FedMeS: Personalized Federated Continual Learning Leveraging Local Memory

Jin Xie¹, Chenqing Zhu¹, Songze Li^{1,2}

¹The Hong Kong University of Science and Technology (Guangzhou) ²The Hong Kong University of Science and Technology {jxie171, czhu032}@connect.hkust-gz.edu.cn, songzeli@ust.hk

Abstract

We focus on the problem of Personalized Feder-1 ated Continual Learning (PFCL): a group of dis-2 tributed clients, each with a sequence of local tasks 3 on arbitrary data distributions, collaborate through 4 a central server to train a personalized model at 5 each client, with the model expected to achieve 6 good performance on all local tasks. We propose 7 a novel PFCL framework called Federated Mem-8 ory Strengthening (FedMeS) to address the chal-9 lenges of client drift and catastrophic forgetting. 10 In FedMeS, each client stores samples from pre-11 vious tasks using a small amount of local mem-12 ory, and leverages this information to both 1) cal-13 ibrate gradient updates in local training; and 2) per-14 form KNN-based Gaussian inference to facilitate 15 local inference. FedMeS is designed to be task-16 oblivious, such that the same inference process is 17 applied to samples from all tasks to achieve good 18 performance. FedMeS is analyzed theoretically 19 and evaluated experimentally. It is shown to out-20 perform all baselines in average accuracy and for-21 getting rate, over various combinations of datasets, 22 task distributions, and client numbers. 23

24 **1** Introduction

Federated learning (FL) [McMahan *et al.*, 2017] is an emerging distributed learning framework that allows for collaborative training of a model across multiple clients while keeping
their raw data locally stored. A typical FL process involves
local training on each client and global model aggregation on
a cloud server, with only model updates or gradients being
shared between clients and server.

Data collected from different clients in an FL system of-32 ten have drastically different distributions. As seen in Figure 33 1(a), this can lead to model parameter divergence and *client* 34 drift [Venkatesha et al., 2022], causing potentially poor per-35 formance for certain clients. The conventional way of train-36 ing a single model is insufficient to fit all the non-IID data, 37 and a personalized model needs to be trained for each partic-38 ipating client [Huang et al., 2021; Fallah et al., 2020], which 39 is known as personalized FL. 40

Another key characteristic in real-world FL systems is 41 that clients are continuously collecting new data (new task) 42 which may exhibit different distributions from previous local 43 data (tasks). Hence, it would be preferable to train a local 44 model that is able to achieve consistently good performance 45 in all local tasks. In an FL system, the problem is solved 46 by federated continual learning (FCL) [Yao and Sun, 2020; 47 Shoham et al., 2019; Yoon et al., 2021]. The phenomenon of 48 a model failing to perform well on previously trained tasks 49 is called *catastrophic forgetting* [Kirkpatrick et al., 2017], 50 which is illustrated in Figure. 1(b). 51



Figure 1: (a) Illustration of the model parameter divergence with non-IID datasets. (b) Illustration of catastrophic forgetting. (c) An overview of a PFCL system in IIoT scenario.

In practical FL systems, the data and task heterogeneity of-52 ten exist both across clients and over time on a single client. 53 For instance, as shown in Figure 1(c), in a IIoT scenario, mul-54 tiple factories manufacturing different products would like to 55 use FL for training defect detection models collaboratively. 56 Other than the difference between the types of products, each 57 factory may experience change of tasks over time due to e.g., 58 change of recipe and upgrade of the production line. Aim-59 ing to address the challenges of client drift and catastrophic 60 forgetting simultaneously, in the paper we focus on the per-61 sonalized federated continual learning (PFCL) problem. In a 62 PFCL system, each client observes a stream of arbitrarily dif-63

ferent tasks, and would like to collaborate through the server
 to train a personalized model, which performs well on all lo cal tasks.

While there have been prior attempts at tackling the FCL problem, like FedWeIT [Yoon *et al.*, 2021], where taskgeneric and task-specific knowledge are shared across clients to decompose the model parameters. However, as demonstrated in [Venkatesha *et al.*, 2022], FedWeIT struggles to handle the issue of client drift caused by data heterogeneity and does not address the scalability of tasks.

In this paper, we propose a novel federated learning 74 PFCL framework called Federated Memory Strengthening 75 (FedMeS). FedMeS utilize small amount of local memory 76 at each client to store information about previous tasks, and 77 78 leverage this memory to assist both the training and inference processes. During training process, the gradients are 79 constantly calibrated against the data samples from previ-80 ous tasks to avoid catastrophic forgetting. A newly designed 81 regularization term adjusted by a loss-based parameter is 82 used to facilitate the training of personalized models using 83 information from the global model. In the inference pro-84 cess, FedMeS directly leverages the memory information in 85 training process to perform KNN-based Gaussian inference, 86 further strengthening the model's personalization capability. 87 Moreover, FedMeS exhibits a major advantage in being task-88 oblivious, meaning that the inference process for test samples 89 from all tasks is identical, and all are expected to achieve high 90 91 performance.

Through extensive experiments with various dataset combinations, task constructions, task distributions, and client numbers, we show that FedMeS uniformly outperforms all baselines in terms of accuracy metrics and forgetting rate. These results highlight the potential of FedMeS for realworld applications and as a basis for future research in the area of PFCL.

99 2 Related Work

100 2.1 Personalized Federated Learning

A lot of work has been done in personalized FL. A simple 101 idea is by deploying a global model and fine tuning param-102 eters through gradient descent on local clients [Cheng et al., 103 2021; Yu et al., 2020b; Zhang et al., 2022]. Meta-learning 104 based FL methods realize model personalization through hy-105 perparameters [Khodak et al., 2019]. PerKNN [Marfog et al., 106 2022] is a special case, where embeddings of training samples 107 are stored for local memorization for KNN-based Gaussian 108 inference. The mainstream design is to interpolate a global 109 model and one local model per client, and the task-specific 110 models are learned both globally and locally [Achituve et al., 111 2021; Shen *et al.*, 2020]. Like using regularization terms on 112 proximal models to help construct personalized information 113 [Li et al., 2021; Marfoq et al., 2021]. 114

115 2.2 Continual Learning

Memory replay methods are widely used in continual learning (CL) to maintain prediction accuracies of past tasks. Generally speaking, a memory buffer is used to store previous

data which are replayed while learning a new task to allevi-119 ate forgetting [Wang et al., 2022; Shim et al., 2021]. Ex-120 perience replay (ER) jointly optimizes the network param-121 eters by interleaving the previous task exemplars with cur-122 rent task data [Riemer et al., 2018; Isele and Cosgun, 2018]. 123 An alternative solution is by constrained optimization. GEM 124 [Lopez-Paz and Ranzato, 2017] and A-GEM [Chaudhry et al., 125 2018] leverage episodic memory to compute previous task 126 gradients to constrain the current update step. Besides replay 127 methods, regularization-based methods [Yu et al., 2020a; 128 Shi et al., 2021] and parameter isolation methods [De Lange 129 et al., 2021] have also been proposed for CL. 130

Federated Continual Learning. Although a lot of work 131 has been done in CL, just a few works have tried to use CL 132 in a federated setting. Besides the FedWeIT, other meth-133 ods, like LFedCon2 [Casado et al., 2020], use traditional 134 classifiers instead of DNN and propose an algorithm dealing 135 with a concept drift based on ensemble retraining. FLwF and 136 FLwF-2T [Usmanova et al., 2021] use a distillation-based 137 approach dealing with catastrophic forgetting in FL scenario 138 and focus on the class-incremental learning scenario. 139

3 Problem Definition

We consider an FL system that consists of m clients and a central server. Over time, each client k (k = 1, ..., m) continually receives private datasets from a sequence of T machine learning tasks. For each task t (t = 1, ..., T), the corresponding dataset at client k is denoted as \mathcal{D}_k^t . We focus on a general non-IID case, where \mathcal{D}_k^t is drawn from some probability distribution \mathcal{P}_k^t , and no particular relationships for \mathcal{P}_k^t across k and t are assumed.

140

The conventional FL problem corresponds to a single-task 149 scenario. For a particular task t, the following objective function is optimized over the global model w 151

$$\min_{\mathbf{w}} \mathcal{G}(F_1(\mathbf{w}; \mathcal{D}_1^t), \dots, F_m(\mathbf{w}; \mathcal{D}_m^t)),$$
(1)

where $F_k(\cdot)$ is a local objective for client k over current to the dataset \mathcal{D}_k^t , and $\mathcal{G}(\cdot)$ is an aggregation function. In [McMahan *et al.*, 2017], $\mathcal{G} = \sum_{i=1}^m p_k F_k(\mathbf{w}, \mathcal{D}_k^t)$ is chosen as a weighted sum with $\sum p_k = 1$.

When task changes over time on each client, the client in-156 tends to obtain an evolving local model which maintains good 157 performance on all its previous tasks. Motivated by this need, 158 we introduce the concept of continual learning [Chaudhry et 159 al., 2018] into the personalized FL framework, and formally 160 define the personalized federated continual learning (PFCL) 161 problem. Specifically, for client k (k = 1, ..., m) with a se-162 quence of local datasets $(\mathcal{D}_k^1, ..., \mathcal{D}_k^T)$, the personalized model w_k^t for task t (t = 1, ..., T) is obtained through 163 164

$$\min_{\substack{w_k^t \\ \mathbf{w}_k^t}} G_k(w_k^t; \mathbf{w}^*) = \mathcal{L}\left(w_k^t; \mathcal{D}_k^t\right) + \frac{\lambda}{2} \left\|w_k^t - \mathbf{w}^*\right\|^2
\text{s.t. } \mathbf{w}^* \in \operatorname*{arg\,min}_{\mathbf{w}} \mathcal{G}\left(\mathcal{L}(\mathbf{w}; \mathcal{D}_1^t), \dots, \mathcal{L}(\mathbf{w}; \mathcal{D}_m^t)\right) \qquad (2)
\mathcal{L}\left(w_k^t; \mathcal{M}_k^t\right) \le \mathcal{L}\left(w_k^{t-1}; \mathcal{M}_k^t\right)$$

Here $\mathcal{L}(w; \mathcal{D})$ is the empirical loss of w on dataset \mathcal{D} . $\mathcal{M}_k^t \subset \bigcup_{i=1}^{t-1} \mathcal{D}_k^i$ is the episodic memory on client k storing samples 166



Figure 2: Overall workflow of FedMeS. Local episodic memory is utilized in both the training and inference processes.

from all previous tasks (i.e., task 1 to t - 1). Each client 167 reserves an equal amount of memory to store some samples 168 from each task. When learning the first task, $\mathcal{M}_k^1 = \emptyset$. Ev-169 ery time when client k finishes learning task t, examples are 170 randomly sampled from \mathcal{D}_k^t and stored in the allocated space. 171 The newly stored examples, together with current \mathcal{M}_k^t , con-172 stitute \mathcal{M}_k^{t+1} as the episodic memory used for next task. The 173 constraint $\mathcal{L}(w_k^t; \mathcal{M}_k^t) \leq \mathcal{L}(w_k^{t-1}; \mathcal{M}_k^t)$ ensures that the 174 model obtained from the new task t has a lower loss than 175 the previous model over the samples of past datasets. This ef-176 fectively alleviates the forgetting of previous tasks on current 177 models. 178

179 4 FedMeS

We propose an algorithm called Federated Memory 180 Strengthening (FedMeS) to solve the PFCL problem defined 181 in (2). The key idea of FedMeS is to flexibly use the lo-182 cal memory on the samples of previous tasks in both the local 183 training and inference processes. In training process, the local 184 memory is used for gradient correction to avoid catastrophic 185 forgetting; in inference process, a KNN algorithm based on 186 the representations of local samples helps to improve the ac-187 curacy of the personalized model. The overall workflow of 188 FedMeS is illustrated in Figure 2. 189

190 4.1 Training Process of FedMeS

As in a conventional FL setting, the training of task t (t = 1, ..., T) proceeds over multiple global iterations between the server and the clients. In each global iteration, the server broadcasts the global model **w** to the clients, waits for clients to upload personalized models w_k^t , and aggregates them to update the global model $\mathbf{w} \leftarrow \frac{1}{m} \sum_{k=1}^m w_k^t$.

¹⁹⁷ During local training process, each client k needs to run ¹⁹⁸ multiple local iterations to update w_k^t . We focus on describing ¹⁹⁹ parameter updates in a single local iteration. As shown in (2), PFCL is a constraint minimizing problem, 200 for which traditional Stochastic Gradient Descent does not 201 directly apply. As shown in [Lopez-Paz and Ranzato, 2017], 202 the constraint $\mathcal{L}(w_k^t; \mathcal{M}_k^t) \leq \mathcal{L}(w_k^{t-1}; \mathcal{M}_k^t)$ is equivalent to 203 the following condition on the inner product of gradients on 204 the current and previous tasks: 205

$$\left\langle \nabla \mathcal{L}\left(w_{k}^{t}; \mathcal{D}_{k}^{t}\right), \nabla \mathcal{L}\left(w_{k}^{t}; \mathcal{M}_{k}^{t}\right) \right\rangle \geq 0.$$
 (3)

By this transformation, it is not necessary to store the old 206 parameters w_k^{t-1} and compute loss on \mathcal{M}_k^t in every iteration; 207 only the inner product needs to be computed and compared. 208 If the inequality in (3) is satisfied, it means that the updates on 209 current task t and the local memory are roughly in the same 210 direction, so the optimization on current task would not neg-211 atively impact the performance of past tasks, and it is safe to 212 update the model along the gradient of current task as follows, 213 for some learning rate η_1 : 214

$$w_k^t = w_k^t - \eta_1 \left(\nabla \mathcal{L}(w_k^t; \mathcal{D}_k^t) + \lambda \| w_k^t - \mathbf{w} \| \right).$$
(4)

When (3) does not hold, client k first adjusts its local weights w_k^t to avoid forgetting, through the following gradient correction step, for some learning rate η_2 : 217

$$w_{k}^{t} = w_{k}^{t} - \eta_{2} \left(\nabla \mathcal{L}(w_{k}^{t}; \mathcal{D}_{k}^{t}) - \frac{\nabla \mathcal{L}(w_{k}^{t}; \mathcal{D}_{k}^{t})^{\top} \nabla \mathcal{L}(w_{k}^{t}; \mathcal{M}_{k}^{t})}{\nabla \mathcal{L}(w_{k}^{t}; \mathcal{M}_{k}^{t})^{\top} \nabla \mathcal{L}(w_{k}^{t}; \mathcal{M}_{k}^{t})} \nabla \mathcal{L}(w_{k}^{t}; \mathcal{M}_{k}^{t}) \right)$$
(5)

This gradient correction occurs *within* the local weights and does not involve global weights (the reason why we do not need term $\lambda \| w_k^t - \mathbf{w} \|$ in (5)).

We note that *only one* of (4) and (5) would be executed in every local training iteration: we update the gradients based on the local objective function in (2) only when (3) holds, which means that the gradient update does not lead to catastrophic forgetting. When (3) is not satisfied, as demonstrated in [Chaudhry *et al.*, 2018], updating w_L^t 226 according to (5) multiple times allows the inner product $\langle \nabla \mathcal{L}(w_k^t; \mathcal{D}_k^t), \nabla \mathcal{L}(w_k^t; \mathcal{M}_k^t) \rangle$ (which is < 0 currently) to gradually approach and eventually exceed zero.

In FedMeS, rather than fixing regularization parameter λ , we propose a loss-based approach for dynamically adjusting λ . Specifically, we set the value of λ as:

$$\lambda = 2 \cdot \text{sigmoid}\left(\frac{1}{\mathcal{L}(\mathbf{w}, \mathcal{D}_k^t)}\right) \tag{6}$$

Intuitively, when $\mathcal{L}(\mathbf{w}, \mathcal{D}_k^t)$ is relatively large, it means the global model \mathbf{w} performs poorly on the current task of client k, and the personalized model w_k^t should deviate from the global model by decreasing λ . On the other hand, a small $\mathcal{L}(\mathbf{w}, \mathcal{D}_k^t)$ would encourage w_k^t to approach the global model \mathbf{w} which correspond to a lager λ . Here the sigmoid function is used to limit the value of λ within [0, 2].

240 4.2 Inference Process of FedMeS

As shown in Figure 2, FedMeS utilizes local memory not only to mitigate catastrophic forgetting during training, but also to improve the inference performance on test samples. Specifically, to perform an inference task after learning task t, a client k first generates a set of representation-label pairs (R-L pairs) from the current local memory as

$$\left\{ \left(P_{w_{k}^{t}}\left(\mathbf{m}_{k}^{i}\right), y_{\mathbf{m}_{k}^{i}}\right) : \left(\mathbf{m}_{k}^{i}, y_{\mathbf{m}_{k}^{i}}\right) \in \mathcal{M}_{k}^{t} \right\}$$
(7)

Here \mathbf{m}_k^i , $i = 1, ..., |\mathcal{M}_k^t|$, is the input of the *i*-th sample in \mathcal{M}_k^t , and $y_{\mathbf{m}_k^i}$ is its label.

Function $P_{w_k^t}(\mathbf{m}_k^t)$ generates an embedding representation of \mathbf{m}_k^i , which for example, could be the output of the last convolutional layer in the case of CNNs, or the output of an arbitary self-attention layer in the case of transformers. Then, for a test sample x (from unknown task), we first find the *K* nearest neighbors of x from the formed R-L pairs:

$$\mathcal{K}^{(K)}(\mathbf{x}) = \left\{ \left(P_{w_k^t}(\mathbf{m}_k^j), y_{\mathbf{m}_k^j} \right) : 1 \le j \le K \right\}$$
(8)

which satisfy

$$dist(P_{w_k^t}(\mathbf{x}), P_{w_k^t}(\mathbf{m}_k^j)) \le dist(P_{w_k^t}(\mathbf{x}), P_{w_k^t}(\mathbf{m}_k^{j+1})) \quad (9)$$

Here $dist(\cdot, \cdot)$ could be any distance metric, and in FedMeS the Euclidean distance is used. Denote $b_{w_k^t}(\mathbf{x})$ as the local estimate of the conditional probability $\mathcal{P}_{\mathcal{M}_k^t}(y|\mathbf{x})$, where $\mathcal{P}_{\mathcal{M}_k^t}$ is the probability distribution of \mathcal{M}_k^t . Then the *K* nearest neighbours found in $\mathcal{K}^{(K)}(\mathbf{x})$ are used to compute $[b_{w_k^t}(\mathbf{x})]_y$ with a Gaussian kernel:

$$[b_{w_k^t}(\mathbf{x})]_y \propto \sum_{j=1}^{K} \mathbf{1}_{y=y_{\mathbf{m}_k^j}} \times \\ \exp\{-dist(P_{w_k^t}(\mathbf{x}), P_{w_k^t}(\mathbf{m}_k^j))\}$$
(10)

Finally, the FedMeS prediction result of x on client k is obtain by the following distribution

$$\operatorname{FedMeS}_{k}(\mathbf{x}) \triangleq \theta_{k} \cdot b_{w_{k}^{t}}(\mathbf{x}) + (1 - \theta_{k})h_{w_{k}^{t}}(\mathbf{x}).$$
(11)

Here $h_{w_k^t}$ is the personalized local model parameterized by w_k^t , and $\theta_k \in (0, 1)$ is a hyperparameter which can be tuned through a local validation or cross-validation. 266

Remark 1. As proved in [Khandelwal et al., 2019], aug-267 menting the model inference with a memorization mechanism 268 (KNN in this case) helps to improve the performance. In 269 [Marfog et al., 2022], local memorization through KNN has 270 been applied to improve the accuracies of local models in 271 personalized FL, for a single task. FedMeS extends the ap-272 plication of this technique on episodic memorization over an 273 arbitrary sequence of tasks, via utilizing a subset of samples 274 from each task. Also, this inference enhancement of FedMeS 275 comes for free, as this memory is readily available from the 276 preceeding training process. 277

Remark 2. Another major advantage of FedMeS is that it 278 is task-oblivious. That is, the same inference process is ap-279 plied for all test samples, and no prior knowledge is needed 280 about which task the sample belongs to. This also reflects 281 the robustness of FedMeS: regardless of the original task, a 282 good inference performance is always guaranteed by a uni-283 fied FedMeS inference process. This is, however, not the case 284 for other task-incremental learning CL methods like in [De-285 lange et al., 2021]. 286

4.3 Convergence Analysis

In this section, we analyze the convergence performance of FedMeS. All proofs are omitted due to page limit. Following assumptions are made to facilitate the analysis. For each communication round r on client k when solving task t, denote $w_k^{(r)}$, $\mathbf{w}^{(r)}$ respectively as the value of w_k^t , \mathbf{w} at round r.

Assumption 1. The loss function $\mathcal{L}(w_k^{(r)})$ is c-strongly convex and L-smooth for k = 1, ..., m.

Assumption 2. The expectation of stochastic gradients of the loss function $\mathcal{L}(w_k^{(r)})$ is uniformly bounded at all devices and all iterations, i.e.: 296

$$\mathbb{E}[\|\nabla \mathcal{L}\left(w_k^{(r)}, \xi_k^r\right)\|^2] \le \sigma^2 \tag{12}$$

287

Assumption 3. The global model converges with rate g(r), i.e., there exists g(r) such that $\lim_{r\to\infty} g(r) = 0$, $\mathbb{E}[\|\mathbf{w}^{(r)} - \mathbf{w}^*\|^2] \leq g(r)$.

First, we discuss the situation where the constraint 303 $\mathcal{L}(w_k^t; \mathcal{M}_k^t) \leq \mathcal{L}(w_k^{t-1}; \mathcal{M}_k^t)$ in (2) is not satisfied, which corresponds to $\langle \nabla \mathcal{L}(w_k^t; \mathcal{D}_k^t), \nabla \mathcal{L}(w_k^t; \mathcal{M}_k^t) \rangle < 0$ in our al-304 305 gorithm. Under this circumstance, FedMeS starts to execute 306 (5) to perform gradient correction. We denote the iteration in-307 dex of repeating (5) as s(s = 1, 2, ...), and $g(w_k^{(s)}, \xi_k^s)$ as the 308 stochastic gradient of $\mathcal{L}(w_k^{(s)}; \mathcal{M}_k^t)$ at iteration s. Following 309 First and second moment limits assumptions in [Bottou et al., 310 2018], we make two assumptions below, 311

Assumption 4. There exists scalars $\mu_G \ge \mu > 0$ such that for all $s \in \mathbb{N}$, 312



Figure 3: (a - d) Average accuracy and Average forgetting rate among all clients in all learned tasks at *x*-th task on Split CIFAR-100 with 10 clients and 20 clients. (e - h) Average accuracy and Average forgetting rate among all clients in all learned tasks at *x*-th task on Split MiniImageNet with 10 clients and 20 clients.

$$\nabla \mathcal{L}\left(w_{k}^{(s)};\mathcal{M}_{k}^{t}\right)^{\top} \mathbb{E}_{\xi_{k}^{s}}\left[g\left(w_{k}^{(s)},\xi_{k}^{s}\right)\right] \geq \mu \|\nabla \mathcal{L}\left(w_{k}^{(s)};\mathcal{M}_{k}^{t}\right)\|_{2}^{2}, \quad (13)$$

$$\|\mathbb{E}_{\xi_k^s}[g\left(w_k^{(s)}, \xi_k^s\right)]\|_2 \le \mu_G \|\nabla \mathcal{L}\left(w_k^{(s)}; \mathcal{M}_k^t\right)\|_2$$

Assumption 5. There exists scalars $M \ge 0$ and $M_V \ge 0$ such that, for all $s \in \mathbb{N}$,

 $\mathbb{V}_{\xi_k^s}[g(w_k^{(s)},\xi_k^s)] \leq M + M_V \|\nabla \mathcal{L}(w_k^{(s)};\mathcal{M}_k^t)\|_2^2$ (14) **Theorem 6.** Under the assumptions above and with the up-

319 dating rule of (5), when $s \ge 2$, we have,

$$\mathbb{E}[\mathcal{L}(w_k^{(s)}; \mathcal{M}_k^t) - \mathcal{L}(w_k^*; \mathcal{M}_k^t)] \le \frac{LM}{2c^2\mu^2}$$
(15)

Using Theorem 6, as long as the constraint $\mathcal{L}(w_k^t; \mathcal{M}_k^t) \leq \mathcal{L}(w_k^{t-1}; \mathcal{M}_k^t)$ is violated, the updating rule of (5) on w_k ensures that $\mathcal{L}(w_k^t; \mathcal{M}_k^t)$ converges to its local optimum. Therefore, after a certain number of iterations $\mathcal{L}(w_k^t; \mathcal{M}_k^t)$ would be less than $\mathcal{L}(w_k^{t-1}; \mathcal{M}_k^t)$ with high probability, satisfying the constraint again.

Then, we analyze the situation where the inequality constraint of $\mathcal{L}(w_k^t; \mathcal{M}_k^t) \leq \mathcal{L}(w_k^{t-1}; \mathcal{M}_k^t)$ is satisfied on client *k*. In this case, the PFCL objective for FedMeS can be simplified as (2). As is proved in Theorem 1 in [Li *et al.*, 2021] for (2), the following theorem holds.

Theorem 7. With Assumptions 1, 2 and 3, there exists a constant C such that for $\lambda \in \mathbb{R}$, w_k^r converges to $w_k^* :=$ arg min $_{w_k^t} G_k(w_k^t; \mathbf{w}^*)$ with rate Cg(r).

By Theorem 6, even if the violation of $\mathcal{L}(w_k^t; \mathcal{M}_k^t) \leq$ 335 $\mathcal{L}\left(w_{k}^{t-1}; \mathcal{M}_{k}^{t}\right)$ unfortunately occurs, as we only optimize on 336 \mathcal{D}_k^t and neglect \mathcal{M}_k^t under the rule of (4), the off-track w_k 337 can always be corrected to meet the constraint. By Theorem 338 7, the local model w_k^t would enjoy the same convergence rate 339 with the global model **w** with a constant multiple gap when $\mathcal{L}(w_k^t; \mathcal{M}_k^t) \leq \mathcal{L}(w_k^{t-1}; \mathcal{M}_k^t)$ is satisfied. To sum up, while 340 341 the gradient correction in (5) may occur several times during 342 the training, the corrected local model would always converge 343 to its optimum. 344

5 Experiments

5.1 Setup

Datasets and models. We select five commonly used pub-347 lic datasets: CIFAR-100 [Krizhevsky et al., 2009], EM-348 INIST (Extended-MNIST) [Cohen et al., 2017], CORe50 349 [Lomonaco and Maltoni, 2017], MiniImageNet-100 [Vinyals 350 et al., 2016] and TinyImageNet-200 [Le and Yang, 2015]. For 351 the purpose of PFCL evaluation, following the dataset split-352 ting method proposed in [Rebuffi et al., 2017] we split these 353 datasets into multiple tasks forming four cross-class datasets: 354 Split CIFAR-100: CIFAR-100 consists of 100 classes, we 355 split them into 10 tasks with 10 classes each. Split EM-356 **INIST**: We utilize 60 of the 62 categories in the original 357 dataset, and split them into 10 tasks with 6 classes each. A 358 total of 120,000 imagines are used. Split CORe50: CORe50 359 is specifically designed for assessing continual learning tech-360 niques and has 50 objects collected in 11 different sessions. 361 We naturally split it into 11 tasks with 50 classes each. Split 362 MiniImageNet: MiniImageNet-100 is commonly used in 363 few-shot learning benchmarks, which consists of 50,000 data 364 points and 10,000 testing points from 100 classes. We split 365 this dataset into 10 tasks with 10 classes each. 366

345

346

Besides, we also design the cross-domain datasets to 367 evaluate the cross domain performance for FedMeS. Fu-368 sion Tasks-40: This benchmark combines images from 369 three distinct datasets: CIFAR-100, MiniImageNet-100, and 370 TinyImageNet-200, resulting in a total of 400 classes. These 371 classes are then divided into 40 non-IID tasks, with each task 372 comprising 10 disjoint classes from the other tasks. This 373 dataset is substantial, with 200,000 images from the three het-374 erogeneous datasets. 375

We use 6-layer CNNs for the Split CIFAR-100 and Split CORe50, 2-layer CNNs for the Split EMINIST, and ResNet-18 [He *et al.*, 2016] for Split MiniImageNet and Fusion Tasks-40. For the task and dataset distributions, each client is assigned a unique task sequence, in which each task consists of randomly selected subset of 2-5 classes, with the goal of ensuring data heterogeneity.

Table 1: Acc_ALL(Acc) and average forgetting rate (FR) over all clients and all learned tasks.

DATASETS	SPLIT CIFAR-100				SPLIT EMNIST			
CLIENT NUMBER	10		20		10		20	
METHODS	ACC	FR	ACC	FR	ACC	FR	ACC	FR
FEDAVG	$.249 \pm .02$	$.35 \pm .05$	$.263 \pm .03$	$.29 \pm .03$	$.450 \pm .02$	$.49 \pm .02$	$.465 \pm .03$	$.48 \pm .02$
Ditto	$.219 \pm .02$	$.37 \pm .03$	$.221 \pm .02$	$.38 \pm .06$	$.388 \pm .03$	$.71 \pm .04$	$.381 \pm .02$	$.72 \pm .03$
FedRep	$.415 \pm .04$	$.13 \pm .02$	$.425 \pm .06$	$.13 \pm .02$	$.723 \pm .04$	$.23 \pm .04$	$.723 \pm .06$	$.24 \pm .06$
FEDAGEM	$.351 \pm .04$	$.14 \pm .03$	$.398 \pm .05$	$.14 \pm .03$	$.817 \pm .04$	$.09 \pm .05$	$.828 \pm .04$	$.16 \pm .04$
FEDWEIT	$.421 \pm .05$	$.06 \pm .03$	$.432 \pm .05$	$.08 \pm .03$	$.867 \pm .03$	$.03 \pm .01$	$.857 \pm .03$	$.03 \pm .02$
FEDMES	$.530 \pm .05$.06 ±.01	.533 ±.04	$.06 \pm .02$.935 ±.01	$.01 \pm .01$.964 ±.01	$\textbf{.01} \pm \textbf{.01}$
DATASETS	SPLIT CORE50				SPLIT MINIIMAGENET			
CLIENT NUMBER	10		20		10		20	
METHODS	ACC	FR	ACC	FR	ACC	FR	ACC	FR
FEDAVG	$.303 \pm .01$	$.60 \pm .02$	$.311 \pm .01$	$.59 \pm .02$	$.271 \pm .02$	$.51 \pm .04$	$.262 \pm .02$	$.39 \pm .05$
Ditto	$.266 \pm .01$	$.79 \pm .02$	$.267 \pm .01$	$.81 \pm .03$	$.264 \pm .03$	$.49 \pm .05$	$.265 \pm .01$	$.51 \pm .04$
FedRep	$.547 \pm .03$	$.34 \pm .02$	$.551 \pm .04$	$.35 \pm .03$	$.410 \pm .03$	$.30 \pm .05$	$.388 \pm .03$	$.24 \pm .02$
FEDAGEM	$.731 \pm .04$	$.18 \pm .03$	$.741 \pm .04$	$.20 \pm .03$	$.504 \pm .05$	$.18 \pm .03$	$.477 \pm .05$	$.17 \pm .05$
FEDWEIT	$.595 \pm .04$	$1.17 \pm .04$	$.589 \pm .05$	$.19 \pm .04$	$.319 \pm .04$	$.17 \pm .03$	$.343 \pm .04$	$.15 \pm .03$
FEDMES	.877 ±.04	$.04 \pm .01$.891 ±.04	$.03 \pm .01$	$.645 \pm .05$.08 ±.02	.659 ±.03	$\textbf{.08} \pm \textbf{.02}$

Metrics. There are mainly two kinds of metrics considered in this paper following [Chaudhry *et al.*, 2018].

• Average Accuracy: We apply four different kinds of 385 average accuracy to evaluate the performance, we define 386 the averaged accuracy of client k among all learned t387 tasks after the training of task t: $A_{k,t} = \frac{1}{t} \sum_{i=1}^{t} a_{t,i}^k$ as accuracy of client k at task t, where $a_{t,i}^k$ (i < t) is the 388 389 test accuracy of task i after the training of task t in client 390 k; averaged accuracy of client k after training all T tasks $Acc_Client_k = \frac{1}{T}\sum_{i=1}^{T}A_{k,i}$; average accuracy among all m clients at t-th task: $Acc_Task_t = \frac{1}{m}\sum_{j=1}^{m}A_{j,t}$; 391 392 393 average accuracy among all m clients in all learned T394 tasks after completing the training process of all tasks: 395 $Acc_ALL = \frac{1}{m} \frac{1}{T} \sum_{j=1}^{m} \sum_{i=1}^{T} A_{j,t}.$ 396

• Forgetting rate: The forgetting rate is the averaged disparity between minimum task accuracy during continuous training, it can measure the performance preventing catastrophic forgetting. For the forgetting rate F_t^k of client k at t-th task, it is defined as $F_t^k = \frac{1}{t-1}\sum_{i=1}^{t-1} \max_{j \in \{1,...,t-1\}} (a_{j,i}^k - a_{t,i}^k)$.

Baselines. Since there is no particular algorithm for PFCL
 problems, we compare our proposed FedMeS with other per sonalized FL and FCL techniques.

- FedAvg: A classical FL method which the server aggregates the models for all clients according to a weighted averaging of model parameters in each clients.
- Ditto: A simple personalized FL method that utilizes a regularization term addressing the accuracy, robustness and fairness in FL while optimizing communication efficiency.
- FedRep: A personalized FL method that learns a divided model with global representation and personalized heads. Only the global representation is communicated between the server and clients, while each client adapts its personalized head locally.
- FedAGEM: This can be seen as a simple federated continual learning method that combines the conventional

A-GEM method with FedAvg, achieved by applying 420 A-GEM as the local training process on the client side. 421

• FedWeIT: state-of-the-art FCL approach based on pa-422 rameter isolation, which uses masks to divide the model 423 parameter into base parameters and task-adaptive pa-424 rameters. The server averages the base parameters and 425 broadcasts the task-adaptive parameters from all clients. 426 Each client then trains all the task-adaptive parameters 427 with the new task's weights based on a regularized ob-428 jective. 429

All the experiments were conducted using PyTorch version 430 1.9 on a single machine equipped with two Intel Xeon 6226R 431 CPUs, 384GB of memory, and four NVIDIA 3090 GPUs. 432 The operating system utilized was Ubuntu 20.04.4. Each experiment is repeated for 5 times. The averages and standard 434 deviations of the above metrics are reported. 435

436

437

5.2 Results

Cross-class Performance

Tables 1 presents Acc_ALL and forgetting rate in four cross-438 class datasets. For every cross-class datasets, our proposed 439 FedMeS method outperforms all baselines in terms of aver-440 age accuracy forgetting rate. It is observed that FedWeIT ex-441 perienced a significant decline in performance when applied 442 to the Split MiniImageNet. This is primarily due to its re-443 quirement of modifying the model structure to decompose the 444 model parameters individually. Specifically, the downsample 445 layers in ResNet-18 contain a relatively small number of es-446 sential parameters, and decomposing these layers negatively 447 impacted the model's accuracy. Additionally, this modifica-448 tion process significantly increases the complexity of imple-449 mentation. In contrast, FedMeS does not require such mod-450 ifications, thus highlighting another advantage of FedMeS 451 in terms of efficient implementation. Figure 3 and Figure 4 452 respectively presents the Acc_Task and Acc_Client for the 453 Split CIFAR-100 datasets. For every task FedMeS achieves 454 highest Acc_Task and lowest forgetting rate, and for each 455 client FedMeS achieves the best Acc_Client. According to 456 these results, we further make the following observations. 457

Catastrophic forgetting. Shown in Figure 3, catastrophic 458 forgetting causes serious limitation for FedAvg and Ditto, 458

as they do not incorporate previous task information in train-460 ing. As a result, their model accuracies are inferior to other 461 methods with much higher forgetting rates. FedRep ex-462 hibits certain robustness against heterogeneity over tasks and 463 clients. However, without a designed mechanism to address 464 catastrophic forgetting, it still subjects to gradual decay in 465 average accuracy as new tasks arrive. The isolation method 466 of FedWeIT to obtain adaptive weights on the clients can-467 not well maintain the knowledge from the previous tasks, re-468 sulting in a lower accuracy than FedMeS in every dataset. 469 470 FedMeS is less affected by catastrophic forgetting due to its use of episodic memory to replay knowledge from previous 471 472 tasks.

Client drift. FedAGEM and FedWeIT fail to effectively 473 address data heterogeneity, resulting in inferior model perfor-474 mance compared to FedMeS. FedWeIT relies on the stored 475 knowledge of all tasks at the server, which may dilute the im-476 pact of individual tasks of each client. In contrast, FedMeS 477 achieves the highest accuracy in all settings thanks to its reli-478 able personalization mechanism and local inference process. 479 The advantage of FedMeS in adapting client drift is more ev-480 ident from Figure 4, where FedMeS is shown to achieve the 481 highest accuracy performance for all clients. Also, FedMeS 482 has the narrowest shade area over all clients, indicating its 483 ability to obtain consistent performance across all clients and 484 all tasks. 485



Figure 4: (a) and (b) Average accuracy of client x among all learned tasks on Split CIFAR-100 with 10 clients and 20 clients. (c) and (d) Average accuracy of client x among all learned tasks on Split MiniImageNet with 10 clients and 20 clients. The shade area is the accuracy range of tasks in each client.

486 Cross-domain Performance

Figure 5 presents the averaged test accuracy of each client 487 across all tasks in the Fusion Tasks-40 dataset. The results 488 indicate that the proposed FedMeS method outperforms the 489 FedWeIT method, with higher average accuracy and lower 490 forgetting rate among all clients and tasks. The poor per-491 formance of certain clients has a significant impact on other 492 clients, and the isolation method employed by FedWeIT to 493 mitigate catastrophic forgetting proves to be ineffective in this 494 experiment with a large number of tasks. In contrast, the pro-495 posed FedMeS method demonstrates consistent performance 496 across all tasks and clients, providing strong evidence for 497 498 its effectiveness in addressing the catastrophic forgetting and client drift issues in the PFCL problem. 499



Figure 5: Averaged test accuracy of each client among all learned t tasks after the training of task t on Fusion Tasks-40.

500

6 Conclusion

This paper has presented an in-depth examination of the chal-501 lenges associated with catastrophic forgetting and client drift 502 in PFCL and proposed the FedMeS framework as a solution 503 to these issues. FedMeS levearges a small reference memory 504 in the local training process to replay knowledge from pre-505 vious tasks to alleviate forgetting; and the same memory is 506 also used for the inference process by applying KNN-based 507 Gaussian inference to further improve model personalization 508 capability. We thoroughly analyzed the convergence behav-509 ior of FedMeS, and performed extensive experiments over 510 various PFCL tasks. For all experiments, FedMeS uniformly 511 outperforms existing techniques in terms of prediction accu-512 racy and forgetting rate. 513

514 Acknowledgement

This work is in part supported by the National Nature 515 Science Foundation of China (NSFC) Grant 62106057, 516 Guangzhou Municipal Science and Technology Guangzhou-517 HKUST(GZ) Joint Project 2023A03J0151 and Project 518 2023A03J0011, Foshan HKUST Projects FSUST20-519 FYTRI04B, and Guangdong Provincial Key Lab of 520 Integrated Communication, Sensing and Computation for 521 Ubiquitous Internet of Things. 522

523 **References**

- [Achituve *et al.*, 2021] Idan Achituve, Aviv Shamsian, Aviv
 Navon, Gal Chechik, and Ethan Fetaya. Personalized fed-
- erated learning with gaussian processes. Advances in Neu-
- *ral Information Processing Systems*, 34:8392–8406, 2021.
- [Bottou *et al.*, 2018] Léon Bottou, Frank E Curtis, and Jorge
 Nocedal. Optimization methods for large-scale machine
 learning. *Siam Review*, 60(2):223–311, 2018.
- ⁵³¹ [Casado *et al.*, 2020] Fernando E Casado, Dylan Lema,
 ⁵³² Roberto Iglesias, Carlos V Regueiro, and Senén Barro.
- Federated and continual learning for classification tasks
 in a society of devices. *arXiv preprint arXiv:2006.07129*,
 2020.
- [Chaudhry *et al.*, 2018] Arslan Chaudhry, Marc'Aurelio
 Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny.
 Efficient lifelong learning with a-gem. *arXiv preprint arXiv:1812.00420*, 2018.
- [Cheng *et al.*, 2021] Gary Cheng, Karan Chadha, and John
 Duchi. Fine-tuning is fine in federated learning. *arXiv preprint arXiv:2108.07313*, 2021.
- [Cohen *et al.*, 2017] Gregory Cohen, Saeed Afshar, Jonathan
 Tapson, and Andre Van Schaik. Emnist: Extending mnist
 to handwritten letters. In 2017 international joint con-*ference on neural networks (IJCNN)*, pages 2921–2926.
 IEEE, 2017.
- [De Lange *et al.*, 2021] Matthias De Lange, Rahaf Aljundi,
 Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. A continual learning
 survey: Defying forgetting in classification tasks. *IEEE transactions on pattern analysis and machine intelligence*,
 44(7):3366–3385, 2021.
- [Delange *et al.*, 2021] Matthias Delange, Rahaf Aljundi,
 Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg
 Slabaugh, and Tinne Tuytelaars. A continual learning
 survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelli-*gence, 2021.
- ⁵⁶⁰ [Fallah *et al.*, 2020] Alireza Fallah, Aryan Mokhtari, and
 ⁵⁶¹ Asuman Ozdaglar. Personalized federated learning with
 ⁵⁶² theoretical guarantees: A model-agnostic meta-learning
 ⁵⁶³ approach. Advances in Neural Information Processing
 ⁵⁶⁴ Systems, 33:3557–3568, 2020.
- [He *et al.*, 2016] Kaiming He, Xiangyu Zhang, Shaoqing
 Ren, and Jian Sun. Deep residual learning for image recog nition. In *Proceedings of the IEEE conference on computer*
- vision and pattern recognition, pages 770–778, 2016.

- [Huang et al., 2021] Yutao Huang, Lingyang Chu, Zirui
 Zhou, Lanjun Wang, Jiangchuan Liu, Jian Pei, and Yong
 Zhang. Personalized cross-silo federated learning on non iid data. In AAAI, pages 7865–7873, 2021.
- [Isele and Cosgun, 2018] David Isele and Akansel Cosgun.573Selective experience replay for lifelong learning. In Pro-
ceedings of the AAAI Conference on Artificial Intelligence,
volume 32, 2018.576
- [Khandelwal *et al.*, 2019] Urvashi Khandelwal, Omer Levy,
 Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Generalization through memorization: Nearest neighbor language models. In *International Conference on Learning Representations*, 2019.
- [Khodak et al., 2019] Mikhail Khodak, Maria-Florina F Balcan, and Ameet S Talwalkar. Adaptive gradient-based meta-learning methods. Advances in Neural Information Processing Systems, 32, 2019.
- [Kirkpatrick *et al.*, 2017] 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.
- [Krizhevsky *et al.*, 2009] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 593
- [Le and Yang, 2015] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015. 597
- [Li *et al.*, 2021] Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning through personalization. In *International Conference on Machine Learning*, pages 6357–6368. PMLR, 2021. 601
- [Lomonaco and Maltoni, 2017] Vincenzo Lomonaco and 602 Davide Maltoni. Core50: a new dataset and benchmark for continuous object recognition. In *Conference on Robot Learning*, pages 17–26. PMLR, 2017. 605
- [Lopez-Paz and Ranzato, 2017] David Lopez-Paz and 606 Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. Advances in neural information 608 processing systems, 30, 2017. 609
- [Marfoq *et al.*, 2021] Othmane Marfoq, Giovanni Neglia, Aurélien Bellet, Laetitia Kameni, and Richard Vidal. Federated multi-task learning under a mixture of distributions. *Advances in Neural Information Processing Systems*, 34:15434–15447, 2021. 614
- [Marfoq et al., 2022] Othmane Marfoq, Giovanni Neglia, 615
 Richard Vidal, and Laetitia Kameni. Personalized federated learning through local memorization. In *International* 617
 Conference on Machine Learning, pages 15070–15092. 618
 PMLR, 2022. 619
- [McMahan *et al.*, 2017] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks

- from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- [Rebuffi *et al.*, 2017] Sylvestre-Alvise Rebuffi, Alexander
 Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl:
- 627 Incremental classifier and representation learning. In Pro-
- ceedings of the IEEE conference on Computer Vision and
 Pattern Recognition, pages 2001–2010, 2017.
- [Riemer *et al.*, 2018] Matthew Riemer, Ignacio Cases,
 Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and
 Gerald Tesauro. Learning to learn without forgetting by
 maximizing transfer and minimizing interference. *arXiv preprint arXiv:1810.11910*, 2018.
- [Shen *et al.*, 2020] Tao Shen, Jie Zhang, Xinkang Jia,
 Fengda Zhang, Gang Huang, Pan Zhou, Kun Kuang, Fei
 Wu, and Chao Wu. Federated mutual learning. *arXiv*
- *preprint arXiv:2006.16765*, 2020.
 [Shi *et al.*, 2021] Yujun Shi, Li Yuan, Yunpeng Chen, and
- Jiashi Feng. Continual learning via bit-level informa tion preserving. In *Proceedings of the IEEE/CVF Confer-*
- ence on Computer Vision and Pattern Recognition, pages
 16674–16683, 2021.
- ⁶⁴⁴ [Shim *et al.*, 2021] Dongsub Shim, Zheda Mai, Jihwan ⁶⁴⁵ Jeong, Scott Sanner, Hyunwoo Kim, and Jongseong Jang.
- 646 Online class-incremental continual learning with adversar-

ial shapley value. In Proceedings of the AAAI Conference

- *on Artificial Intelligence*, volume 35, pages 9630–9638, 2021.
- [Shoham *et al.*, 2019] 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.
- [Usmanova *et al.*, 2021] Anastasiia Usmanova, François
 Portet, Philippe Lalanda, and German Vega. A distillationbased approach integrating continual learning and federated learning for pervasive services. *arXiv preprint arXiv:2109.04197*, 2021.
- ⁶⁶⁰ [Venkatesha *et al.*, 2022] Yeshwanth Venkatesha, Youngeun
 ⁶⁶¹ Kim, Hyoungseob Park, Yuhang Li, and Priyadarshini
 ⁶⁶² Panda. Addressing client drift in federated continual
 ⁶⁶³ learning with adaptive optimization. Available at SSRN
 ⁶⁶⁴ 4188586, 2022.
- [Vinyals *et al.*, 2016] Oriol Vinyals, Charles Blundell, Timo thy Lillicrap, Daan Wierstra, et al. Matching networks for
 one shot learning. *Advances in neural information pro-*
- cessing systems, 29, 2016.
- ⁶⁶⁹ [Wang *et al.*, 2022] Liyuan Wang, Xingxing Zhang, Kuo
 ⁶⁷⁰ Yang, Longhui Yu, Chongxuan Li, Lanqing Hong, Shifeng
- ⁶⁷¹ Zhang, Zhenguo Li, Yi Zhong, and Jun Zhu. Memory re-
- play with data compression for continual learning. *arXiv preprint arXiv:2202.06592*, 2022.
- 674 [Yao and Sun, 2020] Xin Yao and Lifeng Sun. Continual lo-
- cal training for better initialization of federated models. In
- 676 2020 IEEE International Conference on Image Processing
- 677 (ICIP), pages 1736–1740. IEEE, 2020.

- [Yoon *et al.*, 2021] Jaehong Yoon, Wonyong Jeong, Giwoong Lee, Eunho Yang, and Sung Ju Hwang. Federated continual learning with weighted inter-client transfer. In *International Conference on Machine Learning*, pages 12073–12086. PMLR, 2021. 682
- [Yu et al., 2020a] Lu Yu, Bartlomiej Twardowski, Xialei683Liu, Luis Herranz, Kai Wang, Yongmei Cheng, Shangling684Jui, and Joost van de Weijer. Semantic drift compensa-
tion for class-incremental learning. In Proceedings of the
IEEE/CVF Conference on Computer Vision and Pattern
Recognition, pages 6982–6991, 2020.688
- [Yu *et al.*, 2020b] Tao Yu, Eugene Bagdasaryan, and Vitaly Shmatikov. Salvaging federated learning by local adaptation. *arXiv preprint arXiv:2002.04758*, 2020. 691
- [Zhang et al., 2022] Lin Zhang, Li Shen, Liang Ding, 692
 Dacheng Tao, and Ling-Yu Duan. Fine-tuning global 693
 model via data-free knowledge distillation for non-iid federated learning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 10174–10183, 2022. 697