
FEDGP: BUFFER-BASED GRADIENT PROJECTION FOR CONTINUAL FEDERATED LEARNING

Shenghong Dai¹, Yicong Chen¹, Jy-yong Sohn², S M Iftekharul Alam³, Ravikumar Balakrishnan³,

Suman Banerjee⁴, Nageen Himayat³, Kangwook Lee¹

¹Department of Electrical and Computer Engineering,
University of Wisconsin-Madison, Madison, WI 53706, USA
{sdai37, ychen2229, kangwook.lee}@wisc.edu

²Department of Applied Statistics, Yonsei University, Seoul, Korea

³Intel Labs, Santa Clara, CA 95054, USA

⁴Department of Computer Sciences, University of Wisconsin-Madison, Madison, WI 53706, USA

ABSTRACT

Continual Federated Learning (CFL) is essential for enabling real-world applications where multiple decentralized clients adaptively learn from continuous data streams. A significant challenge in CFL is mitigating catastrophic forgetting, where models lose previously acquired knowledge when learning new information. Existing works on this issue either make unrealistic assumptions about the availability of task boundaries or face difficulties due to the constraints of device storage capacities and the heterogeneous nature of data distributions among clients. To address these limitations, we introduce a buffer-based Gradient Projection method (FedGP). This method tackles catastrophic forgetting by leveraging local buffer samples and aggregated buffer gradients, thus preserving knowledge across multiple clients. Our method is compatible with existing CFL techniques, enhancing their performance in the CFL context. Our experiments on standard benchmarks show consistent performance improvements across diverse scenarios. For example, on a task-incremental learning setting with CIFAR100, our method can help increase the accuracy up to 27%. Our code is available at <https://github.com/shenghongdai/FedGP/>.

1 INTRODUCTION

Federated Learning (FL) is a machine learning technique that facilitates collaborative model training among a large number of users while keeping data decentralized for privacy and efficient communication. In real-world applications, models trained via FL need the flexibility to continuously adapt to new data streams without forgetting past knowledge. This is critical in a variety of scenarios, such as autonomous vehicles, which must adapt to changes in the surroundings like new buildings or vehicle types without losing proficiency in previously encountered contexts. These real-world considerations make it essential to integrate FL with continual learning (CL) (Shmelkov et al., 2017; Chaudhry et al., 2018; Thrun, 1995; Aljundi et al., 2017; Chen & Liu, 2018; Aljundi et al., 2018), thereby giving rise to the concept of Continual Federated Learning (CFL).

The biggest challenge in CFL, as in CL, is *catastrophic forgetting*, where the model gradually shifts its focus from old data to new data and unintentionally discards previously acquired knowledge. Initial attempts to mitigate catastrophic forgetting in CFL incorporated existing CL solutions at each client

A preliminary version of this work was presented at the Federated Learning Systems (FLSys) Workshop @ Sixth Conference on Machine Learning and Systems, June 2023.

of FL, such as replaying previous task data or penalizing the updates of weights that are crucial for preserving the knowledge from earlier tasks. However, recent works (Bakman et al., 2023; Ma et al., 2022; Yoon et al., 2021) have observed that this naïve approach cannot fully mitigate the problem due to two reasons: (i) small-scale devices participating in FL only have limited buffer size to store the data from previous tasks, (ii) data distributions are not identical across clients in FL. Moreover, existing methods developed for CFL suffer from several limitations. These include scalability issues as the number of tasks grows (Yoon et al., 2021; Venkatesha et al., 2022), the need for significant effort in generating or collecting surrogate data (Ma et al., 2022), and significant communication overhead (Yao & Sun, 2020). A crucial constraint shared by all these methods is that they require explicit task boundaries. Mitigating catastrophic forgetting in practical scenarios where fixed task boundaries are absent throughout the training process, known as *general continual learning* (Buzzega et al., 2020), remains an important open question.

To address these existing challenges of CFL, we introduce a method called buffer-based Gradient Projection, which we dub FedGP. Our approach, illustrated in Fig. 1, involves two key components:

1. **Global Buffer Gradients:** Each client k computes the local buffer gradient g_{ref}^k of the global model with respect to its local buffer data. All local buffer gradients are then securely averaged to obtain aggregated buffer gradient g_{ref} .
2. **Local Gradient Projection:** In the next round, each client k updates its local model such that the direction for the model update does not conflict with aggregated buffer gradient g_{ref} from the previous round, ensuring each client preserves past information from all clients.

Importantly, our FedGP method is designed to be fully compatible with (i) general continual learning settings, when task boundary is unknown, and (ii) secure aggregation techniques (Bonawitz et al., 2017). Secure aggregation ensures that while clients share gradients or model updates, the individual data remains private (Bakman et al., 2023).

Our contributions: We introduce a new method for CFL, called FedGP. This method utilizes information from previous tasks across clients to effectively mitigate catastrophic forgetting, without having access to task boundaries. Furthermore, FedGP can seamlessly integrate with existing CFL techniques to enhance performance.

We conduct comprehensive experiments to demonstrate the effectiveness of FedGP across various standard image classification benchmarks and a text classification task. FedGP consistently improves accuracy and reduces forgetting on top of existing CFL baselines across diverse benchmark datasets. Further, we evaluate the robustness of our method considering various buffer sizes, communication frequency, asynchronous environments, and different numbers of tasks and users. This evaluation also includes an ablation study to examine each contributions of FedGP key components.

2 RELATED WORK

Prior work related to our paper falls into three categories: Continual Learning (CL), Federated Learning (FL), and Continual Federated Learning (CFL).

2.1 CONTINUAL LEARNING (CL)

CL addresses the problem of learning multiple tasks consecutively using a single model. Catastrophic forgetting (McCloskey & Cohen, 1989; Ratcliff, 1990; French, 1999), where a classifier trained for

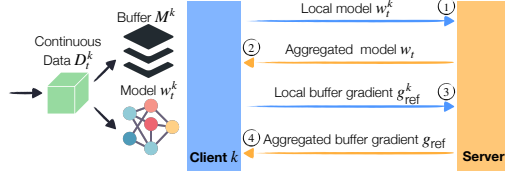


Figure 1: An overview of our proposed method, FedGP. In each round t , client k receives data D_t^k and trains a local model w_t^k . To address catastrophic forgetting, a portion of the incoming data is stored in a buffer \mathcal{M}^k . Given the aggregated model w_t provided by the central server, each client computes the gradient with respect to w_t using its buffer data \mathcal{M}^k . The server securely aggregates the local buffer gradients from all clients to obtain an aggregated buffer gradient g_{ref} , which will guide the local model update for each client k in the subsequent round.

a current task performs poorly on previous tasks, is a major challenge. Existing approaches can be categorized into regularization-based, architecture-based, and replay-based methods.

Regularization-based methods Some CL methods add a regularization term in the loss used for the model update; they penalize the updates on weights that are important for previous tasks. EWC (Kirkpatrick et al., 2017), SI (Zenke et al., 2017), Riemannian Walk (Chaudhry et al., 2018) are methods within this category. EWC uses Fisher information matrix to evaluate the importance of parameters for previous tasks. Besides, LwF (Li & Hoiem, 2017) leverages knowledge distillation to preserve outputs on previous tasks while learning the current task.

Architecture-based methods A class of CL methods assigns a subset of model parameters to each task, so that different tasks are learned by different parameters. This class of methods is also known as parameter isolation methods. Some methods including PNN (Rusu et al., 2016) and DEN (Yoon et al., 2017) uses dynamic architectures where the architecture changes dynamically as the number of tasks increases. These methods have issues where the number of required parameters grows linearly with the number of tasks. To tackle this issue, fixed network are used in the recent methods including PackNet (Mallya & Lazebnik, 2018), HAT (Serra et al., 2018) and PathNet (Fernando et al., 2017). SupSup (Wortsman et al., 2020) and DualNet (Pham et al., 2021) are the latest SOTA methods.

Replay-based methods To avoid catastrophic forgetting, a class of CL methods employs a replay buffer to save a small portion of the data seen in previous tasks and reuse it in the training of subsequent tasks. One of the early works in this area is ER (Ratcliff, 1990; Robins, 1995). In more contemporary studies, iCaRL (Rebuffi et al., 2017) stores exemplars of data from previous tasks and adds distillation loss for old exemplars to mitigate the forgetting issue. Deep Generative Replay (Shin et al., 2017) retains the memories of the previous tasks by loading the synthetic data generated by GANs without replaying the actual data for the previous tasks. GSS (Aljundi et al., 2019) optimally selects data for replay buffer by maximizing the diversity of samples in terms of the gradient in the parameter space. GEM (Lopez-Paz & Ranzato, 2017) and its variant A-GEM (Chaudhry et al., 2019) leverage an episodic memory that stores part of seen samples for each task to prevent forgetting old knowledge. Similarly, OGD (Farajtabar et al., 2020) stores gradients as opposed to actual data, providing a reference in projection. More recent work include GDumb (Prabhu et al., 2020), BiC (Wu et al., 2019), DER++, and Co²L (Cha et al., 2021). Despite its simplicity, replay-based techniques have shown great performances on multiple benchmarks (Mai et al., 2022; Parisi et al., 2019). FedGP leverages a replay-based method that alleviates forgetting by reusing some data from previous tasks.

General continual learning Prior works on CL often rely on the information about the task boundaries. For example, some regularization-based methods store network responses at task boundaries; architecture-based methods update the model architecture after one task is finished; some replay-based methods perform specific steps specifically at task boundaries. However, when dealing with some data in practical settings, task boundaries are not clearly defined. This scenario, where sequential tasks are learned continuously without explicit knowledge of task boundaries, is referred to as *general continual learning* (Buzzega et al., 2020; Aljundi et al., 2019; Chaudhry et al., 2019). To address general continual learning, replay-based methods can utilize reservoir sampling (Vitter, 1985), which allows sampling throughout the training rather than relying on task boundaries. In our work, we specifically focus on general continual learning with reservoir sampling, particularly in the context of *federated learning* setups.

2.2 FEDERATED LEARNING (FL)

FL enables collaborative training of a model with improved data privacy (Kairouz et al., 2021; Lim et al., 2020; Zhao et al., 2018; Konečný et al., 2016). FedAvg (McMahan et al., 2017) is a widely used FL algorithm, but most existing methods (Li et al., 2020; Shoham et al., 2019; Karimireddy et al., 2020; Li et al., 2019; Mohri et al., 2019) assume static data distribution over time, ignoring temporal dynamics.

2.3 CONTINUAL FEDERATED LEARNING (CFL)

CFL tackles the problem of learning multiple consecutive tasks in the FL setup. FedProx (Li et al., 2020) and FedCurv (Shoham et al., 2019) aim to preserve previously learned tasks, while FedWeIT (Yoon et al., 2021) and NetTailor (Venkatesha et al., 2022) prevent interference between

irrelevant tasks. Other methods including CFed (Ma et al., 2022), FedCL (Yao & Sun, 2020), and GLFC (Dong et al., 2022) use surrogate datasets, importance weights, or class-aware techniques to distill the knowledge obtained from previous tasks. However, existing CFL methods suffer from several limitations, *e.g.*, not scalable as the number of tasks increases (Yoon et al., 2021; Venkatesha et al., 2022), requiring a surrogate dataset (Ma et al., 2022) or additional communication overhead (Yao & Sun, 2020), and not applicable to general continual setting that does not have fixed task boundaries. Our method effectively addresses these limitations, enhancing the efficacy of CFL.

3 PRELIMINARIES

We focus on finding a single classifier f (having model parameter w) that performs well on T tasks. We assume that at time slot $t \in [T]$, the classifier is only allowed to be trained for task t , where we define $[N] := \{1, \dots, N\}$ for a positive integer N . We assume the feature-label samples (x_t, y_t) for task t are drawn from an unknown distribution D_t . The optimization problem for CL at time $\tau \in [T]$ is written as

$$\min_w \sum_{t=1}^{\tau} \mathbb{E}_{(x_t, y_t) \sim D_t} [\ell(y_t, f(x_t; w))], \quad (1)$$

where ℓ is the loss function, and $f(x_t; w)$ is the output of classifier f with parameter w , for input x_t . We consider a practical scenario where we do not have enough storage to save all the data seen for the previous task ($t < \tau$); instead, we employ a replay buffer \mathcal{M} that selectively stores a subset of data. We use the buffer data as a proxy to summarize past samples and refine the model updates. We constrain the model updates in a way that the average loss for the data in buffer \mathcal{M} does not increase. Given the model $w_{\tau-1}$ trained on previous tasks, the constrained optimization problem at time $\tau \in [T]$ is represented as:

$$\begin{aligned} \min_w \mathbb{E}_{(x_\tau, y_\tau) \sim D_\tau} [\ell(y_\tau, f(x_\tau; w))] \\ \text{s.t. } \mathbb{E}_{(x_b, y_b) \sim D_b} [\ell(y_b, f(x_b; w)) \leq \ell(y_b, f(x_b; w_{\tau-1}))], \end{aligned} \quad (2)$$

where D_b is a uniform distribution over the samples in buffer \mathcal{M} , and (x_b, y_b) are sampled from this distribution D_b . The optimization problem in Eq. 2 can be reformulated for various CL methods as below. First, some methods including DER (Buzzega et al., 2020) use regularization techniques to find the model parameter w that minimizes the loss with respect to the local replay buffer \mathcal{M} as well as current samples. For a given regularization coefficient γ , the optimization problem for CL with replay buffers at time $\tau \in [T]$ is:

$$\min_w \mathbb{E}_{(x_\tau, y_\tau) \sim D_\tau} [\ell(y_\tau, f(x_\tau; w))] + \gamma \mathbb{E}_{(x_b, y_b) \sim D_b} [\ell(y_b, f(x_b; w))]. \quad (3)$$

Second, some other methods, including A-GEM (Chaudhry et al., 2019), attempt to approximately implement the constraints of Eq. 2 by considering the gradients with respect to the current/buffer data. Specifically, the constraint promotes the alignment of the gradient with respect to the current batch of data (x_τ, y_τ) and that for the buffer data (x_b, y_b) sampled from the distribution D_b . This optimization problem at time $\tau \in [T]$ is formulated as:

$$\begin{aligned} \min_w \mathbb{E}_{(x_\tau, y_\tau) \sim D_\tau} [\ell(y_\tau, f(x_\tau; w))] \\ \text{s.t. } \mathbb{E}_{(x_\tau, y_\tau) \sim D_\tau, (x_b, y_b) \sim D_b} [\langle \nabla_w [\ell(y_\tau, f(x_\tau; w))], \nabla_w [\ell(y_b, f(x_b; w))] \rangle] \geq 0 \end{aligned} \quad (4)$$

For the continual *federated* learning (CFL) setup where the data is owned by K clients, we use the superscript $k \in [K]$ to denote each client, *i.e.*, client k samples the data from D_t^k at time t and employs a local replay buffer \mathcal{M}^k . In the case of using FedAvg (McMahan et al., 2017), each round of the CFL is operated as follows. First, each client $k \in [K]$ performs multiple iterations of local updates with D_t^k with the assistance of replay buffer \mathcal{M}^k . Second, once the local training is completed, each client sends the model updates to the central server. Finally, the central server aggregates the model updates and transmits them back to clients.

4 FEDGP

We introduce a method FedGP that is compatible with various CFL techniques, significantly enhancing their performance in the CFL context. Our approach draws inspiration from A-GEM (Chaudhry

et al., 2019), which projects the gradient with respect to its own historical data. Building upon this idea, we utilize the global buffer gradient, which is the average buffer gradient across all clients, as a reference to project the local gradient. This allows us to take advantage of the collective experience of multiple clients and mitigate the risk of forgetting previously learned knowledge in FL scenarios.

Algorithm 1 FedAvg ServerUpdate with FedGP

Initialize random w^k , and set $\mathcal{M}^k = \{\}$, $g_{\text{ref}} = \text{None}$
for each task $t = 1$ **to** T **do**
 for each communication $r = 1$ **to** R **do**
 $w^k \leftarrow \text{ClientUpdate}(t, w^k, g_{\text{ref}}), \forall k$
 $w \leftarrow \text{SecAgg}(w^k)$
 $g_{\text{ref}}^k \leftarrow \text{ComputeBufferGrad}(w, \mathcal{M}^k), \forall k$
 $g_{\text{ref}} \leftarrow \text{SecAgg}(g_{\text{ref}}^k)$
 end for
end for
Return w , the final global model

Algorithm 2 ClientUpdate(t, w, g_{ref}) at client k

Input: Task index t , model w , buffer gradient g_{ref} , batch size β
Load the dataset \mathcal{D}_t^k , local buffer \mathcal{M}^k
Initialize $n = 0$ at the first task
for each batch $\{(x_i, y_i)\}_{i=1}^\beta$ in \mathcal{D}_t^k **do**
 $g = \nabla_w \left[\frac{1}{\beta} \sum_{i=1}^\beta \ell(y_i, f(x_i; w)) \right]$
 $\tilde{g} \leftarrow g - \text{proj}_{g_{\text{ref}}} g \cdot \mathbf{1}(g_{\text{ref}}^\top g \leq 0)$
 $w \leftarrow w - \alpha \tilde{g}$ for some learning rate α
 ReservoirSampling($\mathcal{M}^k, \{(x_i, y_i)\}_{i=1}^\beta, n$)
 $n \leftarrow n + \beta$
end for
Return w to server

Algorithm 3 ComputeBufferGrad(w, \mathcal{M}^k)

Input: global model w , local buffer \mathcal{M}^k
 $(x_1, y_1) \dots (x_m, y_m) \leftarrow$ random samples from \mathcal{M}^k
 $g = \frac{1}{m} \sum_{i=1}^m \nabla_w [\ell(y_i, f(x_i; w))]$
Return g to server

storing a subset of samples for previous tasks, and the data \mathcal{D}_t^k for the current task. For each batch $\{(x_i, y_i)\}_{i=1}^\beta$ in \mathcal{D}_t^k , the client computes the batch gradient g for the model w . The client then compares the direction of g with the direction of the global buffer gradient g_{ref} received from the server. When the angle between g and g_{ref} is greater than 90° , it implies that while using the direction of g as a reference for gradient descent may improve performance on the current task, but at the cost of degrading performance on previous tasks. To retain the knowledge on the previous tasks, we do the following: whenever g and g_{ref} are having a negative inner product, we project the gradient g onto the global buffer gradient (which can be considered as a reference) g_{ref} and remove this component from g , i.e., define

$$\tilde{g} = g - \frac{g^\top g_{\text{ref}}}{g_{\text{ref}}^\top g_{\text{ref}}} g_{\text{ref}} \cdot \mathbf{1}(g_{\text{ref}}^\top g \leq 0), \quad (5)$$

following the idea suggested in (Chaudhry et al., 2019). As illustrated in Fig. 2, this projection helps prevent the model updates along the direction that is harming the performance on previous tasks.

As a replay-based method, FedGP maintains a local buffer on each client, which is a memory buffer storing a subset of data sampled from old tasks. The local buffer at client k is denoted by \mathcal{M}^k . As the continuous data is loaded to the client, it keeps updating the buffer so that \mathcal{M}^k becomes a good representative of old tasks.

Algorithm 1 provides the overview of our method in CFL setup, including the process of sharing information (model and buffer gradient) between the server and each client. For each new round $r \in [R]$, the server first aggregates the local models w^k from client $k \in [K]$, getting a global model w . Afterwards, the server aggregates the local buffer gradient g_{ref}^k (the gradient computed on the global model w with respect to the local buffer \mathcal{M}^k) from client $k \in [K]$ to obtain a global buffer gradient g_{ref} . It is worth noting that the term ‘‘aggregation’’ in this context refers to the averaging of locally computed values across all clients. Such aggregation can be securely performed by the central server using secure aggregation (Bonawitz et al., 2017), which is denoted as ‘‘SecAgg’’ in Algorithm 1. Note that here we have two functions used at the client side, ClientUpdate and ComputeBufferGrad, which are given in Algorithm 2 and 3, respectively. ClientUpdate shows how client k updates its local model for task t . The client first loads the global model w and the global buffer gradient g_{ref} which are received from the server in the previous round. It also loads the local buffer \mathcal{M}^k

After gradient projection, the client updates its local model w by applying the gradient descent step with the updated gradient \tilde{g} . Finally, the client updates the contents of the buffer \mathcal{M}^k by using the reservoir sampling (Vitter, 1985) written in Algorithm 4 in the Appendix. Reservoir sampling selects a random sample of $|\mathcal{M}^k|$ elements from a local input stream, while ensuring that each element has an equal probability of being included in the sample. One of the advantages of this method is that it does not require any prior knowledge of the data stream. Once the updated local models $\{w^k\}_{k=1}^K$ are transmitted to the server, the global model w is securely updated on the server side, and transmitted to each client. Then, each client k computes the local buffer gradient (*i.e.*, the gradient of the model w with respect to the samples in the local buffer \mathcal{M}^k) as shown in Algorithm 3 `ComputeBufferGrad`.

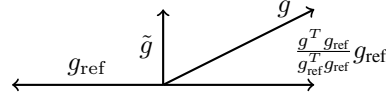


Figure 2: Illustration of the gradient projection in Eq. 5. If the angle between the gradient update g and global buffer gradient (considered as a reference) g_{ref} is larger than 90° , we project g on g_{ref} to minimize the interference and merely update along the directions of \tilde{g} that is orthogonal to g_{ref} .

After each client computes the local buffer gradient g_{ref}^k , the server allows the use of secure aggregation to combine these local buffer gradients and update the global buffer gradient g_{ref} . Secure aggregation is a well-established technique in FL that ensures the server learns nothing about individual clients’ data beyond their aggregated sum. This compatibility with secure aggregation enhances the privacy safeguards of our proposed approach, effectively minimizing the risk of data leakage from individual clients. The aforementioned process takes place between each communication and is repeated R times within each task. After traversing T tasks, the final global model w is obtained, as shown in Algorithm 1. Note that the pseudocode describes the FL+FedGP process. FedGP is designed to be compatible with various CFL techniques. Details of the combination are elaborated upon in the Appendix C.

5 EXPERIMENTS

In this section, we assess the efficacy of our method, FedGP, in combination with various CFL baselines, under non-IID data distribution across clients. To evaluate these methods, we conduct experiments on image classification tasks for benchmark datasets including rotated-MNIST (Lopez-Paz & Ranzato, 2017), permuted-MNIST (Goodfellow et al., 2013), sequential-CIFAR10, and sequential-CIFAR100 (Lopez-Paz & Ranzato, 2017) datasets, as well as a text classification task (Mehta et al., 2021) on sequential-YahooQA dataset (Zhang et al., 2015). We also explore FedGP on an object detection task on a streaming CARLA dataset (Dai et al., 2023; Dosovitskiy et al., 2017) in Appendix B. We have further explored an ablation study, which examines each element of our approach, highlighting their roles in performance enhancement. All experiments were conducted on a Linux workstation equipped with 8 NVIDIA GeForce RTX 2080Ti GPUs and averaged across five runs, each using a different seed. For further details and additional results, please refer to Appendix A.

5.1 IMAGE CLASSIFICATION

5.1.1 SETTINGS

Evaluation Datasets. We evaluate our approach on three CL scenarios: domain incremental learning (domain-IL), class incremental learning (class-IL), and task incremental learning (task-IL). For domain-IL, the data distribution of each class changes across different tasks. We use the rotated-MNIST and permuted-MNIST datasets for domain-IL, where each task rotates the training digits by a random angle or applies a random permutation. We create $T = 10$ tasks for domain-IL experiments.

For class-IL and task-IL, we partition the set of classes into disjoint subsets and assign each subset to a particular task. For instance, in our image classification experiments for class-IL and task-IL, we divide the CIFAR-100 dataset (with $C = 100$ classes) into $T = 10$ subsets, each of which contains the samples for $C/T = 10$ classes. Each task $t \in [T]$ is defined as the classification of images from each subset $t \in [T]$. The difference between class-IL and task-IL is that in the task-IL setup, we assume the task identity t is given at inference time. That is, the model f predicts among the $C/T = 10$ classes corresponding to task t . The class-IL and task-IL settings for CIFAR-10 are defined by splitting the CIFAR-10 dataset into $T = 5$ tasks, with each task having two unique classes.

In the FL setup, we assume that the data distribution is non-IID across the different clients. Once we define the data for each task, we split it among K clients in a non-IID manner. For the rotated-MNIST or permuted-MNIST dataset, each client receives samples for two MNIST digits. To create a sequential-CIFAR10 or sequential-CIFAR100 dataset, we use the Latent Dirichlet Allocation (LDA) (Hsu et al., 2019). This algorithm partitions the dataset among multiple clients by assigning samples of each class to different clients based on the probability distribution $p \sim \text{Dir}(\alpha)$, where $\alpha = 0.3$. Communication of models and buffer gradients occurs whenever all clients complete E local epochs training.

Architecture and Hyperparameters. For the rotated-MNIST and permuted-MNIST dataset, we use a simple CNN architecture (McMahan et al., 2017), and split the dataset into $K = 10$ clients. Each client performs local training for $E = 1$ epoch between communications, and we set the number of communication rounds as $R = 20$ for each task. For the sequential-CIFAR10 and sequential-CIFAR100 dataset, we use a ResNet18 architecture, and divide the dataset into $K = 10$ clients. Each client trains for $E = 5$ epochs between communications, and uses $R = 20$ rounds of communication for each task. During local training, Stochastic Gradient Descent (SGD) is employed with a learning rate of 0.01 for MNIST and 0.1 for CIFAR datasets. Unless otherwise noted, our studies used a 200 buffer size, a negligible storage concern for edge device like iPhone.

Baselines. We evaluate the performance improvement of FedGP on three types of baselines: 1) *FL*, the basic FedAvg which trains only on the current task without considering performance on previous tasks; 2) *FL+CL*, which is FedAvg (FL) with continual learning solutions applied to clients; and 3) *CFL*, which represents the existing Continual Federated Learning methods.

CL methods include A-GEM (Chaudhry et al., 2019), which aligns model gradients for buffer and incoming data; GPM (Saha et al., 2021), using the network representation/activations approximated by top singular vectors as the old tasks’ reference vector; DER (Buzzega et al., 2020), utilizing network output logits for past experience distillation; iCaRL (Rebuffi et al., 2017), which adds current task nearest-mean-of-exemplars to a memory buffer via herding and counters representation deterioration with a self-distillation loss term; and L2P (Wang et al., 2022), a state-of-the-art approach that instructs pre-trained models to sequentially learn tasks using prompt pool memory spaces.

CFL methods we tested are FedCurv (Shoham et al., 2019), which avoids updating past task-critical weights; FedProx (Li et al., 2020), introducing a proximal weight for global model alignment; CFed (Ma et al., 2022), using surrogate dataset-based knowledge distillation; and GLFC (Dong et al., 2022), a tripartite method to counteract forgetting: 1) clients retain old task data and incorporate it with a normalization factor during new task training, 2) clients save the prior model, compute the KL divergence loss between the new and old model outputs (from the last layer), and 3) an additional proxy server is used to gather perturbed gradients for client sample generation. FOT (Bakman et al., 2023) also projects the gradients on the subspace specified by previous tasks. Details are in Appendix A.6. As for FedWeIT (Yoon et al., 2021), we believe it might not serve as a suitable benchmark given its focus on personalized federated learning (FL) without a global model accuracy to contrast with.

Note that CFed, GLFC, GPM, and iCaRL require task boundaries during training. They exploit task changes to snapshot the network, with iCaRL further relying on these for memory buffer updates. For specific parameters and implementation unique to each method, please refer to Appendix A.9.

Performance Metrics. We assess the performance of the global model on the test dataset, which is a union of the test data for all previous tasks. The average accuracy (measured after training on task t) is denoted as $\text{Acc}_t = \frac{1}{t} \sum_{i=1}^t a_{t,i}$, where $a_{t,i}$ is evaluated on task i after training up to task t . Additionally, we measure a performance metric called *forgetting* in Appendix A.2, which is defined as the difference between the best accuracy obtained throughout the training and the current accuracy (Chaudhry et al., 2018). This metric measures the model’s ability to retain knowledge of previous tasks while learning new ones. The forgetting at task t is defined as: $\text{Fgt}_t = \frac{1}{t-1} \sum_{i=1}^{t-1} \max_{j=1, \dots, t-1} (a_{j,i} - a_{t,i})$. We also computed the Backward transfer (BWT) and Forward transfer (FWT) metrics (Lopez-Paz & Ranzato, 2017). See Appendix A.7 for details.

Table 1: Average accuracy Acc_T (%) on standard benchmark datasets. ‘-’ indicates experiments we were unable to run, because of compatibility issues (e.g. GLFC and iCaRL in Domain-IL) or the absence of surrogate (e.g. CFeD on MNIST). The results, averaged over 5 random seeds, demonstrate the benefits of our proposed method in combination with all baselines. A buffer size of 200 is utilized whenever methods require it. Note that FL+L2P needs additional pretrained ViT.

Method	rotated-MNIST (Domain-IL)		sequential-CIFAR10 (Class-IL)		sequential-CIFAR10 (Task-IL)	
	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP
FL (McMahan et al., 2017)	68.02 \pm 3.1	79.46 \pm 4.1 (\uparrow 11.44)	17.44 \pm 1.3	18.02 \pm 0.6 (\uparrow 0.58)	70.58 \pm 4.0	80.83 \pm 2.0 (\uparrow 10.25)
FL+A-GEM (Chaudhry et al., 2019)	68.34 \pm 5.6	74.74 \pm 2.3 (\uparrow 6.40)	17.82 \pm 0.9	19.44 \pm 0.9 (\uparrow 1.62)	77.14 \pm 3.1	83.16 \pm 1.6 (\uparrow 6.02)
FL+GPM (Saha et al., 2021)	74.42 \pm 6.4	-	17.59 \pm 0.4	-	74.50 \pm 3.6	-
FL+DER (Buzzega et al., 2020)	57.73 \pm 3.6	81.33 \pm 3.3 (\uparrow 23.60)	18.44 \pm 3.7	30.94 \pm 3.8 (\uparrow 12.50)	69.34 \pm 3.2	77.99 \pm 0.8 (\uparrow 8.65)
FL+iCaRL (Rebuffi et al., 2017)	-	-	28.54 \pm 3.8	33.92 \pm 3.0 (\uparrow 5.38)	80.85 \pm 2.9	80.09 \pm 4.1 (\downarrow 0.76)
FL+L2P (Wang et al., 2022)	80.90 \pm 3.3	85.05 \pm 0.7 (\uparrow 4.15)	28.61 \pm 1.0	81.86 \pm 7.2 (\uparrow 53.25)	98.49 \pm 0.1	98.63 \pm 0.3 (\uparrow 0.14)
FedCurv (Shoham et al., 2019)	68.21 \pm 2.6	80.53 \pm 4.3 (\uparrow 12.32)	17.36 \pm 0.7	17.86 \pm 0.5 (\uparrow 0.50)	67.77 \pm 1.4	81.28 \pm 1.1 (\uparrow 13.51)
FedProx (Li et al., 2020)	67.79 \pm 3.2	78.74 \pm 4.1 (\uparrow 10.95)	16.67 \pm 2.7	17.97 \pm 0.8 (\uparrow 1.30)	69.57 \pm 6.5	81.23 \pm 1.3 (\uparrow 11.66)
CFeD (Ma et al., 2022)	-	-	16.30 \pm 4.6	24.07 \pm 8.5 (\uparrow 7.77)	77.35 \pm 4.6	79.30 \pm 5.7 (\uparrow 1.95)
GLFC (Dong et al., 2022)	-	-	41.42 \pm 1.3	41.61 \pm 1.3 (\uparrow 0.19)	81.84 \pm 2.1	82.87 \pm 1.0 (\uparrow 1.03)
	permuted-MNIST (Domain-IL)		sequential-CIFAR100 (Class-IL)		sequential-CIFAR100 (Task-IL)	
FL	25.92 \pm 2.1	34.23 \pm 2.7 (\uparrow 8.31)	8.76 \pm 0.1	17.08 \pm 1.8 (\uparrow 8.32)	47.74 \pm 1.2	74.71 \pm 0.9 (\uparrow 26.97)
FL+A-GEM	33.43 \pm 1.4	39.09 \pm 3.5 (\uparrow 5.66)	8.90 \pm 0.1	19.53 \pm 1.3 (\uparrow 10.63)	63.84 \pm 0.8	74.84 \pm 0.5 (\uparrow 11.00)
FL+GPM	31.92 \pm 3.4	-	8.18 \pm 0.1	-	54.48 \pm 1.4	-
FL+DER	19.79 \pm 1.7	38.81 \pm 2.0 (\uparrow 19.02)	13.32 \pm 1.6	22.96 \pm 3.6 (\uparrow 9.64)	57.71 \pm 1.2	65.57 \pm 1.9 (\uparrow 7.86)
FL+iCaRL	-	-	21.76 \pm 1.1	27.44 \pm 1.2 (\uparrow 5.68)	69.91 \pm 0.7	72.83 \pm 0.5 (\uparrow 2.92)
FL+L2P	66.98 \pm 4.6	69.15 \pm 3.1 (\uparrow 2.17)	23.12 \pm 1.7	46.16 \pm 0.4 (\uparrow 23.04)	94.46 \pm 0.4	94.91 \pm 0.2 (\uparrow 0.45)
FedCurv	26.00 \pm 2.4	35.21 \pm 5.1 (\uparrow 9.21)	8.92 \pm 0.1	16.67 \pm 0.9 (\uparrow 7.76)	49.14 \pm 1.6	74.64 \pm 0.7 (\uparrow 25.49)
FedProx	25.92 \pm 2.5	35.60 \pm 4.7 (\uparrow 9.68)	8.75 \pm 0.2	16.92 \pm 1.4 (\uparrow 8.17)	47.05 \pm 3.2	73.95 \pm 0.8 (\uparrow 26.89)
CFeD	-	-	13.76 \pm 1.2	26.66 \pm 0.3 (\uparrow 12.9)	51.41 \pm 1.0	72.20 \pm 0.9 (\uparrow 20.79)
GLFC	-	-	13.18 \pm 0.4	13.47 \pm 0.7 (\uparrow 0.29)	49.78 \pm 0.8	49.20 \pm 1.2 (\downarrow 0.58)

5.1.2 RESULTS

Table 1 presents the average accuracy Acc_T of various methods on image classification benchmark datasets measured upon completion of the final task T . For each setting, we compare the performance of an existing method with/without FedGP. We observe that the proposed methods (represented by “w/ FedGP”) nearly always improves the base methods (“w/o FedGP”) across the different datasets and scenarios, as seen from the upward arrows indicating performance improvements. Additional results are obtained for forgetting performance Fgt_T given in Table 11 in Appendix A.2. Moreover, the performance of FedGP is analyzed progressively across tasks in Appendix A.1.

Remarkably, even a simple integration of the basic baseline, FL, with FedGP surpassed the performance of most FL+CL and CFL baselines. For instance, in the sequential-CIFAR100 experiment, FL with FedGP (17.08% class-IL, 74.71% task-IL) outperformed a majority of the baselines. Specifically, it exceeds the performance of the two advanced CFL baselines: GLFC (13.18% class-IL, 49.78% task-IL) and CFeD (13.76% class-IL, 51.41% task-IL). This underscores the substantial capability of our method in the CFL setting. Importantly, FedGP can achieve competitive performance even without utilizing information about task boundaries, unlike CFeD, GLFC, GPM, and iCaRL.

We also note that the FL+L2P method consistently exhibited the highest accuracy, largely due to the utilization of a pretrained Vision Transformer (ViT) (Dosovitskiy et al., 2020; Zhang et al., 2022), which helps mitigate the catastrophic forgetting. This is why we wrote the numbers in gray with a caveat in the caption. Yet, our approach still managed to achieve significant performance augmentation on top of it.

Effect of Buffer size. Table 2 reports the performances of baseline CL methods (A-GEM and DER) with/without FedGP for different buffer sizes, ranging from 200 to 5120. For all different datasets and all IL settings, increasing the buffer size further improves the advantage of applying FedGP, by providing more data for replay and mitigating forgetting. However, a finite buffer cannot maintain the entire history of data. In Fig. 3 we reported the effect of buffer size on the accuracy of old tasks. At the end of each task, we measured the accuracy of the trained model with respect to the test data for task 1. We tested on sequential-CIFAR100 dataset, and considered task incremental learning (task-IL) setup. One can observe that when the buffer size B is small, the accuracy drops as the model is trained on new tasks. On the other hand, when $B \geq 100$, the task-IL accuracy for task 1 is

Table 2: Impact of the buffer size on Acc_T (%)

Buffer Size	Method	rotated-MNIST (<i>Domain-IL</i>)		sequential-CIFAR100 (<i>Class-IL</i>)		sequential-CIFAR100 (<i>Task-IL</i>)	
		w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP
200	FL+A-GEM	68.34 \pm 5.6	74.74 \pm 2.3 (\uparrow 6.40)	8.90 \pm 0.1	19.53 \pm 1.3 (\uparrow 10.63)	63.84 \pm 0.8	74.84 \pm 0.5 (\uparrow 11.00)
500		70.18 \pm 8.7	78.74 \pm 3.2 (\uparrow 8.56)	8.87 \pm 0.1	25.89 \pm 0.9 (\uparrow 17.02)	64.38 \pm 1.4	79.35 \pm 0.5 (\uparrow 14.97)
5120		69.97 \pm 3.2	79.17 \pm 4.3 (\uparrow 9.20)	8.85 \pm 0.1	33.30 \pm 2.5 (\uparrow 24.45)	64.99 \pm 1.5	84.52 \pm 0.3 (\uparrow 19.53)
200	FL+DER	57.73 \pm 3.6	87.13 \pm 1.1 (\uparrow 29.40)	13.32 \pm 1.6	22.96 \pm 3.6 (\uparrow 9.64)	57.71 \pm 1.2	65.57 \pm 1.9 (\uparrow 7.86)
500		60.00 \pm 7.2	88.83 \pm 1.6 (\uparrow 28.83)	15.44 \pm 1.5	34.87 \pm 1.7 (\uparrow 19.43)	60.79 \pm 1.2	73.53 \pm 1.1 (\uparrow 12.74)
5120		58.63 \pm 3.9	89.46 \pm 1.2 (\uparrow 30.83)	18.89 \pm 1.0	45.76 \pm 3.8 (\uparrow 26.87)	62.77 \pm 1.5	83.41 \pm 1.3 (\uparrow 20.64)

Table 3: Effect of communication on accuracy (%) performance, with values in brackets indicating differences from the FL baseline.

Method	R-MNIST	P-MNIST	S-CIFAR10		S-CIFAR100	
	<i>Domain-IL</i>	<i>Domain-IL</i>	<i>Class-IL</i>	<i>Task-IL</i>	<i>Class-IL</i>	<i>Task-IL</i>
FL	68.02 \pm 3.1	27.49 \pm 2.0	17.44 \pm 1.3	70.58 \pm 4.0	8.76 \pm 0.1	47.74 \pm 1.2
FL w/ FedGP (2 \times comm overhead)	79.46 \pm 4.1 (\uparrow 11.44)	35.91 \pm 4.0 (\uparrow 8.42)	18.02 \pm 0.6 (\uparrow 0.58)	80.83 \pm 2.0 (\uparrow 10.25)	17.08 \pm 1.8 (\uparrow 8.32)	74.71 \pm 0.9 (\uparrow 26.97)
FL w/ FedGP (equalized comm overhead)	75.63 \pm 3.9 (\uparrow 7.61)	34.96 \pm 3.2 (\uparrow 7.47)	16.65 \pm 1.0 (\downarrow 0.79)	78.79 \pm 2.8 (\uparrow 8.21)	13.62 \pm 0.6 (\uparrow 4.86)	73.96 \pm 0.4 (\uparrow 26.22)
FL w/ FedGP (0.5 \times comm overhead)	76.05 \pm 4.0 (\uparrow 8.03)	29.75 \pm 4.6 (\uparrow 2.26)	14.30 \pm 1.3 (\downarrow 3.14)	66.90 \pm 3.6 (\downarrow 3.68)	13.09 \pm 0.5 (\uparrow 4.33)	69.96 \pm 0.6 (\uparrow 22.22)
FL w/ FedGP (0.2 \times comm overhead)	70.59 \pm 4.7 (\uparrow 2.57)	15.51 \pm 2.7 (\downarrow 11.98)	13.37 \pm 2.6 (\downarrow 4.07)	59.75 \pm 6.4 (\downarrow 10.83)	13.59 \pm 0.9 (\uparrow 4.83)	59.31 \pm 1.6 (\uparrow 11.57)

maintained throughout the process. Note that training with our default setting $B = 200$ does not hurt the accuracy for task 1 throughout the continual learning process.

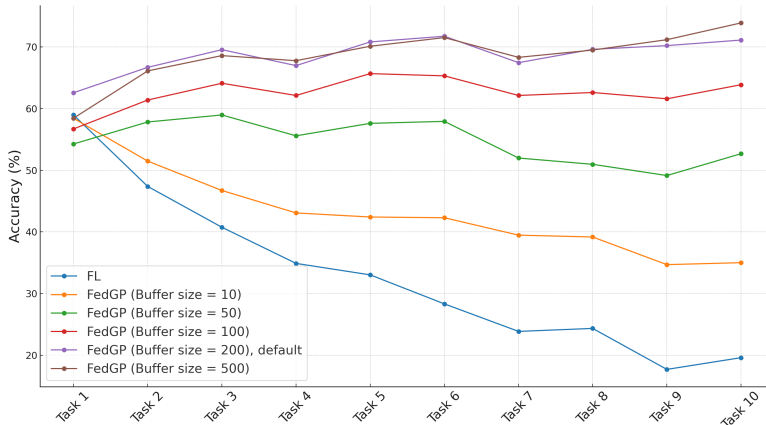


Figure 3: Accuracy (%) for Task 1, under the Task-IL setting on sequential-CIFAR100. We tested on FedGP with different buffer sizes.

We are assuming that every client has the same buffer size. If the buffer sizes are not equal during model training, clients with bigger buffers might add more diverse data, which could make the model biased. A possible solution is to use a reweighting algorithm, which we plan to explore in the future.

Effect of communication frequency. Compared with baseline methods, FedGP has extra communication overhead for transmitting the buffer gradients from each client to the server. This means that the required amount of communication is doubled for FedGP. We consider a variant of FedGP which updates the model and buffer gradient less frequently (i.e., reduce the communication rounds for each task), which has reduced communication than the vanilla FedGP. Table 3 reports the performance for different datasets, when the communication overhead is set to 2x, 1x, 0.5x and 0.2x. First, in most cases, FedGP with equalized (1x) communication overhead is outperforming FL. In addition, for most of tested datasets including R-MNIST, P-MNIST and S-CIFAR100, FedGP outperforms FL with at most 0.5x communication overhead. This means that FedGP enjoys a higher performance with less communication, in the standard benchmark datasets for continual federated learning.

Effect of computation overhead. Computation overhead is also an important aspect to consider and we have conducted experiment on the actual wall-clock time measurements. Taking a CIFAR100 experiment as an example, the running time for 200 epochs for FedAvg on our device is 4068.97s. When FedGP, which is built on top of FedAvg, was used, it ran for an additional 293.26s. This indicates that it ran 7.2% longer over the same 200 epochs. The time consumed by FedGP can be divided into two parts: (i) computing the global reference gradient after each FedAvg, and (ii) projecting the gradient. In the above experiment, the reference gradient computation was done 200 times, taking a total of 49.07s, and the gradient projection was performed on 109,471 batches, which is 68.38% of the total batches, taking a total of 244.19s. Overall, FedGP is computationally efficient and affordable for edge devices.

Table 4: Acc_T (%) for asynchronous task boundaries on the sequential-CIFAR100 dataset.

Method	<i>Class-IL</i>	<i>Task-IL</i>
FL	16.22 \pm 1.2	59.04 \pm 1.7
FL+A-GEM	16.92 \pm 1.0	69.41 \pm 1.3
FL+A-GEM+FedGP	30.74 \pm 1.5	77.70 \pm 0.4
FL+DER	31.95 \pm 2.6	68.28 \pm 1.5
FL+DER+FedGP	36.29 \pm 1.0	72.02 \pm 0.7

Effect of the privacy leakage. Regarding the data leakage issue, we empirically evaluate the privacy leakage of the FedGP framework and various other CFL baselines. For this purpose, we are employing the Mutual Information Neural Estimator (MINE) (Belghazi et al., 2018). This approach allows us to estimate the mutual information between individual user model updates and the aggregated model update across all users at each global training round. This methodology closely follows the precedent work (Elkordy et al., 2022). Results will be added soon.

Asynchronous task boundaries. In our previous experiments, we assumed synchronous task boundaries where clients finish tasks at the same time. However, in many real-world scenarios, different clients finish each task asynchronously. Motivated by this practical setting, we conducted experiments in an asynchronous task boundary setting on sequential-CIFAR100. For every $R = 20$ communications, instead of traversing all data allocated for the current task, each client traverses exactly 500 samples allocated to it, irrespective of whether these samples come from the same task. Consequently, some clients might finish a task faster than others and move on to the next task. Thus, during each global communication, clients could be working on different tasks. This setup more closely aligns with our general continual learning settings, when the task boundary is unknown. Table 4 shows the accuracy of each method averaged over T tasks after finishing all training, under the asynchronous setting. Similar to the synchronous case, FedGP improves the accuracy of baseline methods including A-GEM and DER. Notably, we have a better performance in the asynchronous setting (see Table 4) compared with the synchronous setting (see Table 1). This might be because, in the asynchronous setting, some clients receive new tasks earlier than others, which allows the model to be exposed to more diverse data for each round, thus reducing the forgetting effect.

Effect of the number of tasks. We have conducted experiments with different number of tasks for each dataset. For CIFAR100, we experimented with task numbers 5 and 10, while for CIFAR10 we tested with task numbers 2 and 5. Our results in Table 5 consistently demonstrate that the FedGP algorithm provides a significant improvement in performance across all these different task numbers. An interesting observation is that as the number of tasks increases, FedGP have better performance improvement to baseline. This is because a higher number of tasks increases the likelihood of data distribution shifts and therefore the problem of catastrophic forgetting becomes more prominent. As such, FedGP, designed to handle this issue, has more opportunities to improve the learning process in such scenarios. This might also partly explain why, in Table 1, FedGP shows a generally higher improvement over the baselines on the sequential-CIFAR100 dataset compared to the sequential-CIFAR10.

Effect of the number of users. We also conducted experiments to assess scalability by increasing the client count to $K = 20$. Table 6 shows the results for $K = 20$ clients. This results demonstrates that FedGP consistently improves the performance of baselines, across different number of clients.

Additionally, we evaluated a real-world scenario where only a random subset of clients participates in training during each round. Detailed information and results are available in Appendix A.4

Table 5: Average accuracy Acc_T (%) across various task numbers.

(# of Task, # of Classes per Task)	sequential-CIFAR10 (<i>Class-IL</i>)		sequential-CIFAR10 (<i>Task-IL</i>)	
	FL	FL w/ FedGP	FL	FL w/ FedGP
(2, 5)	43.53 \pm 0.8	44.05 \pm 0.8 (\uparrow0.52)	75.54 \pm 0.6	77.52 \pm 0.8 (\uparrow1.98)
(5, 2)	17.44 \pm 1.3	18.02 \pm 0.6 (\uparrow0.6)	70.58 \pm 4.0	80.83 \pm 2.0 (\uparrow10.25)
	sequential-CIFAR100 (<i>Class-IL</i>)		sequential-CIFAR100 (<i>Task-IL</i>)	
	FL	FL w/ FedGP	FL	FL w/ FedGP
(5, 20)	16.49 \pm 0.3	22.71 \pm 0.9 (\uparrow6.22)	50.60 \pm 0.9	69.41 \pm 0.8 (\uparrow18.81)
(10, 10)	8.76 \pm 0.1	17.08 \pm 1.8 (\uparrow8.32)	47.74 \pm 1.2	74.71 \pm 0.9 (\uparrow26.97)

Table 6: The Acc_T (%) performance measured when we have $K = 20$ users. Similar to the results for $K = 10$ in Table 1, our method improves the performance of baselines.

Method	rotated-MNIST (<i>Domain-IL</i>)		sequential-CIFAR10 (<i>Class-IL</i>)		sequential-CIFAR10 (<i>Task-IL</i>)	
	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP
FL	62.45 \pm 8.5	76.01 \pm 4.6 (\uparrow13.56)	16.44 \pm 1.4	15.82 \pm 1.7 (\downarrow0.62)	68.18 \pm 5.3	73.45 \pm 4.3 (\uparrow5.27)
FedCurv	62.57 \pm 8.3	76.46 \pm 4.1 (\uparrow13.89)	17.31 \pm 0.6	14.64 \pm 3.1 (\downarrow2.67)	67.33 \pm 3.3	70.31 \pm 3.7 (\uparrow2.98)
FedProx	62.14 \pm 8.6	75.84 \pm 4.4 (\uparrow13.70)	16.37 \pm 1.1	16.15 \pm 1.3 (\downarrow0.22)	66.24 \pm 1.4	74.79 \pm 3.9 (\uparrow8.55)
FL+A-GEM	67.66 \pm 8.0	78.10 \pm 3.6 (\uparrow10.44)	16.15 \pm 1.9	17.36 \pm 0.8 (\uparrow1.21)	72.39 \pm 3.4	80.61 \pm 2.6 (\uparrow8.22)
FL+DER	57.33 \pm 3.2	87.84 \pm 1.5 (\uparrow30.51)	17.13 \pm 2.3	19.18 \pm 3.7 (\uparrow2.05)	70.82 \pm 1.9	77.04 \pm 2.5 (\uparrow6.22)
Method	permuted-MNIST (<i>Domain-IL</i>)		sequential-CIFAR100 (<i>Class-IL</i>)		sequential-CIFAR100 (<i>Task-IL</i>)	
	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP
FL	20.26 \pm 1.6	20.67 \pm 4.7 (\uparrow0.41)	8.61 \pm 0.1	17.47 \pm 1.1 (\uparrow8.86)	50.00 \pm 1.6	76.29 \pm 0.8 (\uparrow26.29)
FedCurv	20.25 \pm 1.9	23.30 \pm 5.7 (\uparrow3.05)	8.93 \pm 0.0	19.42 \pm 1.1 (\uparrow10.49)	49.83 \pm 1.4	79.58 \pm 0.6 (\uparrow29.75)
FedProx	20.19 \pm 1.4	23.78 \pm 5.2 (\uparrow3.59)	8.88 \pm 0.1	18.86 \pm 1.0 (\uparrow9.98)	50.86 \pm 1.2	78.19 \pm 0.9 (\uparrow27.33)
FL+A-GEM	24.43 \pm 2.1	23.29 \pm 3.8 (\downarrow1.14)	8.62 \pm 0.1	19.58 \pm 1.2 (\uparrow10.96)	63.02 \pm 0.6	76.23 \pm 0.6 (\uparrow13.21)
FL+DER	17.89 \pm 1.3	46.17 \pm 3.0 (\uparrow28.28)	11.53 \pm 0.5	26.64 \pm 2.8 (\uparrow15.11)	57.00 \pm 1.4	69.42 \pm 1.0 (\uparrow12.42)

5.2 TEXT CLASSIFICATION

In addition to image classification, we also extended the evaluation of our method on text classification task (Mehta et al., 2021). For this purpose, we utilized the YahooQA (Zhang et al., 2015) dataset which comprises texts (questions and answers), and user-generated labels representing 10 different topics. Similar to the approach taken with the CIFAR10 dataset, we partitioned the YahooQA dataset into 5 tasks, where each task consisted of two distinct classes. Within each task, we used LDA to partition data across 10 clients in a non-IID manner. To conduct the experiment, we employed a pretrained DistilBERT (Sanh et al., 2019) with linear classification layer. We freeze the DistilBERT model and only fine-tune the additional linear layer. The results of this experiment can be found in Table 7. We can observe that FedGP consistently enhances the accuracy (Acc_T) over baselines, particularly in class-IL scenarios.

Table 7: Average classification accuracy Acc_T (%) on split-YahooQA dataset.

Method	sequential-YahooQA (<i>Class-IL</i>)		sequential-YahooQA (<i>Task-IL</i>)	
	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP
FL	17.86 \pm 0.6	30.67 \pm 4.4 (\uparrow12.81)	80.87 \pm 1.2	88.04 \pm 1.4 (\uparrow7.17)
FL+A-GEM	20.86 \pm 0.3	47.02 \pm 1.9 (\uparrow26.16)	87.29 \pm 1.3	90.20 \pm 0.2 (\uparrow2.91)
FL+DER	43.64 \pm 2.1	54.28 \pm 1.3 (\uparrow10.64)	89.57 \pm 0.2	90.48 \pm 0.2 (\uparrow0.91)

5.3 ABLATIONS ON OUR ALGORITHM DESIGN

We have conducted experiments for ablations on our algorithm design, which is decomposed into (1) gradient manipulation algorithm, and (2) buffer updating algorithm. Regarding the gradient manipulation algorithm, we tested on 8 different methods that use the reference gradient to manipulate the gradient. Regarding the buffer updating algorithm, we compared three different methods.

5.3.1 GRADIENT MANIPULATION

First, we considered different ways of manipulating the gradient g , given the reference gradient g_{ref} . In the below table, we compared four different methods of updating g :

- Average: define $g \leftarrow (g + g_{\text{ref}})/2$.
- Rotate: rotate g towards g_{ref} while keeping the magnitude.
- Project: project g to the space that is orthogonal to g_{ref} .
- Project & Scale: apply Project and then scale up the vector such that the magnitude is identical to the original g

Recall that our FedGP applies Project method only when the angle between g and g_{ref} is larger than 90 degree, i.e., when the reference gradient g_{ref} (measured for the previous tasks) and the gradient g (measured for the current task) conflicts to each other. Our intuition for such choice is, it is better to manipulate g if the direction favorable for current task is conflicting with the direction favorable for previous tasks. To support that this choice is meaningful, we compared two ways of deciding when we manipulate the gradients, denoted below:

- (> 90): update the gradient g only when the $\text{angle}(g, g_{\text{ref}}) > 90$
- (Always): update the gradient g always

We compared the performances of above choices in Table 8, for sequential-CIFAR100 dataset. One can confirm that our FedGP (denoted by Project (> 90) in the table) far outperforms all other combinations, showing that our design (doing projection for conflicting case only) is the right choice. If we check each component (Project and (> 90)) independently, one can check that choosing Project outperforms Average, Rotate and Project & Scale in most cases, and choosing (> 90) outperforms Always in all cases.

Table 8: Effect of different gradient manipulation methods on the accuracy (%) of FedGP, tested on sequential-CIFAR100

Method	Class-IL	Task-IL
FL	8.76 \pm 0.1	47.74 \pm 1.2
Average (Always)	7.26 \pm 1.95	35.96 \pm 3.23
Average (> 90)	7.79 \pm 0.65	36.57 \pm 1.55
Rotate (Always)	7.59 \pm 0.89	36.15 \pm 2.83
Rotate (> 90)	8.41 \pm 0.78	38.97 \pm 1.83
Project & Scale (Always)	8.77 \pm 0.09	32.96 \pm 1.10
Project & Scale (> 90)	12.30 \pm 0.65	73.61 \pm 0.75
Project (Always)	8.90 \pm 0.08	34.00 \pm 1.98
Project (> 90), ours	17.08\pm1.8	74.71\pm0.9

We also tested whether doing the projection is helpful in all cases when $\text{angle}(g, g_{\text{ref}}) > 90$. We considered applying the projection for $p\%$ of the cases having $\text{angle}(g, g_{\text{ref}}) > 90$, for $p = 10, 50, 80$ and 100. Note that $p = 100\%$ case reduces to our FedGP.

Table 9 shows the effect of projection rate $p\%$ on the accuracy, tested on sequential-CIFAR100 dataset. In both class-IL and task-IL settings, increasing p always improves the accuracy of the FedGP method. This supports that our way of projection is suitable for the continual federated learning setup.

5.3.2 BUFFER UPDATING

In Table 10, we compared three different buffer updating algorithms:

- Sliding Window Sampling: replaces the earliest data point in the buffer when new data arrives

Table 9: Effect of projection rate $p\%$ on the accuracy (%) of FedGP, tested on sequential-CIFAR100

Method	Class-IL	Task-IL
FL, $p = 0\%$	$8.76^{\pm 0.1}$	$47.74^{\pm 1.2}$
FedGP, $p = 10\%$	$8.82^{\pm 0.07}$	$54.90^{\pm 1.61}$
FedGP, $p = 50\%$	$8.91^{\pm 0.07}$	$67.89^{\pm 0.67}$
FedGP, $p = 80\%$	$10.36^{\pm 0.42}$	$72.73^{\pm 0.74}$
FedGP, $p = 100\%$ (ours)	$17.08^{\pm 1.8}$	$74.71^{\pm 0.9}$

- Random Sampling: randomly replaces a data point in the buffer with incoming new data
- Reservoir Sampling (Ours): given n (the number of observed samples up to now) and $|\mathcal{M}^k|$ (the buffer size), we do the following
 - when $n \leq |\mathcal{M}^k|$, we put the current sample in the buffer
 - when $n > |\mathcal{M}^k|$, with probability $\frac{|\mathcal{M}^k|}{n} < 1$, we replace a sample in the buffer with the current sample

Note that when the number of incoming data is n , those n samples have the same probability of getting into the buffer, for the Reservoir Sampling method used in our paper. Thus, when Reservoir Sampling is used, the buffer contains approximately equal number of samples for each task (when each task has the same number of samples), throughout the continual learning process. This is the underlying intuition why we choose such buffer updating algorithm. To support this claim, we report the performance of different sampling methods in Table 10. We can confirm that our sampling method is outperforming other sampling methods.

We present the pseudocode for the `ReservoirSampling` algorithm in Algorithm 4. We now prove that for reservoir sampling, the probability of a sample contained in the buffer is $\frac{|\mathcal{M}^k|}{n}$. We prove this by induction; suppose this statement holds when $n - 1$ samples are observed, and we show that this holds when one more sample is observed. Note that the probability of a sample contained in the buffer can be computed as $\frac{|\mathcal{M}^k|}{n-1} * (1 - \frac{|\mathcal{M}^k|}{n} * \frac{1}{|\mathcal{M}^k|}) = \frac{|\mathcal{M}^k|}{n}$, where

- $\frac{|\mathcal{M}^k|}{n-1}$ is the probability of the sample initially contained in the buffer;
- $(1 - \frac{|\mathcal{M}^k|}{n} * \frac{1}{|\mathcal{M}^k|})$ is the probability of a sample not being kicked out of the buffer;
- $\frac{|\mathcal{M}^k|}{n} * \frac{1}{|\mathcal{M}^k|}$ is the probability of a sample being kicked out of the buffer

Algorithm 4 `ReservoirSampling`($\mathcal{M}^k, (x, y), n$) (Vitter, 1985) at client k

Input: local buffer \mathcal{M}^k , incoming data (x, y) and the number of observed samples n
if $n \leq |\mathcal{M}^k|$ **then**
 Add data (x, y) into local buffer \mathcal{M}^k
else
 $i \leftarrow \text{Uniform}\{1, 2, \dots, n\}$
if $i \leq |\mathcal{M}^k|$ **then**
 $\mathcal{M}^k[i] \leftarrow (x, y)$
end if
end if
 Return \mathcal{M}^k , the updated local buffer

Fig. 4 compares the sample distribution across different tasks, for different buffer updating algorithms (uniform-random sampling, reservoir sampling and sliding window sampling). It is important to note that the depicted distributions are representative under the condition where the buffer size is equal to the number of samples per task. This clearly shows the difference of buffer updating algorithms, which guided to the different performances reported in Table 10.

Table 10: Effect of different buffer updating algorithms on the accuracy (%) of FedGP, tested on sequential-CIFAR100

Method	Class-IL	Task-IL
FL	8.76 ± 0.1	47.74 ± 1.2
Sliding Window Sampling	8.82 ± 0.15	46.16 ± 2.38
Random Sampling	9.72 ± 0.10	54.82 ± 1.58
Reservoir Sampling (ours)	17.08 ± 1.8	74.71 ± 0.9

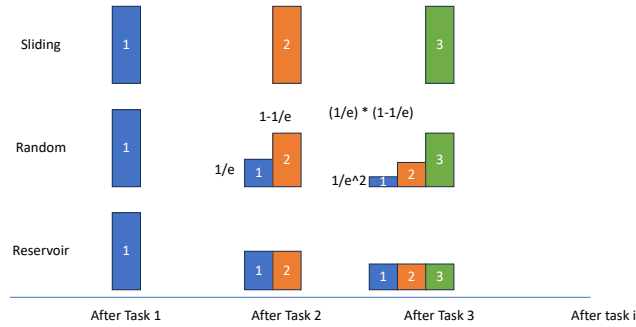


Figure 4: The buffer sample distribution across different tasks for different buffer updating algorithms (uniform-random, sliding window, our reservoir). This is illustrated on the scenario where the buffer size is equal to the number of samples per task.

All in all, the reservoir sampling method used in FedGP allows us to have balanced sample distribution across different tasks, thus allowing us to mitigate the catastrophic forgetting and to improve the accuracy in the continual federated learning setting.

6 CONCLUSION

In this paper, we present FedGP, a novel method of using buffer data for mitigating the catastrophic forgetting issues in CFL. Specifically, we use the gradient projection method to prevent model updates that harm the performance on previous tasks. Our empirical results on benchmark datasets (rotated-MNIST, permuted-MNIST, sequential-CIFAR10 and sequential-CIFAR100) and on a text classification dataset show that FedGP improves the performance of existing CL and CFL methods.

ACKNOWLEDGEMENTS

This work is funded by Intel/NSF joint grant CNS-2003129. Suman Banerjee and Shenghong Dai were also supported in part through the following US National Science Foundation(NSF) grants: CNS-1838733, CNS-1719336, CNS-1647152, and CNS-1629833, and an award from the US Department of Commerce with award number 70NANB21H043. Kangwook Lee was also supported by NSF Award DMS-2023239, and the University of Wisconsin-Madison Office of the Vice Chancellor for Research and Graduate Education with funding from the Wisconsin Alumni Research Foundation. Yicong Chen was supported by the 2023 Wisconsin Hilldale Undergraduate Research Fellowship.

REFERENCES

- Rahaf Aljundi, Punarjay Chakravarty, and Tinne Tuytelaars. Expert Gate: Lifelong learning with a network of experts. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3366–3375, 2017.
- Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 139–154, 2018.
- Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. *Advances in neural information processing systems*, 32, 2019.
- Yavuz Faruk Bakman, Duygu Nur Yaldiz, Yahya H Ezzeldin, and Salman Avestimehr. Federated orthogonal training: Mitigating global catastrophic forgetting in continual federated learning. *arXiv preprint arXiv:2309.01289*, 2023.
- Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. Mutual information neural estimation. In *International conference on machine learning*, pp. 531–540. PMLR, 2018.
- Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for privacy-preserving machine learning. In *proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1175–1191, 2017.
- Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. *Advances in neural information processing systems*, 33:15920–15930, 2020.
- Hyuntak Cha, Jaeho Lee, and Jinwoo Shin. Co²L: Contrastive continual learning. In *Proceedings of the IEEE/CVF International conference on computer vision*, pp. 9516–9525, 2021.
- Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 532–547, 2018.
- Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with A-GEM. In *ICLR*, 2019.
- Zhiyuan Chen and Bing Liu. Lifelong machine learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 12(3):1–207, 2018.
- Shenghong Dai, SM Iftekhharul Alam, Ravikumar Balakrishnan, Kangwook Lee, Suman Banerjee, and Nageen Himayat. Online federated learning based object detection across autonomous vehicles in a virtual world. In *2023 IEEE 20th Consumer Communications & Networking Conference (CCNC)*, pp. 919–920. IEEE, 2023.
- 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*, pp. 10164–10173, 2022.
- Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. CARLA: An open urban driving simulator. In *Conference on robot learning*, pp. 1–16. PMLR, 2017.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Ahmed Roushdy Elkordy, Jiang Zhang, Yahya H Ezzeldin, Konstantinos Psounis, and Salman Avestimehr. How much privacy does federated learning with secure aggregation guarantee? *arXiv preprint arXiv:2208.02304*, 2022.

-
- Mehrdad Farajtabar, Navid Azizan, Alex Mott, and Ang Li. Orthogonal gradient descent for continual learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 3762–3773. PMLR, 2020.
- Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A Rusu, Alexander Pritzel, and Daan Wierstra. PathNet: Evolution channels gradient descent in super neural networks. *arXiv preprint arXiv:1701.08734*, 2017.
- Robert M French. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4):128–135, 1999.
- Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*, 2013.
- Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.
- Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pp. 5132–5143. PMLR, 2020.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. Fair resource allocation in federated learning. *arXiv preprint arXiv:1905.10497*, 2019.
- 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.
- Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- Wei Yang Bryan Lim, Nguyen Cong Luong, Dinh Thai Hoang, Yutao Jiao, Ying-Chang Liang, Qiang Yang, Dusit Niyato, and Chunyan Miao. Federated learning in mobile edge networks: A comprehensive survey. *IEEE Communications Surveys & Tutorials*, 22(3):2031–2063, 2020.
- David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30, 2017.
- 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*, volume 3, 2022.
- Zheda Mai, Ruiwen Li, Jihwan Jeong, David Quispe, Hyunwoo Kim, and Scott Sanner. Online continual learning in image classification: An empirical survey. *Neurocomputing*, 469:28–51, 2022.
- Arun Mallya and Svetlana Lazebnik. PackNet: Adding multiple tasks to a single network by iterative pruning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 7765–7773, 2018.

-
- Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165. Elsevier, 1989.
- 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*, pp. 1273–1282. PMLR, 2017.
- Sanket Vaibhav Mehta, Darshan Patil, Sarath Chandar, and Emma Strubell. An empirical investigation of the role of pre-training in lifelong learning. *arXiv preprint arXiv:2112.09153*, 2021.
- Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In *International Conference on Machine Learning*, pp. 4615–4625. PMLR, 2019.
- German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural networks*, 113:54–71, 2019.
- Quang Pham, Chenghao Liu, and Steven Hoi. DualNet: Continual learning, fast and slow. *Advances in Neural Information Processing Systems*, 34:16131–16144, 2021.
- Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. GDumb: A simple approach that questions our progress in continual learning. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pp. 524–540. Springer, 2020.
- Roger Ratcliff. Connectionist models of recognition memory: constraints imposed by learning and forgetting functions. *Psychological review*, 97(2):285, 1990.
- 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*, pp. 2001–2010, 2017.
- Joseph Redmon and Ali Farhadi. YOLO9000: better, faster, stronger. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7263–7271, 2017.
- G Anthony Reina, Alexey Gruzdev, Patrick Foley, Olga Perepelkina, Mansi Sharma, Igor Davidyuk, Ilya Trushkin, Maksim Radionov, Aleksandr Mokrov, Dmitry Agapov, et al. OpenFL: An open-source framework for federated learning. *arXiv preprint arXiv:2105.06413*, 2021.
- Anthony Robins. Catastrophic forgetting, rehearsal and pseudorehearsal. *Connection Science*, 7(2): 123–146, 1995.
- Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- Gobinda Saha, Isha Garg, and Kaushik Roy. Gradient projection memory for continual learning. *arXiv preprint arXiv:2103.09762*, 2021.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *International Conference on Machine Learning*, pp. 4548–4557. PMLR, 2018.
- Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. *Advances in neural information processing systems*, 30, 2017.
- Konstantin Shmelkov, Cordelia Schmid, and Karteek Alahari. Incremental learning of object detectors without catastrophic forgetting. In *Proceedings of the IEEE international conference on computer vision*, pp. 3400–3409, 2017.
- Neta Shoham, Tomer Avidor, Aviv Keren, Nadav Israel, Daniel Benditkis, Liron Mor-Yosef, and Itai Zeitak. Overcoming forgetting in federated learning on non-IID data. *arXiv preprint arXiv:1910.07796*, 2019.

-
- Sebastian Thrun. Is learning the n-th thing any easier than learning the first? *Advances in neural information processing systems*, 8, 1995.
- Yeshwanth Venkatesha, Youngeun Kim, Hyoungseob Park, Yuhang Li, and Priyadarshini Panda. Addressing client drift in federated continual learning with adaptive optimization. *Available at SSRN 4188586*, 2022.
- Jeffrey S Vitter. Random sampling with a reservoir. *ACM Transactions on Mathematical Software (TOMS)*, 11(1):37–57, 1985.
- Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 139–149, 2022.
- Mitchell Wortsman, Vivek Ramanujan, Rosanne Liu, Aniruddha Kembhavi, Mohammad Rastegari, Jason Yosinski, and Ali Farhadi. Supermasks in superposition. *Advances in Neural Information Processing Systems*, 33:15173–15184, 2020.
- Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 374–382, 2019.
- Xin Yao and Lifeng Sun. Continual local training for better initialization of federated models. In *2020 IEEE International Conference on Image Processing (ICIP)*, pp. 1736–1740. IEEE, 2020.
- Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dynamically expandable networks. *arXiv preprint arXiv:1708.01547*, 2017.
- 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*, pp. 12073–12086. PMLR, 2021.
- Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *International Conference on Machine Learning*, pp. 3987–3995. PMLR, 2017.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28, 2015.
- Zizhao Zhang, Han Zhang, Long Zhao, Ting Chen, Serkan Ö Arik, and Tomas Pfister. Nested hierarchical transformer: Towards accurate, data-efficient and interpretable visual understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 3417–3425, 2022.
- Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning with non-IID data. *arXiv preprint arXiv:1806.00582*, 2018.

A SUPPLEMENTARY RESULTS

In this section, we furnish additional experimental outcomes that serve to further bolster the findings of our primary investigation.

A.1 PROGRESSIVE PERFORMANCE OF FEDGP ACROSS TASKS

Fig. 5 depicts the average accuracy Acc_t measured at task $t = 1, 2, \dots, 10$ and the average forgetting Fgt_t measured at task $t = 2, 3, \dots, 10$. The accuracy of FedAvg rapidly drops as different tasks are given to the model, as expected. FedCurv and FedProx perform similarly to FedAvg, while A-GEM and DER partially alleviate forgetting, resulting in higher accuracies and reduced forgetting compared to FedAvg. Combining these baselines with FedGP lead to significant performance improvements, which allows the solid lines in the accuracy plot consistently remain at the top. For example, for the experiment on task-IL for sequential-CIFAR100, the accuracy measured at task 5 (denoted by Acc_5) is 55.37% for FedProx, while 71.12% for FedProx+FedGP. These results demonstrate that FedGP effectively mitigates forgetting and enhances existing methods in CFL.

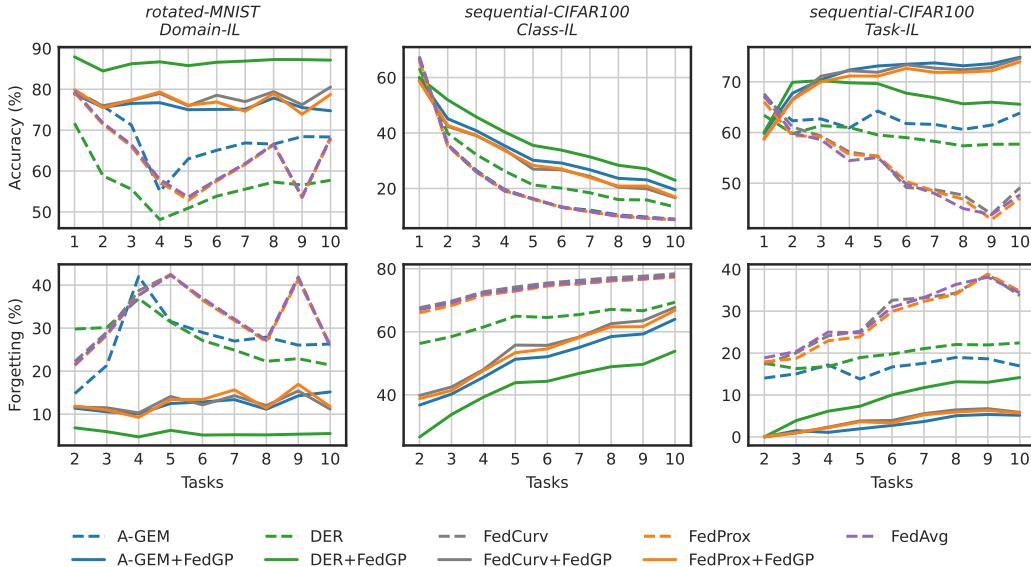


Figure 5: Evaluating accuracy (\uparrow) and forgetting (\downarrow) in multiple datasets with and without FedGP using a buffer size of 200. The solid lines indicate the results obtained with our method, while the dotted lines represent the results obtained without our method. The results show a significant improvement in accuracy as well as reduced forgetting for all settings.

A.2 FORGETTING ANALYSIS ACROSS DATASETS

We present the complementary information to Table 1 in Table 11, illustrating the extent of Fgt_T observed across multiple benchmark datasets. Our method exhibits exceptional effectiveness in mitigating forgetting. Remarkably, it demonstrates consistent performance across all datasets and baselines, making it a versatile solution.

In line with the presentation of forgetting in Table 11, we present the forgetting analysis when the number of clients is set to 20 in Table 12. Notably, our method exhibits consistent and impressive performance across varying numbers of users. It consistently proves its effectiveness regardless of the specific user count, showcasing its robustness and reliability.

Table 11: Average forgetting Fgt_T (%) (lower is better) on benchmark datasets at the final task T .

		rotated-MNIST (<i>Domain-IL</i>)		sequential-CIFAR10 (<i>Class-IL</i>)		sequential-CIFAR10 (<i>Task-IL</i>)	
Method	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP	
FL	25.98 \pm 3.2	11.66 \pm 2.7 (↓14.32)	80.69 \pm 3.6	78.62 \pm 4.3 (↓2.07)	15.37 \pm 4.8	4.49 \pm 1.9 (↓10.88)	
FedCurv	25.80 \pm 2.4	11.18 \pm 2.7 (↓14.62)	80.90 \pm 6.6	79.85 \pm 3.9 (↓1.05)	19.37 \pm 4.8	4.77 \pm 1.6 (↓14.60)	
FedProx	25.74 \pm 3.1	11.76 \pm 2.9 (↓13.98)	84.35 \pm 2.4	80.24 \pm 2.5 (↓4.11)	18.24 \pm 4.9	4.17 \pm 1.0 (↓14.07)	
FL+A-GEM	26.30 \pm 5.7	15.18 \pm 2.4 (↓11.12)	82.18 \pm 6.6	80.38 \pm 2.5 (↓1.80)	10.00 \pm 3.0	4.15 \pm 0.7 (↓5.85)	
FL+DER	21.42 \pm 4.0	5.51 \pm 1.2 (↓15.91)	60.98 \pm 14.6	47.88 \pm 7.2 (↓13.10)	6.34 \pm 4.9	2.73 \pm 1.3 (↓3.61)	
		permuted-MNIST (<i>Domain-IL</i>)		sequential-CIFAR100 (<i>Class-IL</i>)		sequential-CIFAR100 (<i>Task-IL</i>)	
FL	43.47 \pm 5.3	21.40 \pm 4.9 (↓22.07)	77.69 \pm 0.5	67.02 \pm 2.3 (↓10.67)	34.38 \pm 1.6	5.39 \pm 0.8 (↓28.99)	
FedCurv	42.88 \pm 5.0	22.85 \pm 3.5 (↓20.03)	78.40 \pm 0.9	67.75 \pm 0.8 (↓10.65)	33.71 \pm 2.2	5.86 \pm 0.7 (↓27.85)	
FedProx	42.59 \pm 5.6	20.77 \pm 5.6 (↓21.82)	77.35 \pm 0.4	66.81 \pm 2.2 (↓10.54)	34.79 \pm 3.6	5.69 \pm 0.9 (↓29.10)	
FL+A-GEM	35.61 \pm 5.3	24.05 \pm 2.4 (↓11.56)	77.97 \pm 0.7	63.99 \pm 2.0 (↓13.98)	16.92 \pm 1.1	5.16 \pm 0.5 (↓11.76)	
FL+DER	45.33 \pm 5.0	34.71 \pm 5.0 (↓10.62)	69.37 \pm 1.7	53.84 \pm 6.7 (↓15.53)	22.43 \pm 0.7	14.16 \pm 1.7 (↓8.27)	

Table 12: The Fgt_T (%) (lower is better) performance measured when we have $K = 20$ users.

		rotated-MNIST (<i>Domain-IL</i>)		sequential-CIFAR10 (<i>Class-IL</i>)		sequential-CIFAR10 (<i>Task-IL</i>)	
Method	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP	
FL	31.00 \pm 9.5	13.45 \pm 3.6 (↓17.55)	82.62 \pm 3.1	73.39 \pm 4.5 (↓9.23)	17.93 \pm 2.7	6.14 \pm 4.9 (↓11.79)	
FedCurv	30.73 \pm 9.3	12.97 \pm 3.8 (↓17.76)	79.55 \pm 3.8	75.38 \pm 5.3 (↓4.17)	18.19 \pm 3.0	9.14 \pm 3.1 (↓9.05)	
FedProx	31.04 \pm 9.7	13.31 \pm 3.4 (↓17.73)	82.94 \pm 1.1	78.67 \pm 4.2 (↓4.27)	20.60 \pm 2.6	8.52 \pm 3.0 (↓12.08)	
FL+A-GEM	25.22 \pm 8.8	11.02 \pm 3.0 (↓14.20)	82.39 \pm 2.4	80.25 \pm 4.1 (↓2.14)	12.29 \pm 2.2	4.00 \pm 2.4 (↓8.29)	
FL+DER	28.93 \pm 6.6	5.18 \pm 1.1 (↓23.75)	55.10 \pm 9.8	60.90 \pm 3.8 (↑5.80)	3.20 \pm 1.6	2.71 \pm 1.7 (↓0.49)	
		permuted-MNIST (<i>Domain-IL</i>)		sequential-CIFAR100 (<i>Class-IL</i>)		sequential-CIFAR100 (<i>Task-IL</i>)	
FL	24.27 \pm 5.2	8.67 \pm 7.0 (↓15.60)	73.05 \pm 0.5	62.71 \pm 0.9 (↓10.34)	27.07 \pm 1.7	2.48 \pm 0.7 (↓24.59)	
FedCurv	24.02 \pm 5.4	8.10 \pm 5.4 (↓15.92)	80.07 \pm 0.5	68.58 \pm 1.1 (↓11.49)	34.63 \pm 1.7	3.48 \pm 0.6 (↓31.15)	
FedProx	23.01 \pm 5.7	5.93 \pm 5.1 (↓17.08)	79.46 \pm 0.5	68.40 \pm 0.9 (↓11.06)	32.82 \pm 1.4	4.13 \pm 0.7 (↓28.69)	
FL+A-GEM	22.12 \pm 4.9	9.45 \pm 5.4 (↓12.67)	72.97 \pm 1.1	60.27 \pm 1.3 (↓12.70)	12.54 \pm 1.3	2.66 \pm 0.2 (↓9.88)	
FL+DER	32.26 \pm 1.1	27.30 \pm 4.2 (↓4.96)	67.07 \pm 0.8	47.74 \pm 3.8 (↓19.33)	19.78 \pm 1.7	8.67 \pm 1.4 (↓11.11)	

A.3 EXTENDED ANALYSIS ON THE INFLUENCE OF BUFFER SIZE

In the main body of our study, we examine the influence of different buffer sizes on the performance metric Acc_T , utilizing rotated-MNIST and sequential-CIFAR100 datasets. To further augment our analysis, we have included two additional datasets in Table 13, incorporating various buffer sizes. By evaluating Acc_T (where higher values indicate better performance), we discovered that our proposed method, referred to as FedGP, consistently enhances the average accuracy across these two datasets.

Table 13: Impact of the buffer size on Acc_T (%)

		permuted-MNIST (<i>Domain-IL</i>)		sequential-CIFAR10 (<i>Class-IL</i>)		sequential-CIFAR10 (<i>Task-IL</i>)	
Buffer Size	Method	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP	w/o FedGP	w/ FedGP
200	FL+A-GEM	33.43 \pm 1.4	39.09 \pm 3.5 (↑5.66)	17.82 \pm 0.9	19.44 \pm 0.9 (↑1.62)	77.14 \pm 3.1	83.16 \pm 1.6 (↑6.02)
500		33.35 \pm 1.0	42.45 \pm 6.9 (↑9.10)	18.39 \pm 0.2	20.34 \pm 0.6 (↑1.95)	78.43 \pm 3.0	85.95 \pm 0.6 (↑7.52)
5120		32.72 \pm 1.4	40.07 \pm 2.5 (↑7.35)	16.41 \pm 2.6	20.64 \pm 2.2 (↑4.23)	73.89 \pm 3.3	86.82 \pm 1.5 (↑12.93)
200	FL+DER	19.79 \pm 1.7	43.43 \pm 0.9 (↑23.64)	18.44 \pm 3.7	30.94 \pm 3.8 (↑12.50)	69.34 \pm 3.2	77.99 \pm 0.8 (↑8.65)
500		19.17 \pm 1.6	43.38 \pm 2.4 (↑24.21)	20.81 \pm 3.6	29.78 \pm 4.3 (↑8.97)	71.17 \pm 1.5	74.98 \pm 3.5 (↑3.81)
5120		18.57 \pm 1.4	44.68 \pm 2.4 (↑26.11)	34.75 \pm 2.2	42.38 \pm 4.5 (↑7.63)	78.22 \pm 2.3	81.94 \pm 1.7 (↑3.72)

A.4 RANDOM SAMPLING

We implement a more realistic federated learning environment by applying uniform sampling techniques to randomly select the participating clients in each round. We conduct experiments on CIFAR100. A total of 50 clients is set up, and during each communication, only a random 50% of the clients participate in training. As can be seen, even in such a scenario, where our algorithm cannot update the reference gradient using the local buffer from all clients, there is still an improvement in performance using our algorithm.

Table 14: Average accuracy Acc_T (%) with 50 clients and 50% client sampling rate, for sequential-CIFAR100

Method	Class-IL	Task-IL
FL	7.46 ± 0.08	43.85 ± 1.33
FL+FedGP	9.34 ± 0.31 ($\uparrow 1.88$)	65.76 ± 0.48 ($\uparrow 21.91$)

A.5 PERFORMANCE ON THE CURRENT TASK

Balancing the retention of old tasks and the learning of new ones is a common challenge in continual learning. It can be difficult to determine the best approach, especially when two tasks are significantly different. This is a challenge faced by many methods in continual learning.

We provided additional experimental results on the performance measured for the current task. The below Fig. 6 shows the Class-IL accuracy of FedGP (with buffer size 200) and FL for S-CIFAR100, where the total number of tasks is set to 10. During the continual learning process, we measured the accuracy of each model for the current task. One can confirm that using FedGP does not hurt the current task accuracy, compared with FL. Note that this shows that FedGP does not impair the performance of the current task, while also alleviating the forgetting in upcoming rounds.

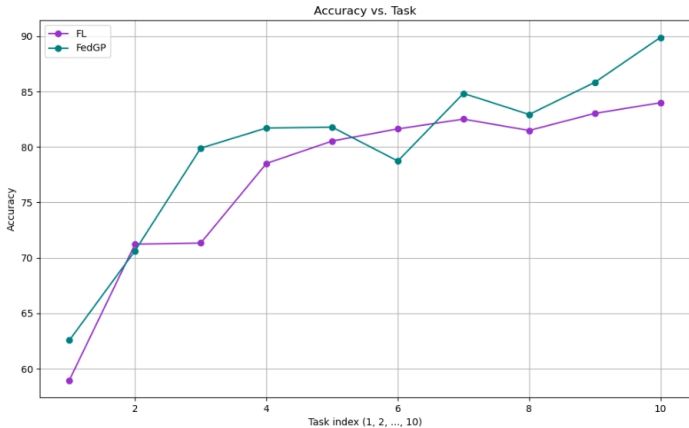


Figure 6: Class-IL Accuracy (%) of current task for FL and FedGP on the S-CIFAR100

A.6 COMPARISON WITH FOT

We compared our method with SOTA paper (Bakman et al., 2023) proposing Federated Orthogonal Training (FOT) algorithm for continual federated learning.

In FOT, at the end of each task, the server aggregates the activations of each local model (computed for local data points) and computes the subspace spanned by the aggregated activations. This subspace is used during the global model update process; the gradient is updated in the direction that is orthogonal to the subspace. Note that FOT has several advantages compared with existing baselines; the privacy leakage is mitigated, the communication cost is reduced, and the solution has theoretical guarantees.

While both FOT and our FedGP project the gradients on the subspace specified by previous tasks, they have two main differences. First, the subspace is defined in a different manner: FOT relies on the representations of local model activations to define S . FedGP, on the other hand, relies on the gradient of model computed on its local buffer data. Second, FOT projects the gradient computed at the server side, while FedGP projects the gradient computed at each client.

The table 15 compares the accuracy of FL, FOT and FedGP (with buffer size 200) for P-MNIST and S-CIFAR100 under the task incremental learning scenario, consistent with the paper’s benchmarks. Assuming that the local buffer is available, FedGP outperforms FOT.

Table 15: Comparative accuracy (%) performance analysis of FL, FOT, and FL+FedGP

Methods	P-MNIST (Domain-IL)	S-CIFAR100 (Task-IL)
FL	25.92±2.1	47.74±1.2
FOT	23.77±1.1	50.57±1.5
FL+FedGP	34.23±2.7	74.71±0.9

A.7 BACKWARD AND FORWARD TRANSFER METRICS

Our method outperforms FedAvg (FL) in both Backward and Forward Transfer metrics across the sequential-CIFAR10 and sequential-CIFAR100 datasets, as shown in the Table 16.

Table 16: Backward and Forward Transfer (\uparrow) Results for sequential-CIFAR100 and sequential-CIFAR10

Metric	Dataset	Methods	Class-IL	Task-IL
Backward	CIFAR100	FL	-78.11	-36.52
Backward	CIFAR100	FL+FedGP	-72.24 (\uparrow 5.87)	-3.68 (\uparrow 32.84)
Backward	CIFAR10	FL	-78.78	-14.48
Backward	CIFAR10	FL+FedGP	-78.55 (\uparrow 0.23)	-0.60 (\uparrow 13.88)
Forward	CIFAR100	FL	16.98	16.98
Forward	CIFAR100	FL+FedGP	17.16 (\uparrow 0.18)	17.48 (\uparrow 0.50)
Forward	CIFAR10	FL	12.75	12.74
Forward	CIFAR10	FL+FedGP	12.98 (\uparrow 0.23)	12.99 (\uparrow 0.25)

A.8 EFFECT OF DIFFERENT CURRICULUM.

We evaluate how the performance of FedGP changes when we shuffle the order of tasks in the continual learning. We randomly shuffle the sequential-CIFAR100 task order and label them as curriculum 1 to 4, as shown in the Table 17. Regardless of the different curriculum, FL+FedGP outperforms FedAvg.

Table 17: Average accuracy Acc_T (%) across randomized curriculum in sequential-CIFAR100.

Curriculum	Methods	Class-IL	Task-IL
1	FL	8.15	46.25
1	FL+FedGP	12.10 (\uparrow 3.95)	72.69 (\uparrow 26.44)
2	FL	8.46	47.56
2	FL+FedGP	14.37 (\uparrow 5.91)	73.19 (\uparrow 25.63)
3	FL	8.82	45.04
3	FL+FedGP	12.58 (\uparrow 3.76)	74.71 (\uparrow 29.67)
4	FL	7.87	43.87
4	FL+FedGP	14.85 (\uparrow 6.98)	73.74 (\uparrow 29.87)

A.9 ADDITIONAL HYPERPARAMETERS FOR SPECIFIC METHODS

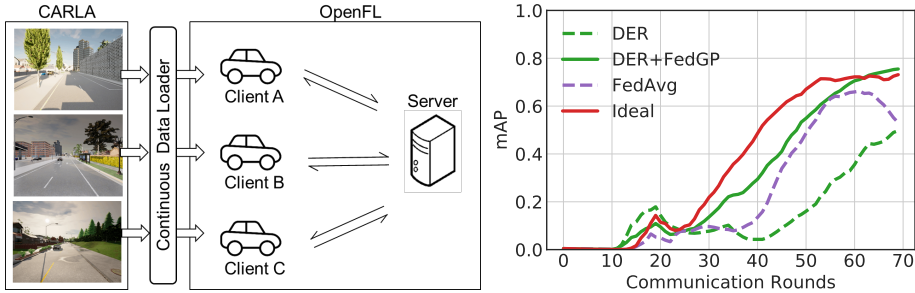
In addition to the hyperparameters discussed in the main paper, additional method-specific hyperparameters are outlined in Table 18.

Table 18: Additional hyperparameters for specific methods.

Method	Parameter	Values
FL+DER	Regularization Coefficient	sequential-CIFAR10 (0.3), Others (1)
FL+L2P	Communication Round R	rotated-MNIST (5), permuted-MNIST (1), sequential-CIFAR10 (20), sequential-CIFAR100 (20)
CFed	Surrogate Dataset	sequential-CIFAR10 (CIFAR100), sequential-CIFAR100 (CIFAR10)

Note: No server distillation included.

B OBJECT DETECTION



(a) Framework for automotive data evaluation. (b) Object detection performance comparison.

Figure 7: (a) The data loader continuously supplies data from CARLA camera outputs to individual FL clients. Each client trains on its local data and updates its buffer to retain old knowledge. (b) The result shows the object detection performance comparison between Ideal, FedAvg, DER, and DER+FedGP on a realistic CARLA dataset.

Here we test FedGP on realistic streaming data (Dai et al., 2023) which leverage two open source tools, an urban driving simulator (CARLA (Dosovitskiy et al., 2017)) and a FL framework (OpenFL (Reina et al., 2021)). As shown in Fig. 7a, CARLA provides OpenFL with a real-time collection of continuous streaming vehicle camera output data and automatic annotation about object detection. This streaming data capture the spatio-temporal dynamics of data generated from real-world applications. After loading data of vehicles from CARLA, OpenFL performs collaborative training over multiple clients.

We evaluate the solutions to the forgetting problem by spawning two vehicles in a virtual town. During the training of the tinyYOLO (Redmon & Farhadi, 2017) object detection model, we use a custom loss that combines classification, detection and confidence losses. In order to quantify the quality of the incremental model trained by various baselines, we report a common metric, namely, mean average precision (mAP). This metric assesses the correspondence between the detected bounding boxes and the ground truth, with higher scores indicating better performance. To calculate mAP, we analyze the prediction results obtained from pre-collected driving snippets of vehicular clients. These driving snippets are gathered by navigating the town over a duration of 3000 simulation seconds.

For those experiments on realistic CARLA streaming data, we compare the performances of Ideal, FedAvg, DER and DER+FedGP. In the Ideal scenario, the client possesses sufficient memory to retain all data from prior tasks, enabling joint training on all stored data. The last two methods are equipped with buffer size of 200. We train for 70 communication rounds and each round continues for about 200 simulation seconds. The results are presented in Fig. 7b. Note that at communication round 60, one client gets on the highway, which incurs a domain shift. One can confirm that the performance of FedAvg degrades in such domain shift scenario, whereas DER and DER+FedGP maintain the accuracy. Moreover, FedGP nearly achieves the performance of the ideal scenario with infinite buffer size, demonstrating the effectiveness of our method.

C CONTINUAL LEARNING METHODS WITH FEDGP

We provided the pseudocode for Algorithm 2 modifications when implementing FL+DER+FedGP and FL+A-GEM+FedGP, respectively presented in Algorithm 5 and Algorithm 6. Other FL+CL and CFL methods are also combined with FedGP in a similar manner.

Algorithm 5 incorporates Dark Experience Replay (DER) into the local update process on client $k \in [K]$. When the server sends the global model w to client k , the client calculates the output logits or pre-softmax response z . In addition, the client samples past data (x', y') and the corresponding logits z' from the buffer \mathcal{M}^k . To address forgetting, the regularization term considers the Euclidean distance between the sampled output logits and the current model’s output logits on buffer data. The gradient g is then refined using this regularization term to minimize the discrepancy between the

Algorithm 5 DER ClientUpdate at client k

Input: Task index t , model w , buffer gradient g_{ref} , batch size β , regularization coefficient λ , and learning rate α
Load the dataset \mathcal{D}_t^k and local buffer \mathcal{M}^k
Initialize $n = 0$ at the first task
for each batch $\{(x_i, y_i)\}_{i=1}^\beta$ in \mathcal{D}_t^k **do**
 Let $X = \{x_i\}_{i=1}^\beta$ and $Y = \{y_i\}_{i=1}^\beta$
 $Z \leftarrow h(X; w)$
 where $f(X; w) := \sigma(h(X; w))$
 Sample (X', Z', Y') from \mathcal{M}^k
 $\ell_{\text{reg}} \leftarrow \lambda \|Z' - h(X'; w)\|_2^2$
 $g = \nabla_w [\ell(Y, f(X; w)) + \ell_{\text{reg}}]$
 $\tilde{g} \leftarrow g - \text{proj}_{g_{\text{ref}}} g \cdot \mathbf{1}(g_{\text{ref}}^\top g \leq 0)$
 $w \leftarrow w - \alpha \tilde{g}$
 ReservoirSampling($\mathcal{M}^k, (X, Z, Y), n$)

 $n \leftarrow n + \beta$
end for
Return w

Algorithm 6 A-GEM ClientUpdate at client k

Input: Task index t , model w , buffer gradient g_{ref} , batch size β
Load the dataset \mathcal{D}_t^k , local buffer \mathcal{M}^k
Initialize $n = 0$ at the first task
for each batch $\{(x_i, y_i)\}_{i=1}^\beta$ in \mathcal{D}_t^k **do**
 $g_c = \nabla_w \left[\frac{1}{\beta} \sum_{i=1}^\beta \ell(y_i, f(x_i; w)) \right]$
 Sample $\{(x'_i, y'_i)\}_{i=1}^\beta$ from \mathcal{M}^k
 $g_b = \nabla_w \left[\frac{1}{\beta} \sum_{i=1}^\beta \ell(y'_i, f(x'_i; w)) \right]$
 $g \leftarrow g_c - \text{proj}_{g_b} g_c \cdot \mathbf{1}(g_b^\top g_c \leq 0)$
 $\tilde{g} \leftarrow g - \text{proj}_{g_{\text{ref}}} g \cdot \mathbf{1}(g_{\text{ref}}^\top g \leq 0)$
 $w \leftarrow w - \alpha \tilde{g}$ for some learning rate α
 ReservoirSampling($\mathcal{M}^k, \{(x_i, y_i)\}_{i=1}^\beta, n$)

 $n \leftarrow n + \beta$
end for
Return w to server

current and past output logits, thereby mitigating forgetting. The following steps are the same as in the main text.

Algorithm 6 combines with A-GEM, applying gradient projection twice. First, the client computes the gradient g_c with respect to the new data from \mathcal{D}_t^k . After replaying previous samples (x', y') stored in the local buffer \mathcal{M}^k , the client computes the gradient g_b with respect to this buffered data. If these gradients differ significantly in terms of their direction, the client projects g_c onto g_b to remove interference.