
Federated Virtual Learning on Heterogeneous Data with Local-global Distillation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Despite Federated Learning (FL)’s trend for learning machine learning models in a
2 distributed manner, it is susceptible to performance drops when training on hetero-
3 geneous data. In addition, FL inevitably faces the challenges of synchronization,
4 efficiency, and privacy. Recently, dataset distillation has been explored in order to
5 improve the efficiency and scalability of FL by creating a smaller, synthetic dataset
6 that retains the performance of a model trained on the local private datasets. *We*
7 *discover that using distilled local datasets can amplify the heterogeneity issue in*
8 *FL*. To address this, we propose a new method, called **F**ederated Virtual Learning
9 on Heterogeneous Data with **L**ocal-**G**lobal **D**istillation (FEDLGD), which trains
10 FL using a smaller synthetic dataset (referred as *virtual data*) created through a
11 combination of local and global dataset distillation. Specifically, to handle syn-
12 chronization and class imbalance, we propose iterative distribution matching to
13 allow clients to have the same amount of balanced *local virtual data*; to harmonize
14 the domain shifts, we use federated gradient matching to distill *global virtual data*
15 that are shared with clients without hindering data privacy to rectify heterogeneous
16 local training via enforcing local-global feature similarity. We experiment on both
17 benchmark and real-world datasets that contain heterogeneous data from different
18 sources, and further scale up to an FL scenario that contains large number of
19 clients with heterogeneous and class imbalance data. Our method outperforms
20 *state-of-the-art* heterogeneous FL algorithms under various settings with a very
21 limited amount of distilled virtual data.

22 1 Introduction

23 Federated Learning (FL) [29] has become a popular solution for different institutions to collaboratively
24 train machine learning models without pooling private data together. Typically, it involves a central
25 server and multiple local clients; then the model is trained via aggregation of local network parameter
26 updates on the server side iteratively. FL is widely accepted in many areas, such as computer vision,
27 natural language processing, and medical image analysis [25, 12, 41].

28 On the one hand, clients with different amounts of data cause asynchronization and affect the efficiency
29 of FL systems. Dataset distillation [39, 5, 46, 44, 45] addresses the issue by only summarizing smaller
30 synthetic datasets from the private local datasets to ensure each client owns the same amount of
31 data. We refer this underexplored strategy as *federated virtual learning*, as the models are trained
32 from synthetic data [40, 10, 16]. These methods have been found to perform better than model-
33 synchronization-based FL approaches while requiring fewer server-client interactions.

34 On the other hand, due to different data collection protocols, data from different clients inevitably
35 face heterogeneity problems with domain shift, which means data may not be independent and
36 identically distributed (iid) among clients. Heterogeneous data distribution among clients becomes a

37 key challenge in FL, as aggregating model parameters from non-iid feature distributions suffers from
 38 client drift [18] and diverges the global model update [26].

39 We observe that using locally distilled datasets can amplify the heterogeneity issue. Figure 1 shows
 40 the tSNE plots of two different datasets, USPS [31] and SynthDigits [9], each considered as a client.
 41 tSNE takes the original and distilled virtual images as input and embeds them into 2D planes. One
 42 can observe that the distribution becomes diverse after distillation.

43 To alleviate the problem of data heterogeneity
 44 in classical FL settings, two main orthogonal
 45 approaches can be taken. *Approach 1* aims
 46 to minimize the difference between the local
 47 and global model parameters to improve conver-
 48 gence [25, 18, 38]. *Approach 2* enforces con-
 49 sistency in local embedded features using an-
 50 chors and regularization loss [37, 47, 42]. The
 51 first approach can be easily applied to distilled
 52 local datasets, while the second approach has
 53 limitations when adapting to federated virtual
 54 learning. Specifically, VHL [37] samples global
 55 anchors from untrained StyleGAN [19] suffers
 56 performance drop when handling amplified het-
 57 erogeneity after dataset distillation. Other meth-
 58 ods, such as those that rely on external global
 59 data [47], or feature sharing from clients [42],
 60 are less practical, as they pose greater data privacy risks compared to classical FL settings¹. *Without
 61 hindering data privacy*, developing strategies following *approach 2* for federated virtual learning on
 62 heterogeneous data remains open questions on 1) *how to set up global anchors for locally distilled
 63 datasets and 2) how to select the proper regularization loss(es)*.

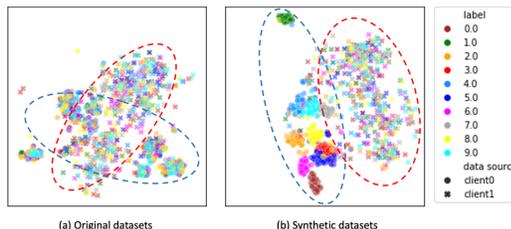


Figure 1: Distilled local datasets can worsen heterogeneity in FL. tSNE plots of (a) original datasets and (b) distilled virtual datasets of USPS (client 0) and SynthDigits (client 1). The two distributions are marked in red and blue. We observe fewer overlapped \circ and \times in (b) compared with (a), indicating higher heterogeneity between two clients after distillation.

64 To this end, we propose FEDLGD, a federated virtual learning method with local and global dis-
 65 tillation. We propose *iterative distribution matching* in local distillation by comparing the feature
 66 distribution of real and synthetic data using an evolving feature extractor. The local distillation results
 67 in smaller sets with balanced class distributions, achieving efficiency and synchronization while
 68 avoiding class imbalance. FEDLGD updates the local model on local distilled synthetic datasets
 69 (named *local virtual data*). We found that training FL with local virtual data can exacerbate hetero-
 70 geneity in feature space if clients’ data has domain shift (Figure. 1). Therefore, unlike previously
 71 proposed federated virtual learning methods that rely solely on local distillation [10, 40, 16], we also
 72 propose a novel and efficient method, *federated gradient matching*, that integrated well with FL to
 73 distill global virtual data as anchors on the server side. This approach aims to alleviate domain shifts
 74 among clients by promoting similarity between local and global features. Note that we only share
 75 local model parameters w.r.t. distilled data. Thus, the privacy of local original data is preserved. We
 76 conclude our contributions as follows:

- 77 • This paper focuses on an important but underexplored FL setting in which local models
 78 are trained on small distilled datasets, which we refer to as *federated virtual learning*. We
 79 design two effective and efficient dataset distillation methods for FL.
- 80 • We are *the first* to reveal that when datasets are distilled from clients’ data with domain shift,
 81 the heterogeneity problem can be *exacerbated* in the federated virtual learning setting.
- 82 • We propose to address the heterogeneity problem by mapping clients to similar features
 83 regularized by gradually updated global virtual data using averaged client gradients.
- 84 • Through comprehensive experiments on benchmark and real-world datasets, we show that
 85 FEDLGD outperforms existing state-of-the-art FL algorithms.

¹Note that FedFA [47], and FedFM [42] are unpublished works proposed concurrently with our work

86 2 Related Work

87 2.1 Dataset Distillation

88 Data distillation aims to improve data efficiency by distilling the most essential feature in a large-
89 scale dataset (e.g., datasets comprising billions of data points) into a certain terse and high-fidelity
90 dataset. For example, Gradient Matching [46] is proposed to make the deep neural network produce
91 similar gradients for both the terse synthetic images and the original large-scale dataset. Besides,
92 [5] proposes matching the model training trajectory between real and synthetic data to guide the
93 update for distillation. Another popular way of conducting data distillation is through Distribution
94 Matching [45]. This strategy instead, attempts to match the distribution of the smaller synthetic
95 dataset with the original large-scale dataset. It significantly improves the distillation efficiency.
96 Moreover, recent studies have justified that data distillation also preserves privacy [7, 4], which is
97 critical in federated learning. In practice, dataset distillation is used in healthcare for medical data
98 sharing for privacy protection [22]. Other modern data distillation strategies can be found here [33].

99 2.2 Heterogeneous Federated Learning

100 FL performance downgrading on non-iid data is a critical challenge. A variety of FL algorithms have
101 been proposed ranging from global aggregation to local optimization to handle this heterogeneous
102 issue. *Global aggregation* improves the global model exchange process for better unitizing the
103 updated client models to create a powerful server model. FedNova [38] notices an imbalance
104 among different local models caused by different levels of training stage (e.g., certain clients train
105 more epochs than others) and tackles such imbalance by normalizing and scaling the local updates
106 accordingly. Meanwhile, FedAvgM [15] applies the momentum to server model aggregation to
107 stabilize the optimization. Furthermore, there are strategies to refine the server model from learning
108 client models such as FedDF [27] and FedFTG [43]. *Local training optimization* aims to explore the
109 local objective to tackle the non-iid issue in FL system. FedProx [25] straightly adds L_2 norm to
110 regularize the client model and previous server model. Scaffold [18] adds the variance reduction term
111 to mitigate the "clients-drift". Also, MOON [24] brings mode-level contrastive learning to maximize
112 the similarity between model representations to stable the local training. There is another line of
113 works [42, 37] proposed to use a global *anchor* to regularize local training. Global anchor can be
114 either a set of virtual global data or global virtual representations in feature space. However, in [37],
115 the empirical global anchor selection may not be suitable for data from every distribution as they
116 don't update the anchor according to the training datasets.

117 2.3 Datasets Distillation for FL

118 Dataset distillation for FL is an emerging topic that has attracted attention due to its benefit for
119 efficient FL systems. It trains model on distilled synthetic datasets, thus we refer it as federated
120 virtual learning. It can help with FL synchronization and improve training efficiency by condensing
121 every client's data into a small set. To the best of our knowledge, there are few published works on
122 distillation in FL. Concurrently with our work, some studies [10, 40, 16] distill datasets locally and
123 share the distilled datasets with other clients/servers. Although privacy is protected against *currently*
124 existing attack models, we consider sharing local distilled data a dangerous move. Furthermore, none
125 of the existing work has addressed the heterogeneity issue.

126 3 Method

127 In this section, we will describe the problem setup, introduce the key technical contributions and
128 rationale of the design for FEDLGD, and explain the overall training pipeline.

129 3.1 Setup for Federated Virtual Learning

130 We start with describing the classical FL setting. Suppose there are N parties who own local datasets
131 (D_1, \dots, D_N) , and the goal of a classical FL system, such as FedAvg [29], is to train a global model

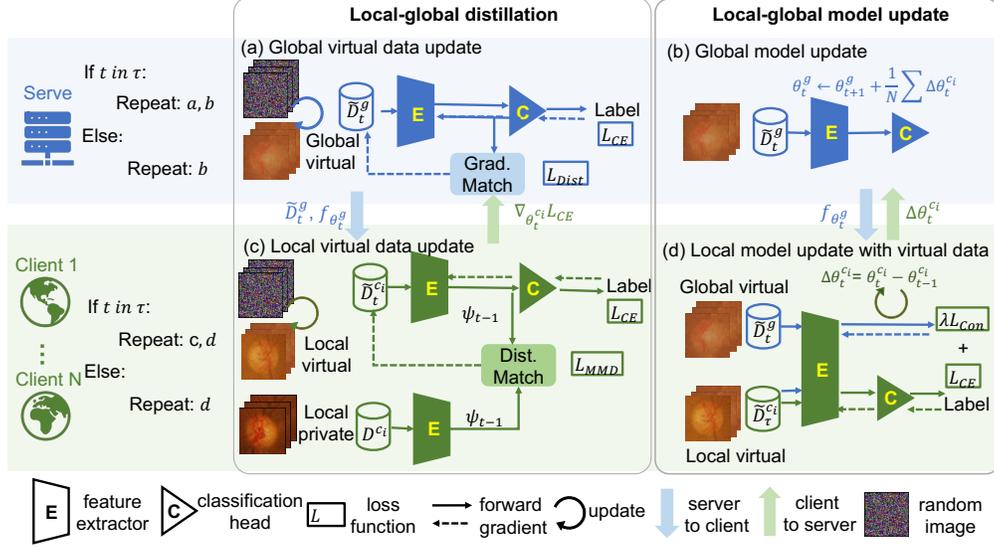


Figure 2: Overview pipeline for FEDLGD. We assume T FL rounds will be performed, among which we will define the selected distillation rounds as $\tau \in [T]$ for local-global iteration. For selected rounds ($t \in \tau$), clients will update local models (d) and refine the local virtual data with the latest network parameters (c), while the server uses aggregated gradients from cross-entropy loss (\mathcal{L}_{CE}) to update global virtual data (a) and update the global model (b). We term this procedure Iterative Local-global Distillation. For the unselected rounds ($t \in T \setminus \tau$), we perform ordinary FL pipeline on local virtual data with regularization loss (\mathcal{L}_{Con}) on global virtual data.

132 with parameters θ on the distributed datasets ($\tilde{D} \equiv \bigcup_{i \in [N]} D_i$). The objective function is written as:

$$\mathcal{L}(\theta) = \sum_{i=1}^N \frac{|D_i|}{|\tilde{D}|} \mathcal{L}_i(\theta), \quad (1)$$

133 where $\mathcal{L}_i(w)$ is the empirical loss of client i .

134 In practice, different clients in FL may have variant amounts of training samples, leading to asynchro-
 135 nized updates. In this work, we focus on a new type of FL training method – federated virtual learning,
 136 that trains on distilled datasets for efficiency and synchronization (discussed in Sec 2.3). Federated
 137 virtual learning synthesizes local virtual data \tilde{D}_i for client i for $i \in [N]$ and form $\tilde{D} \equiv \bigcup_{i \in [N]} \tilde{D}_i$.
 138 Typically, $|\tilde{D}_i| \ll |D_i|$ and $|\tilde{D}_i| = |\tilde{D}_j|$. A basic setup for federated virtual learning is to replace D_i
 139 with \tilde{D}_i in Eq (1), namely FL model is trained on the virtual datasets. As suggested in FedDM [40],
 140 the clients should not share gradients w.r.t. the original data for privacy concern.

141 3.2 Overall Pipeline

142 The overall pipeline of our proposed method contains three stages, including 1) *initialization*, 2)
 143 *iterative local-global distillation*, and 3) *federated virtual learning*. We depict the overview of
 144 FEDLGD pipeline in Figure 2. However, FL is inevitably affected by several challenges, including
 145 synchronization, efficiency, privacy, and heterogeneity. Specifically, we outline FEDLGD as follows:

146 We begin with the initialization of the clients' local virtual data \tilde{D}^c by performing initial rounds of
 147 distribution matching (DM) [45]. Meanwhile, the server will initialize global virtual data \tilde{D}^g
 148 and network parameters θ_0^g . In this stage, we generate the same amount of class-balanced virtual data for
 149 each client and server.

150 Then, we will refine our local and global virtual data using our proposed *local-global* distillation
 151 strategies in Sec. 3.3.1 and 3.3.2. This step is performed for a few selected iterations (e.g. $\tau =$
 152 $\{0, 5, 10\}$) to update θ using \mathcal{L}_{CE} (Eq 3), \tilde{D}^g using \mathcal{L}_{Dist} (Eq 5), and \tilde{D}^c using \mathcal{L}_{MMD} (Eq 2) in early
 153 training epochs. For each selected iterations, the server and clients will update their virtual data for a
 154 few distillation steps.

155 Finally, after refining local and global virtual data \tilde{D}^g and \tilde{D}^c , we continue federated virtual learning
 156 in stage 3 on local virtual data \tilde{D}^c using $\mathcal{L}_{\text{total}}$ (Eq. 3), with \tilde{D}^g as regularization anchor to calculate
 157 \mathcal{L}_{Con} (Eq. 4). We provide implementation details, an algorithm box, and an anonymous link to our
 158 code in the Appendix.

159 3.3 FL with Local-Global Dataset Distillation

160 3.3.1 Local Data Distillation

161 Our purpose is to decrease the number of local data to achieve efficient training to meet the following
 162 goals. First of all, we hope to synthesize virtual data conditional on class labels to achieve class-
 163 balanced virtual datasets. Second, we hope to distill local data that is best suited for the classification
 164 task. Last but not least, the process should be efficient due to the limited computational resource
 165 locally. To this end, we design Iterative Distribution Matching to fulfill our purpose.

166 **Iterative distribution matching.** We aim to gradually improve distillation quality during FL training.
 167 To begin with, we split a model into two parts, feature extractor ψ (shown as E in figure 2) and
 168 classification head h (shown as C in figure 2). The whole classification model is defined as $f^\theta = h \circ \psi$.
 169 The high-level idea of distribution matching can be described as follows. Given a feature extractor
 170 $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$, we want to generate \tilde{D} so that $P_\psi(D) \approx P_\psi(\tilde{D})$ where P is the distribution in
 171 feature space. To distill local data during FL efficiently that best fits our task, we intend to use
 172 the up-to-date server model’s feature extractor as our kernel function to distill better virtual data.
 173 Since we can’t obtain ground truth distribution of local data, we utilize empirical maximum mean
 174 discrepancy (MMD) [11] as our loss function for local virtual distillation:

$$\mathcal{L}_{\text{MMD}} = \sum_k^K \left\| \frac{1}{|D_k^c|} \sum_{i=1}^{|D_k^c|} \psi^t(x_i) - \frac{1}{|\tilde{D}_k^{c,t}|} \sum_{j=1}^{|\tilde{D}_k^{c,t}|} \psi^t(\tilde{x}_j^t) \right\|^2, \quad (2)$$

175 where ψ^t and $\tilde{D}^{c,t}$ are the server feature extractor and local virtual data from the latest global
 176 t . Following [46, 45], we apply the differentiable Siamese augmentation on virtual data \tilde{D}^c . K is the
 177 total number of classes, and we sum over MMD loss calculated per class $k \in [K]$. In such a way, we
 178 can generate balanced local virtual data by optimizing the same number of virtual data per class.

179 Although such an efficient distillation strategy is inspired by DM [45], we highlight the key difference
 180 that DM uses randomly initialized deep neural networks to extract features, whereas we use trained
 181 FL models with task-specific supervised loss. We believe *iterative updating* on the clients’ data using
 182 the up-to-date network parameters can generate better task-specific local virtual data. Our intuition
 183 comes from the recent success of the empirical neural tangent kernel for data distribution learning and
 184 matching [30, 8]. Especially, the feature extractor of the model trained with FEDLGD could obtain
 185 feature information from other clients, which further harmonize the domain shift between clients.
 186 We apply DM [45] to the baseline FL methods and demonstrate the effectiveness of our proposed
 187 iterative strategy in Sec. 4. Furthermore, note that FEDLGD only requires a few hundreds of local
 188 distillations steps using the local model’s feature distribution, which is more computationally efficient
 189 than other bi-level dataset distillation methods [46, 5].

190 **Harmonizing local heterogeneity with global anchors.** Data collected in different sites may have
 191 different distributions due to different collecting protocols and populations. Such heterogeneity will
 192 degrade the performance of FL. Worse yet, we found increased data heterogeneity among clients
 193 when federatively training with distilled local virtual data (see Figure 1). We aim to alleviate the
 194 dataset shift by adding a regularization term in feature space to our total loss function for local model
 195 updating, which is inspired by [37, 20]:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{CE}}(\tilde{D}^g, \tilde{D}^c; \theta) + \lambda \mathcal{L}_{\text{Con}}(\tilde{D}^g, \tilde{D}^c), \quad (3)$$

196 and

$$\mathcal{L}_{\text{Con}} = \sum_{i \in I} \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(z_g \cdot z_p / \tau_{\text{temp}})}{\sum_{a \in A(i)} \exp(z_g \cdot z_a / \tau_{\text{temp}})}, \quad (4)$$

197 where \mathcal{L}_{CE} is the cross-entropy measured on the virtual data, and \mathcal{L}_{Con} is the supervised contrastive
 198 loss where I is the collection of all indices, $A(i)$ indicates all the local and global virtual data indices
 199 without i (i.e. $A(i) \equiv I \setminus \{i\}$), $z = \psi(x)$ is the output of feature extractor, $P(i)$ represents the set of

200 images belonging to the same class y_i without data i , and τ_{temp} is a scalar temperature parameter. In
 201 such a way, global virtual data can be served for calibration, where z_g is from \tilde{D}^g as an anchor, and
 202 z_p and z_a are from \tilde{D}^c . At this point, a critical problem arises: *What global virtual data shall we use?*

203 3.3.2 Global Data Distillation

204 Here, we provide an affirmative solution to the question of generating global virtual data that can
 205 be naturally incorporated into FL pipeline. Although distribution-based matching is efficient, local
 206 clients may not share their features due to privacy concerns. Therefore, we propose to leverage local
 207 clients’ averaged gradients to distill global virtual data and utilize it in Eq. (4). We term our global
 208 data distillation method as *Federated Gradient Matching*.

209 **Federated gradient matching.** The concept of gradient-based dataset distillation is to minimize the
 210 distance between gradients from model parameters trained by original data and distilled data. It is
 211 usually considered as a learning-to-learn problem because the procedure consists of model updates
 212 and distilled data updates. Zhao *et al.* [46] studies gradient matching in the centralized setting via
 213 bi-level optimization that iteratively optimizes the virtual data and model parameters. However, the
 214 implementation in [46] is not appropriate for our specific context because there are two fundamental
 215 differences in our settings: 1) for model updating, the gradient-distilled dataset is on the server and
 216 will not directly optimize the targeted task; 2) for virtual data update, the ‘optimal’ model comes
 217 from the optimized local model aggregation. These two steps can naturally be embedded in local
 218 model updating and global virtual data distillation from the aggregated local gradients. First, we
 219 utilize the distance loss \mathcal{L}_{Dist} [46] for gradient matching:

$$\mathcal{L}_{Dist} = Dist(\nabla_{\theta} \mathcal{L}_{CE}^{\tilde{D}^g}(\theta), \nabla_{\theta} \mathcal{L}_{CE}^{\tilde{D}^c}(\theta)) \quad (5)$$

220 where \tilde{D}^c and \tilde{D}^g denote local and global virtual data, $\nabla_{\theta} \mathcal{L}_{CE}^{\tilde{D}^c}$ is the average client gradient. Then,
 221 our proposed federated gradient matching optimize as follows:

$$\min_{\tilde{D}^g} \mathcal{L}_{Dist}(\theta) \quad \text{subject to} \quad \theta = \frac{1}{N} \theta^{c_i^*},$$

222 where $\theta^{c_i^*} = \arg \min_{\theta} \mathcal{L}_i(\tilde{D}^c)$ is the optimal local model weights of client i at a certain round t .

223 Noting that compared with FedAvg [29], there is no additional client information shared for global
 224 distillation. We also note the approach seems similar to the gradient inversion attack [49] but we
 225 consider averaged gradients w.r.t. local virtual data, and the method potentially defenses inference
 226 attack better (Appendix D.8), which is also implied by [40, 7]. Privacy preservation can be further
 227 improved by employing differential privacy [1], but this is not the main focus of our work.

228 4 Experiment

229 To evaluate FEDLGD, we consider the FL setting in which clients obtain data from different domains
 230 while performing the same task. Specifically, we compare with multiple baselines on benchmark
 231 datasets DIGITS (Sec. 4.2), where each client has data from completely different open-sourced
 232 datasets. The experiment is designed to show that FEDLGD can effectively mitigate large domain
 233 shifts. Additionally, we evaluate the performance of FEDLGD on another benchmark dataset,
 234 CIFAR10C [14], which collects data from different corruptions yielding data distribution shift and
 235 contains a large number of clients, so that we can investigate varied client sampling in FL. The
 236 experiment aims to show FEDLGD’s feasibility on large-scale FL environments. We also validate the
 237 performance under medical datasets, RETINA, in Appendix B.

238 4.1 Training and Evaluation Setup

239 **Model architecture.** We conduct the ablation study to explore the effect of different deep neural
 240 networks’ performance under FEDLGD. Specifically, we adapt ResNet18 [13] and ConvNet [46]
 241 in our study. To achieve the optimal performance, we apply the same architecture to perform both
 242 the local distillation task and the classification task, as this combination is justified to have the best
 243 output [46, 45]. The detailed model architectures are presented in Appendix D.4.

244 **Comparison methods.** We compare the performance of downstream classification tasks using state-of-
 245 the-art (SOTA) FL algorithms, FedAvg [29], FedProx [26], FedNova [38], Scaffold [18], MOON [24],

Table 1: Test accuracy for DIGITS under different images per class (IPC) and model architectures. R and C stand for ResNet18 and ConvNet, respectively, and we set IPC to 10 and 50. There are five clients (MNIST, SVHN, USPS, SynthDigits, and MNIST-M) containing data from different domains. ‘Average’ is the unweighted test accuracy average of all the clients. The best performance under different models is highlighted using **bold**. The best results on ConvNet are marked in **red** and in black for ResNet18.

DIGITS		MNIST		SVHN		USPS		SynthDigits		MNIST-M		Average	
IPC		10	50	10	50	10	50	10	50	10	50	10	50
FedAvg	R	73.0	92.5	20.5	48.9	83.0	89.7	13.6	28.0	37.8	72.3	45.6	66.3
	C	94.0	96.1	65.9	71.7	91.0	92.9	55.5	69.1	73.2	83.3	75.9	82.6
FedProx	R	72.6	92.5	19.7	48.4	81.5	90.1	13.2	27.9	37.3	67.9	44.8	65.3
	C	93.9	96.1	66.0	71.5	90.9	92.9	55.4	69.0	73.7	83.3	76.0	82.5
FedNova	R	75.5	92.3	17.3	50.6	80.3	90.1	11.4	30.5	38.3	67.9	44.6	66.3
	C	94.2	96.2	65.5	73.1	90.6	93.0	56.2	69.1	74.6	83.7	76.2	83.0
Scaffold	R	75.8	93.4	16.4	53.8	79.3	91.3	11.2	34.2	38.3	70.8	44.2	68.7
	C	94.1	96.3	64.9	73.3	90.6	93.4	56.0	70.1	74.6	84.7	76.0	83.6
MOON	R	15.5	80.4	15.9	14.2	25.0	82.4	10.0	11.5	11.0	35.4	15.5	44.8
	C	85.0	95.5	49.2	70.5	83.4	92.0	31.5	67.2	56.9	82.3	61.2	81.5
VHL	R	87.8	95.9	29.5	67.0	88.0	93.5	18.2	60.7	52.2	85.7	55.1	80.5
	C	95.0	96.9	68.6	75.2	92.2	94.4	60.7	72.3	76.1	83.7	78.5	84.5
FEDLGD	R	92.9	96.7	46.9	73.3	89.1	93.9	27.9	72.9	70.8	85.2	65.5	84.4
	C	95.8	97.1	68.2	77.3	92.4	94.6	67.4	78.5	79.4	86.1	80.6	86.7

246 and VHL [37]². We directly use local virtual data from our initialization stage for FL methods other
 247 than ours. We perform classification on client’s testing set and report the test accuracies.

248 **FL training setup.** We use the SGD optimizer with a learning rate of 10^{-2} for DIGITS and CIFAR10C.
 249 If not specified, our default setting for local model update epochs is 1, total update rounds is 100,
 250 the batch size for local training is 32, and the number of virtual data update iterations ($|\tau|$) is 10.
 251 The numbers of default virtual data distillation steps for clients and server are set to 100 and 500,
 252 respectively. Since we only have a few clients for DIGITS and RETINA experiments, we will select all
 253 the clients for each iteration, while the client selection for CIFAR10C experiments will be specified in
 254 Sec. 4.3. The experiments are run on NVIDIA GeForce RTX 3090 Graphics cards with PyTorch.

255 **Proper Initialization for Distillation.** We propose to initialize the distilled data using statistics
 256 from local data to take care of both privacy concerns and model performance. Specifically, each
 257 client calculates the statistics of its own data for each class, denoted as μ_i^c, σ_i^c , and then initializes the
 258 distillation images per class, $x \sim \mathcal{N}(\mu_i^c, \sigma_i^c)$, where c and i represent each client and categorical label.
 259 The server only needs to aggregate the statistics and initializes the virtual data as $x \sim \mathcal{N}(\mu_i^g, \sigma_i^g)$. In
 260 this way, no real data is shared with any participant in the FL system. The comparison results using
 261 different initialization methods proposed in previous works [46, 45] can be found in Appendix C.

262 4.2 DIGITS Experiment

263 **Datasets.** We use the following datasets for our benchmark experiments: DIGITS = {MNIST [21],
 264 SVHN [31], USPS [17], SynthDigits [9], MNIST-M [9]}. Each dataset in DIGITS contains hand-
 265 written, real street and synthetic digit images of 0, 1, \dots , 9. As a result, we have 5 clients in the
 266 experiments, and image size is 28×28 .

267 **Comparison with baselines under various conditions.** To validate the effectiveness of FEDLGD,
 268 we first compare it with the alternative FL methods varying on two important factors: Image-per-class
 269 (IPC) and different deep neural network architectures (arch). We use $IPC \in \{10, 50\}$ and $arch \in$
 270 $\{ResNet18(R), ConvNet(C)\}$ to examine the performance of SOTA models and FEDLGD using
 271 distilled DIGITS. Note that we fix $IPC = 10$ for global virtual data and vary IPC for local virtual data.
 272 Table 1 shows the test accuracies of DIGITS experiments. In addition to testing with original test sets,
 273 we also show the unweighted averaged test accuracy. One can observe that for each FL algorithm,
 274 ConvNet(C) always has the best performance under all IPCs. The observation is consistent with [45]
 275 as more complex architectures may cause over-fitting in training virtual data. It is also shown that
 276 using $IPC = 50$ always outperforms $IPC = 10$ as expected since more data are available for training.
 277 Overall, FEDLGD outperforms other SOTA methods, where on average accuracy, FEDLGD increases
 278 the best test accuracy results among the baseline methods of 2.1% ($IPC = 10$, $arch = C$), 10.4% (IPC

²The detailed information of the methods can be found in Appendix E.

Table 2: Averaged test accuracy for CIFAR10C with ConvNet.

CIFAR10C		FedAvg		FedProx		FedNova		Scaffold		MOON		VHL		FEDLGD	
IPC	Client ratio	10	50	10	50	10	50	10	50	10	50	10	50	10	50
	0.2	27.0	44.9	27.0	44.9	26.7	34.1	27.0	44.9	20.5	31.3	21.8	45.0	32.9	46.8
	0.5	29.8	51.4	29.8	51.4	29.6	45.9	30.6	51.6	23.8	43.2	29.3	51.7	39.5	52.8
	1	33.0	54.9	33.0	54.9	30.0	53.2	33.8	54.5	26.4	51.6	34.4	55.2	47.6	57.4

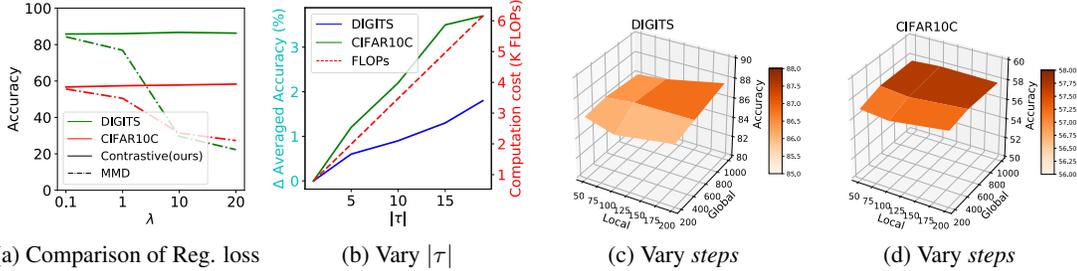


Figure 3: (a) Comparison between different regularization losses and their weighting in total loss (λ). One can observe that supervised contrastive loss gives us better and more stable performance with different coefficient choices. (b) The trade-off between $|\tau|$ and computation cost. One can observe that the model performance improves with the increasing $|\tau|$, which is a trade-off between computation cost and model performance. Vary data updating *steps* for (c) DIGITS and (d) CIFAR10C. One can observe that FEDLGD yields consistent performance, and the accuracy tends to improve with an increasing number of local and global steps.

279 =10, arch = R), 2.2% (IPC = 50, arch = C) and 3.9% (IPC =50, arch = R). VHL [37] is the closest
 280 strategy to FEDLGD and achieves the best performance among the baseline methods, indicating that
 281 the feature alignment solutions are promising for handling heterogeneity in federated virtual learning.
 282 However, VHL is still worse than FEDLGD, and the performance may result from the differences in
 283 synthesizing global virtual data. VHL [37] uses untrained StyleGAN [19] to generate global virtual
 284 data without further updating. On the contrary, we update our global virtual data during FL training.

285 4.3 CIFAR10C Experiment

286 **Datasets.** We conduct real-world FL experiments on CIFAR10C³, where, like previous studies [24],
 287 we apply Dirichlet distribution with $\alpha = 0.5$ to generate 3 partitions on each distorted Cifar10-C [14],
 288 resulting in 57 clients each with class imbalanced non-IID datasets. In addition, we apply random
 289 client selection with ratio = 0.2, 0.5, and 1 and set image size as 28×28 .

290 **Comparison with baselines under different client sampling ratios.** The objective of the experiment
 291 is to test FEDLGD under popular FL questions: class imbalance, large number of clients, different
 292 client sample ratios, and data heterogeneity. One benefit of federated virtual learning is that we can
 293 easily handle class imbalance by distilling the same number (IPC) of virtual data. We will vary IPC
 294 and fix the model architecture to ConvNet since it is validated that ConvNet yields better performance
 295 in virtual training [46, 45]. One can observe from Table 2 that FEDLGD consistently achieves the
 296 best performance under different IPC and client sampling ratios. We would like to point out that
 297 when IPC=10, the performance boosts are significant, which indicates that FEDLGD is well-suited
 298 for FL conditions when there is a large group of clients and each of them has a limited number of
 299 data.

300 4.4 Ablation studies for FEDLGD

301 The success of FEDLGD relies on the novel design of local-global data distillation, where the
 302 selection of regularization loss and the number of iterations for data distillation plays a key role. In
 303 this section, we study the choice of regularization loss and its weighting (λ) in the total loss function.
 304 Recall that among the total FL training epochs, we perform local-global distillation on the selected
 305 τ iterations, and within each selected iteration, the server and clients will perform data updating

³Cifar10-C is a collection of augmented Cifar10 that applies 19 different corruptions.

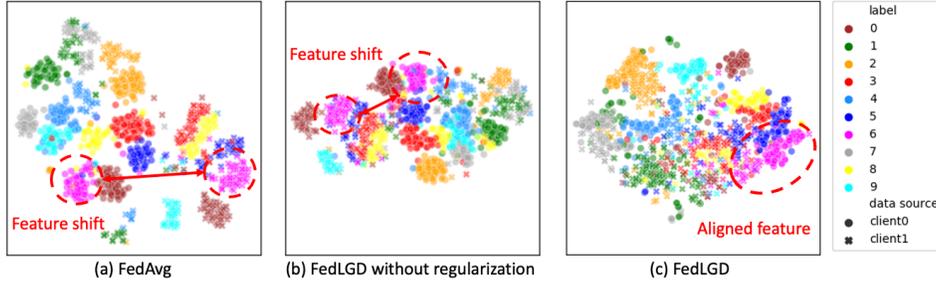


Figure 4: tSNE plots on feature space for FedAvg, FEDLGD without regularization, and FEDLGD. One can observe regularizing training with our global virtual data can rectify feature shift among different clients.

306 for some pre-defined *steps*. The effect of local-global distillation *iterations* and data updating *steps*
 307 will also be discussed. We also perform additional ablation studies such as computation cost and
 308 communication overhead in Appendix C

309 **Effect of regularization loss.** FEDLGD uses supervised contrastive loss \mathcal{L}_{Con} as a regularization
 310 term to encourage local and global virtual data embedding into a similar feature space. To demonstrate
 311 the effectiveness of the regularization term in FEDLGD, we perform ablation studies to replace \mathcal{L}_{Con}
 312 with an alternative distribution similarity measurement, MMD loss, with different λ 's ranging from
 313 0.1 to 20. Figure 3a shows the average test accuracy. Using supervised contrastive loss gives us better
 314 and more stable performance with different coefficient choices.

315 To explain the effect of our proposed regularization loss on feature representations, we embed the
 316 latent features before fully-connected layers to a 2D space using tSNE [28] shown in Figure 4. For
 317 the model trained with FedAvg (Figure 4(a)), features from two clients (\times and \circ) are closer to their
 318 own distribution regardless of the labels (colors). In Figure 4(b), we perform virtual FL training but
 319 without the regularization term (Eq. 4). Figure 4(c) shows FEDLGD, and one can observe that data
 320 from different clients with the same label are grouped together. This indicates that our regularization
 321 with global virtual data is useful for learning homogeneous feature representations.

322 **Analysis of distillation iterations ($|\tau|$).** Figure 3b shows the improved averaged test accuracy if
 323 we increase the number of distillation iterations with FEDLGD. The base accuracy for DIGITS and
 324 CIFAR10C are 85.8 and 55.2, respectively. We fix local and global update *steps* to 100 and 500,
 325 and the selected iterations (τ) are defined as arithmetic sequences with $d = 5$ (i.e., $\tau = \{0, 5, \dots\}$).
 326 One can observe that the model performance improves with the increasing $|\tau|$. This is because we
 327 obtain better virtual data with more local-global distillation iterations, which is a trade-off between
 328 computation cost and model performance. We select $|\tau| = 10$ for efficiency trade-off.

329 **Robustness on virtual data update steps.** In Figure 3c and Figure 3d, we fix $|\tau| = 10$, and vary
 330 (local, global) data updating steps. One can observe that FEDLGD yields stable performance, and the
 331 accuracy slightly improves with an increasing number of local and global steps. Nevertheless, the
 332 results are all the best when comparing with the baselines. It is also worth noting that there is still
 333 trade-off between *steps* and computation cost (See Appendix).

334 5 Conclusion

335 In this paper, we introduce a new approach for FL, called FEDLGD. It utilizes virtual data on both
 336 client and server sides to train FL models. We are the first to reveal that FL on local virtual data
 337 can increase heterogeneity. Furthermore, we propose iterative distribution matching and federated
 338 gradient matching to iteratively update local and global virtual data, and apply global virtual regu-
 339 larization to effectively harmonize domain shift. Our experiments on benchmark and real medical
 340 datasets show that FEDLGD outperforms current state-of-the-art methods in heterogeneous settings.
 341 Furthermore, FEDLGD can be combined with other heterogenous FL methods such as FedProx [26]
 342 and Scaffold [18] to further improve its performance. The potential limitation lies in the additional
 343 communication and computation cost in data distillation, but we show that the trade-off is acceptable
 344 and can be mitigated by decreasing distillation *iterations* and *steps*. Our future direction will be
 345 investigating privacy-preserving data generation. We believe that this work sheds light on how to
 346 effectively mitigate data heterogeneity from a dataset distillation perspective and will inspire future
 347 work to enhance FL performance, privacy, and efficiency.

348 **References**

- 349 [1] Abadi, M., Chu, A., Goodfellow, I., McMahan, H.B., Mironov, I., Talwar, K., Zhang, L.: Deep
350 learning with differential privacy. In: Proceedings of the 2016 ACM SIGSAC conference on
351 computer and communications security. pp. 308–318 (2016)
- 352 [2] Batista, F.J.F., Diaz-Aleman, T., Sigut, J., Alayon, S., Arnay, R., Angel-Pereira, D.: Rim-one dl:
353 A unified retinal image database for assessing glaucoma using deep learning. *Image Analysis &
354 Stereology* **39**(3), 161–167 (2020)
- 355 [3] Carlini, N., Chien, S., Nasr, M., Song, S., Terzis, A., Tramer, F.: Membership inference attacks
356 from first principles. In: 2022 IEEE Symposium on Security and Privacy (SP). pp. 1897–1914.
357 IEEE (2022)
- 358 [4] Carlini, N., Feldman, V., Nasr, M.: No free lunch in" privacy for free: How does dataset
359 condensation help privacy". arXiv preprint arXiv:2209.14987 (2022)
- 360 [5] Cazenavette, G., Wang, T., Torralba, A., Efros, A.A., Zhu, J.Y.: Dataset distillation by matching
361 training trajectories. In: Proceedings of the IEEE/CVF Conference on Computer Vision and
362 Pattern Recognition. pp. 4750–4759 (2022)
- 363 [6] Diaz-Pinto, A., Morales, S., Naranjo, V., Köhler, T., Mossi, J.M., Navea, A.: Cnns for automatic
364 glaucoma assessment using fundus images: an extensive validation. *Biomedical engineering
365 online* **18**(1), 1–19 (2019)
- 366 [7] Dong, T., Zhao, B., Lyu, L.: Privacy for free: How does dataset condensation help privacy?
367 arXiv preprint arXiv:2206.00240 (2022)
- 368 [8] Franceschi, J.Y., De Bézenac, E., Ayed, I., Chen, M., Lamprier, S., Gallinari, P.: A neural
369 tangent kernel perspective of gans. In: International Conference on Machine Learning. pp.
370 6660–6704. PMLR (2022)
- 371 [9] Ganin, Y., Lempitsky, V.: Unsupervised domain adaptation by backpropagation. In: International
372 conference on machine learning. pp. 1180–1189. PMLR (2015)
- 373 [10] Goetz, J., Tewari, A.: Federated learning via synthetic data. arXiv preprint arXiv:2008.04489
374 (2020)
- 375 [11] Gretton, A., Borgwardt, K.M., Rasch, M.J., Schölkopf, B., Smola, A.: A kernel two-sample test.
376 *The Journal of Machine Learning Research* **13**(1), 723–773 (2012)
- 377 [12] Hard, A., Rao, K., Mathews, R., Ramaswamy, S., Beaufays, F., Augenstein, S., Eichner, H.,
378 Kiddon, C., Ramage, D.: Federated learning for mobile keyboard prediction. arXiv preprint
379 arXiv:1811.03604 (2018)
- 380 [13] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Pro-
381 ceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778
382 (2016)
- 383 [14] Hendrycks, D., Dietterich, T.: Benchmarking neural network robustness to common corruptions
384 and perturbations. arXiv preprint arXiv:1903.12261 (2019)
- 385 [15] Hsu, T.M.H., Qi, H., Brown, M.: Measuring the effects of non-identical data distribution for
386 federated visual classification. arXiv preprint arXiv:1909.06335 (2019)
- 387 [16] Hu, S., Goetz, J., Malik, K., Zhan, H., Liu, Z., Liu, Y.: Fedsynth: Gradient compression via
388 synthetic data in federated learning. arXiv preprint arXiv:2204.01273 (2022)
- 389 [17] Hull, J.J.: A database for handwritten text recognition research. *IEEE Transactions on pattern
390 analysis and machine intelligence* **16**(5), 550–554 (1994)
- 391 [18] Karimireddy, S.P., Kale, S., Mohri, M., Reddi, S., Stich, S., Suresh, A.T.: Scaffold: Stochastic
392 controlled averaging for federated learning. In: International Conference on Machine Learning.
393 pp. 5132–5143. PMLR (2020)

- 394 [19] Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4401–4410 (2019)
- 395
396
- 397 [20] Khosla, P., Teterwak, P., Wang, C., Sarna, A., Tian, Y., Isola, P., Maschinot, A., Liu, C., Krishnan, D.: Supervised contrastive learning. *Advances in Neural Information Processing Systems* **33**, 18661–18673 (2020)
- 398
399
- 400 [21] LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**(11), 2278–2324 (1998)
- 401
- 402 [22] Li, G., Togo, R., Ogawa, T., Haseyama, M.: Dataset distillation for medical dataset sharing. *arXiv preprint arXiv:2209.14603* (2022)
- 403
- 404 [23] Li, Q., Diao, Y., Chen, Q., He, B.: Federated learning on non-iid data silos: An experimental study. *arXiv preprint arXiv:2102.02079* (2021)
- 405
- 406 [24] Li, Q., He, B., Song, D.: Model-contrastive federated learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10713–10722 (2021)
- 407
- 408 [25] Li, T., Sahu, A.K., Zaheer, M., Sanjabi, M., Talwalkar, A., Smith, V.: Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems* **2**, 429–450 (2020)
- 409
- 410 [26] Li, X., Huang, K., Yang, W., Wang, S., Zhang, Z.: On the convergence of fedavg on non-iid data. *International Conference on Learning Representations* (2020)
- 411
- 412 [27] Lin, T., Kong, L., Stich, S.U., Jaggi, M.: Ensemble distillation for robust model fusion in federated learning. *Advances in Neural Information Processing Systems* **33**, 2351–2363 (2020)
- 413
- 414 [28] Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. *Journal of machine learning research* **9**(11) (2008)
- 415
- 416 [29] McMahan, B., Moore, E., Ramage, D., Hampson, S., y Arcas, B.A.: Communication-efficient learning of deep networks from decentralized data. In: *Artificial intelligence and statistics*. pp. 1273–1282. PMLR (2017)
- 417
418
- 419 [30] Mohamadi, M.A., Sutherland, D.J.: A fast, well-founded approximation to the empirical neural tangent kernel. *arXiv preprint arXiv:2206.12543* (2022)
- 420
- 421 [31] Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., Ng, A.Y.: Reading digits in natural images with unsupervised feature learning (2011)
- 422
- 423 [32] Orlando, J.I., Fu, H., Breda, J.B., van Keer, K., Bathula, D.R., Diaz-Pinto, A., Fang, R., Heng, P.A., Kim, J., Lee, J., et al.: Refuge challenge: A unified framework for evaluating automated methods for glaucoma assessment from fundus photographs. *Medical image analysis* **59**, 101570 (2020)
- 424
425
426
- 427 [33] Sachdeva, N., McAuley, J.: Data distillation: A survey. *arXiv preprint arXiv:2301.04272* (2023)
- 428 [34] Schuster, A.K., Erb, C., Hoffmann, E.M., Dietlein, T., Pfeiffer, N.: The diagnosis and treatment of glaucoma. *Deutsches Ärzteblatt International* **117**(13), 225 (2020)
- 429
- 430 [35] Shokri, R., Stronati, M., Song, C., Shmatikov, V.: Membership inference attacks against machine learning models. In: *2017 IEEE symposium on security and privacy (SP)*. pp. 3–18. IEEE (2017)
- 431
432
- 433 [36] Sivaswamy, J., Krishnadas, S., Joshi, G.D., Jain, M., Tabish, A.U.S.: Drishti-gs: Retinal image dataset for optic nerve head (onh) segmentation. In: *2014 IEEE 11th international symposium on biomedical imaging (ISBI)*. pp. 53–56. IEEE (2014)
- 434
435
- 436 [37] Tang, Z., Zhang, Y., Shi, S., He, X., Han, B., Chu, X.: Virtual homogeneity learning: Defending against data heterogeneity in federated learning. *arXiv preprint arXiv:2206.02465* (2022)
- 437

- 438 [38] Wang, J., Liu, Q., Liang, H., Joshi, G., Poor, H.V.: Tackling the objective inconsistency problem
439 in heterogeneous federated optimization. *Advances in neural information processing systems*
440 **33**, 7611–7623 (2020)
- 441 [39] Wang, T., Zhu, J.Y., Torralba, A., Efros, A.A.: Dataset distillation. *arXiv preprint*
442 *arXiv:1811.10959* (2018)
- 443 [40] Xiong, Y., Wang, R., Cheng, M., Yu, F., Hsieh, C.J.: Feddm: Iterative distribution matching for
444 communication-efficient federated learning. *arXiv preprint arXiv:2207.09653* (2022)
- 445 [41] Xu, J., Glicksberg, B.S., Su, C., Walker, P., Bian, J., Wang, F.: Federated learning for healthcare
446 informatics. *Journal of Healthcare Informatics Research* **5**, 1–19 (2021)
- 447 [42] Ye, R., Ni, Z., Xu, C., Wang, J., Chen, S., Eldar, Y.C.: Fedfm: Anchor-based feature matching
448 for data heterogeneity in federated learning. *arXiv preprint arXiv:2210.07615* (2022)
- 449 [43] Zhang, L., Shen, L., Ding, L., Tao, D., Duan, L.Y.: Fine-tuning global model via data-free knowl-
450 edge distillation for non-iid federated learning. In: *Proceedings of the IEEE/CVF Conference*
451 *on Computer Vision and Pattern Recognition*. pp. 10174–10183 (2022)
- 452 [44] Zhao, B., Bilen, H.: Dataset condensation with differentiable siamese augmentation. In: *Inter-*
453 *national Conference on Machine Learning*. pp. 12674–12685. PMLR (2021)
- 454 [45] Zhao, B., Bilen, H.: Dataset condensation with distribution matching. In: *Proceedings of the*
455 *IEEE/CVF Winter Conference on Applications of Computer Vision*. pp. 6514–6523 (2023)
- 456 [46] Zhao, B., Mopuri, K.R., Bilen, H.: Dataset condensation with gradient matching. *ICLR* **1**(2), 3
457 (2021)
- 458 [47] Zhou, T., Zhang, J., Tsang, D.: Fedfa: Federated learning with feature anchors to align feature
459 and classifier for heterogeneous data. *arXiv preprint arXiv:2211.09299* (2022)
- 460 [48] Zhu, H., Xu, J., Liu, S., Jin, Y.: Federated learning on non-iid data: A survey. *Neurocomputing*
461 **465**, 371–390 (2021)
- 462 [49] Zhu, L., Liu, Z., Han, S.: Deep leakage from gradients. *Advances in neural information*
463 *processing systems* **32** (2019)