not MNIST	10	2	5	11,339	3,239	1,619
Fashion MNIST	10	2	5	42,700	12,200	6,100
Total	293	50	248	278,078	39,723	79,449

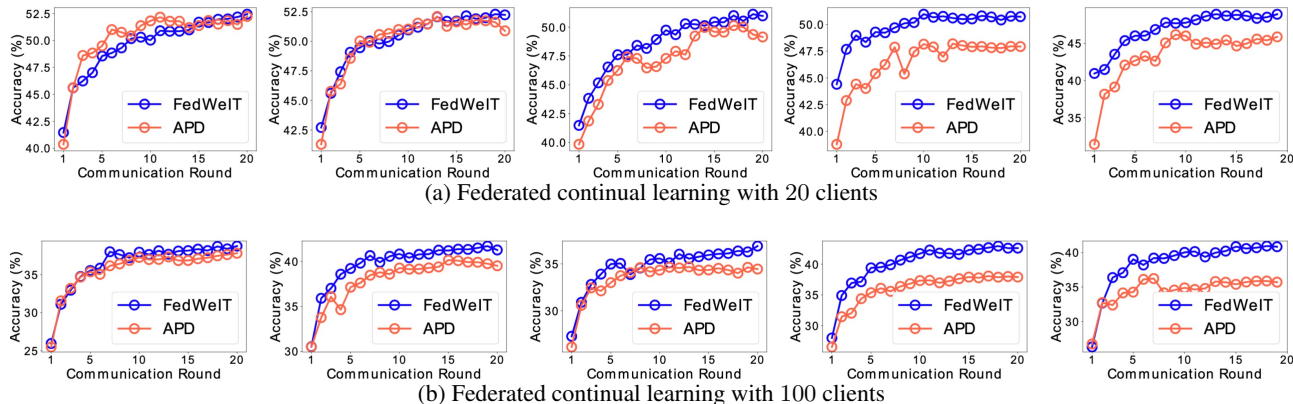


Figure A.10. Task adaptation comparison with FedWeIT and APD using 20 clients and 100 clients. We visualize the last 5 tasks out of 10 tasks per client. *Overlapped-CIFAR-100* dataset are used after splitting instances according to the number of clients (20 and 100).

Table B.6. Experimental results on the *Overlapped-CIFAR-100* dataset with 20 tasks. All results are the mean accuracies over 5 clients, averaged over 3 individual trials.

Overlapped-CIFAR-100 with 20 tasks			
Methods	Accuracy	M Size	C2S/S2C Cost
FedProx	29.76 ± 0.39	0.061 GB	1.22 / 1.22 GB
FedProx-EWC	27.80 ± 0.58	0.061 GB	1.22 / 1.22 GB
FedProx-APD	43.80 ± 0.76	0.093 GB	1.22 / 1.22 GB
FedWeIT	46.78 ± 0.14	0.092 GB	0.37 / 1.07 GB

B. Additional Experimental Results

We further include a quantitative analysis about the communication round frequency and additional experimental results across the number of clients.

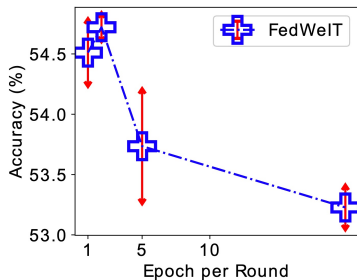
B.1. Effect of the Communication Frequency

We provide an analysis on the effect of the communication frequency by comparing the performance of the model, measured by the number of training epochs per communication

round. We run the 4 different FedWeIT with 1, 2, 5, and 20 training epochs per round. Figure A.11 shows the performance of our FedWeIT variants. As clients frequently update the model parameters through the communication with the central server, the model gets higher performance while maintaining smaller network capacity since the model with a frequent communication efficiently updates the model parameters as transferring the inter-client knowledge. However, it requires much heavier communication costs than the model with sparser communication. For example, the model trained for 1 epochs at each round may need to about 16.9 times larger entire communication cost than the model trained for 20 epochs at each round. Hence, there is a trade-off between model performance of federated continual learning and communication efficiency, whereas FedWeIT variants consistently outperform (federated) continual learning baselines.

B.2. Ablation Study with Model Components

We perform an ablation study to analyze the role of each component of our FedWeIT. We compare the performance of four different variations of our model. **w/o B communi-**



Overlapped-CIFAR-100				
Methods	Accuracy	Model Size	C2S/S2C Cost	Epochs / Round
FedWeIT	55.16 \pm 0.19	0.075 GB	0.37 / 1.07 GB	1
FedWeIT	55.18 \pm 0.08	0.077 GB	0.19 / 0.53 GB	2
FedWeIT	53.73 \pm 0.44	0.083 GB	0.08 / 0.22 GB	5
FedWeIT	53.22 \pm 0.14	0.088 GB	0.02 / 0.07 GB	20

Figure A.11. **Number of Epochs per Round** We show error bars over the number of training epochs per communication rounds on *Overlapped-CIFAR-100* with 5 clients. All models transmit full of local base parameters and highly sparse task-adaptive parameters. All results are the mean accuracy over 5 clients and we run 3 individual trials. Red arrows at each point describes the standard deviation of the performance.

Table B.7. Ablation studies to analyze the effectiveness of parameter decomposition on WeIT. All experiments performed on NonIID-50 dataset.

NonIID-50			
Methods	Acc.	M Size	C2S/S2C Cost
FedWeIT	84.11%	0.078 GB	0.37 / 1.07 GB
w/o B comm.	77.88%	0.070 GB	0.01 / 0.01 GB
w/o A comm.	79.21%	0.079 GB	0.37 / 1.04 GB
w/o A	65.66%	0.061 GB	0.37 / 1.04 GB
w/o m	78.71%	0.087 GB	1.23 / 1.25 GB

cation describes the model that does not transfer the base parameter **B** and only communicates task-adaptive ones. **w/o A communication** is the model that does not communicate task-adaptive parameters. **w/o A** is the model which trains the model only with sparse transmission of local base parameter, and **w/o m** is the model without the sparse vector mask. As shown in Table B.7, without communicating **B** or **A**, the model yields significantly lower performance compared to the full model since they do not benefit from *inter-client knowledge transfer*. The model **w/o A** obtains very low performance due to catastrophic forgetting, and the model **w/o** sparse mask **m** achieves lower accuracy with larger capacity and cost, which demonstrates the importance of performing selective transmission.

B.3. Ablation Study with Regularization Terms

We also perform the additional analysis by eliminating the proposed regularization terms in Table B.8. As described in Section 3.3, without ℓ_1 term, the method achieves even better performance but requires significantly larger memory. Without ℓ_2 term, the method suffers from the forgetting.

Table B.8. Ablation Study on Knowledge Transfer (NonIID-50)

Method	Avg. Accuracy	Memory size	BwT
FedWeIT	84.43% (\pm 0.50)	68.93 MB	-0.0014
FedWeIT w/o ℓ_1	87.12% (\pm 0.24)	354.41 MB	-0.0007
FedWeIT w/o ℓ_2	56.76% (\pm 0.84)	63.44 MB	-0.3203

B.4. Forgetting Analysis

In Figure B.12, we present forgetting performances of our baseline models, such as Local-EWC/APD, FedAvg-EWC/APD, and FedProx-EWC,APD, and our method, *FedWeIT*. As shown in the figure, EWC based methods, including local continual learning and combinations of FedAvg and FedProx, shows performance degradation while learning on new tasks. For example, the performance of the FedAvg-EWC (red with inversed triangular markers) is rapidly dropped from Task 1 to Task 10, which visualized in the left most plot on the first row. On the other hand, both APD based method and our method shows compelling ability to prevent *catastrophic forgetting* regardless how many tasks and what tasks the clients learn afterwards. We provide corresponding results in Table 2 in the main document.

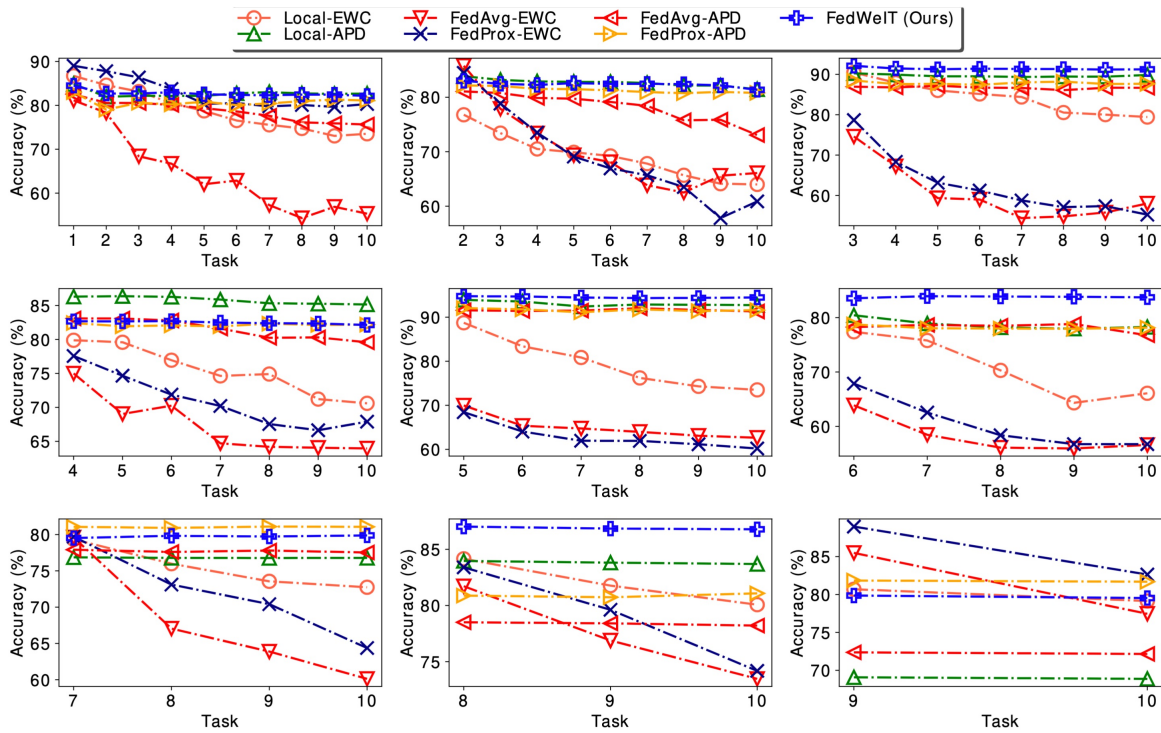


Figure B.12. **Forgetting Analysis** Performance change over the increasing number of tasks for all tasks except the last task (1st to 9th) during federated continual learning on *NonIID-50*. We observe that our method does not suffer from task forgetting on any tasks.