
Parameterized Knowledge Transfer for Personalized Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In recent years, personalized federated learning (pFL) has attracted increasing
2 attention for its potential in dealing with statistical heterogeneity among clients.
3 However, the state-of-the-art pFL methods rely on model parameters aggregation
4 at the server side, which require all models to have the same structure and size, and
5 thus limits the application for more heterogeneous scenarios. To deal with such
6 model constraints, we exploit the potentials of heterogeneous model settings and
7 propose a novel training framework to employ personalized models for different
8 clients. Specifically, we formulate the aggregation procedure in original pFL into a
9 personalized group knowledge transfer training algorithm, namely, KT-pFL, which
10 enables each client to maintain a personalized soft prediction at the server side to
11 guide the others' local training. KT-pFL updates the personalized soft prediction of
12 each client by a linear combination of all local soft predictions using a knowledge
13 coefficient matrix, which can adaptively reinforce the collaboration among clients
14 who own similar data distribution. Furthermore, to quantify the contributions
15 of each client to others' personalized training, the knowledge coefficient matrix
16 is parameterized so that it can be trained simultaneously with the models. The
17 knowledge coefficient matrix and the model parameters are alternatively updated
18 in each round following the gradient descent way. Extensive experiments on
19 various datasets (EMNIST, Fashion_MNIST, CIFAR-10) are conducted under
20 different settings (heterogeneous models and data distributions). It is demonstrated
21 that the proposed framework is the first federated learning paradigm that realizes
22 personalized model training via parameterized group knowledge transfer while
23 achieving significant performance gain comparing with state-of-the-art algorithms.

24 1 Introduction

25 Federated Learning (FL) [1] has emerged as an efficient paradigm to collaboratively train a shared
26 machine learning model among multiple clients without directly accessing their private data. By
27 periodically aggregating parameters from the clients for global model updating, it can converge
28 to high accuracy and strong generalization. FL has shown its capability to protect user privacy
29 while there remains a crucial challenge that significantly degrades the learning performance, i.e.,
30 statistic heterogeneity in users' local datasets. Given the Non-Independent and Identically Distributed
31 (Non-IID) user data, the trained global model often cannot be generalized well over each client [2–5].

32 To deal with the above issues, employing personalized models appears to be an effective solution in
33 FL, i.e., personalized federated learning (pFL). Recent works regarding pFL include regularization-
34 based methods [6–8] (pFedMe [6], L2SGD [7], FedAMP [8]), meta-learning-based Per-FedAvg [9]
35 and cluster-based IFCA [10, 11]. However, in order to aggregate the parameters from all clients, it is
36 inevitable for them to have identical model structure and size. Such constraints would prevent status

37 quo pFL methods from further application in practical scenarios, where clients are often willing to
38 own unique models, i.e., with customized neural architectures to adapt to heterogeneous capacities in
39 computation, communication and storage space, etc.

40 Motivated by the paradigm of Knowledge Distillation (KD) [12–16] that knowledge can be transferred
41 from a neural network to another via exchanging soft predictions instead of using the whole model
42 parameters, we seek to develop a novel training framework that can accommodate heterogeneous
43 model structures for each client and achieve personalized knowledge transfer in each FL training
44 round. To this end, we formulate the aggregation phase in FL to a personalized group knowledge
45 transfer training algorithm dubbed KT-pFL, whose main idea is to allow each client to maintain
46 a personalized soft prediction at the server that can be updated by a linear combination of all
47 clients’ local soft predictions using a knowledge coefficient matrix. The principle of doing so is to
48 reinforce the collaboration between clients with similar data distributions. Furthermore, to quantify
49 the contribution of each client to other’s personalized soft prediction, we parameterize the knowledge
50 coefficient matrix so that it can be trained simultaneously with the models following an alternating
51 way in each iteration round.

52 We show that KT-pFL not only breaks down the barriers of homogeneous model restriction, which
53 requires to transfer the entire parameters set in each round, whose data volume is much larger
54 than that of the soft prediction, but also improves the training efficiency by using a parameterized
55 update mechanism. Experimental results on different datasets and models show that our method
56 can significantly improve the training efficiency and reduce the communication overhead. Our
57 contributions are:

- 58 • To the best of our knowledge, this paper is the first to study the personalized knowledge transfer
59 in FL. We propose a novel training framework, namely, KT-pFL, that maintains a personalized soft
60 prediction for each client in the server to transfer knowledge among all clients.
- 61 • To encourage clients with similar data distribution to collaborate with each other during the training
62 process, we propose the ‘knowledge coefficient matrix’ to identify the contribution from one client to
63 others’ local training. To show the efficiency of the parameterized method, we compared KT-pFL
64 with two non-parameterized learning methods, i.e., TopK-pFL and Sim-pFL, which calculate the
65 knowledge coefficient matrix on the cosine similarity between different model parameters.
- 66 • We provide theoretical performance guarantee for KT-pFL and conduct extensive experiments over
67 various deep learning models and datasets. The efficiency superiority of KT-pFL is demonstrated by
68 comparing our proposed training framework with traditional pFL methods.

69 2 Related Work

70 **Personalized Federated Learning.** Recently, various approaches have been proposed to realize
71 personalized FL, which can be categorized into three types according to the number of global models
72 applied in the server, i.e., single global model, multiple global models, no global model.

73 Single-global-model-based methods extended from the conventional FL algorithms like FedAvg [1]
74 combine the optimization of the local models and global model, which consist of four different kinds
75 of approaches: local fine-tuning [17–19], regularization [6, 7, 20], mixture [11, 21] and meta learning
76 [9, 22]. The straightforward approach for personalization is the local fine-tuning, where all clients
77 collaboratively train a global model at beginning and then customize the well-trained global model
78 using local data with extra gradient descent steps. Another pFL way is to constrain the deviation
79 between local model and global model by adding a regularized term in the objective function (e.g.,
80 pFedMe [6], L2SGD [7, 20]). Instead of adding regularized loss function, Mansour et al. [11]
81 introduce a general approach for FL by mixing the global and local models, and three different
82 personalization methods are proposed with generalization guarantees, namely, client clustering, data
83 interpolation, and model interpolation. The third approach is also used in [21] for the adaptive
84 personalized federated learning (ApFL) algorithm, which aims to find the optimal combination of the
85 global model and the local models. Besides, meta-learning can be used in FL for better personalization
86 by generating a highly-adaptable global model that can help to boost the model training with a few
87 data samples. Inspired by Model-Agnostic Meta-Learning (MAML) [23], Fallah et al. [9] propose to
88 learn personalized models for each client with convergence guarantees (Per-FedAvg).

89 All of the above mentioned pFL methods apply a single global model, and only allow limited
90 customization to the model at the client side. Therefore, some researchers [8, 10, 11] propose to

91 train multiple global models at the server, where similar clients are clustered into several groups and
 92 different models are trained for each group. Mansour et al. [11] provide generalization guarantees
 93 for multiple global models, while Ghosh [10] take a further step to establish the convergence theory
 94 regarding the training loss function. FedAMP [8] can be regraded as a special case of the clustered-
 95 based method that each client owns a personalized global model in the server side.

96 As a contrast, some research proposals haven't applied any global model to deal with the heterogeneity
 97 problem. Taking an example in [24], the authors use multi-task learning framework to design a new
 98 training method called MOCHA, to address the statistical and systematic heterogeneity challenges.
 99 However, multi-task learning based methods need to solve the data-local quadratic subproblem for
 100 every client, which is computationally intensive. It requires all clients to participate in every training
 101 round, which is inapplicable for large-scale clients.

102 **Knowledge Distillation.** Knowledge Distillation (KD) [25] is introduced to enhance the performance
 103 of a small student network by transferring the knowledge from a large teacher network. The principle
 104 is to update the student model to approximate the soft prediction output of the teacher model [26–30].
 105 As KD is independent with model structure, some literature [31–33] are proposed to take advantage
 106 of such flexibility to implement collaborative learning via KD, i.e., distilling the knowledge of an
 107 ensemble of teacher models to a student model [32]. Most of these schemes construct an ensembling
 108 teacher by simply averaging the teachers' soft predictions, or by heuristically combining the output of
 109 the teacher models, which are hard to achieve optimal combination of teachers. In our framework, KD
 110 is used in a more efficient way that the weights of the clients' soft predictions are updated together
 111 with the model parameters during every FL training iterations.

112 3 Problem Formulation

113 We aim to collaboratively train personalized models for a set of clients applying different model
 114 structures in FL. Consider supervised learning whose goal is to learn a function that maps every input
 115 data to the correct class out of \mathcal{C} possible options. We assume that there are N clients, and each
 116 client n can only access to his private dataset $\mathbb{D}_n := \{x_i^n, y_i\}$, where x_i is the i -th input data sample,
 117 y_i is the corresponding label of x_i , $y_i \in \{1, 2, \dots, C\}$. The number of data samples in dataset \mathbb{D}_n
 118 is denoted by D_n . $\mathbb{D} = \{\mathbb{D}_1, \mathbb{D}_2, \dots, \mathbb{D}_N\}$, $N = \sum_{n=1}^N D_n$. In conventional FL, the goal of the
 119 learning system is to learn a global model \mathbf{w} that minimizes the total empirical loss over the entire
 120 dataset \mathbb{D} :

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) := \sum_{n=1}^N \frac{D_n}{D} \mathcal{L}_n(\mathbf{w}), \text{ where } \mathcal{L}_n(\mathbf{w}) = \frac{1}{D_n} \sum_{i=1}^{D_n} \mathcal{L}_{CE}(\mathbf{w}; x_i, y_i), \quad (1)$$

121 where $\mathcal{L}_n(\mathbf{w})$ is the n -th client's local loss function that measures the local empirical risk over the
 122 private dataset \mathbb{D}_n and \mathcal{L}_{CE} is the cross-entropy loss function that measures the difference between
 123 the predicted values and the ground truth labels.

124 However, this formulation requires all local clients to have a unified model structure, which cannot
 125 be extended to more general cases where each client applies a unique model. Therefore, we need to
 126 reformulate the above optimization problem to breaks the barrier of homogeneous model structure.
 127 Besides, given the Non-IID clients' datasets, it is inappropriate to minimize the total empirical loss.
 128 To that end, we propose the following training framework:

129 **Definition 1.** Let $s(\mathbf{w}^n, \hat{x})$ denote the *collaborative knowledge* from client n , and \hat{x} denote a data
 130 sample from a public dataset \mathbb{D}_r that all clients can access to. Define the personalized loss function
 131 of client n as

$$\mathcal{L}_{per,n}(\mathbf{w}^n) := \mathcal{L}_n(\mathbf{w}^n) + \lambda \sum_{\hat{x} \in \mathbb{D}_r} \mathcal{L}_{KL} \left(\sum_{m=1}^N c_{mn} \cdot s(\mathbf{w}^m, \hat{x}), s(\mathbf{w}^n, \hat{x}) \right), \quad (2)$$

132 where $\lambda > 0$ is a hyper-parameters, \mathcal{L}_{KL} stands for Kullback–Leibler (KL) Divergence function
 133 and is added to the loss function to transfer personalized knowledge from a teacher to another. c_{mn}
 134 is the knowledge coefficient which is used to estimate the contribution from client m to n . The
 135 second term in (2) allows each client to build his own personalized aggregated knowledge in the
 136 server and enhance the collaboration effect between clients with large c . The concept of *collaborative*
 137 *knowledge* can either refer to soft predictions or model parameters depending on the definition in

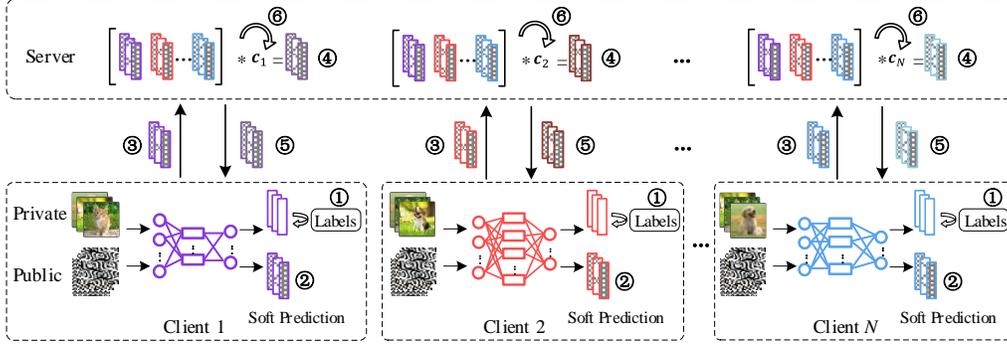


Figure 1: Illustration of the KT-pFL framework. The workflow includes 6 steps: ① local training on private data; ②, ③ each client outputs the local soft prediction on public data and sends it to the server; ④ the server calculates each client’s personalized soft prediction via a linear combination of local soft predictions and knowledge coefficient matrix; ⑤ each client downloads the personalized soft prediction to perform distillation phase; ⑥ the server updates the knowledge coefficient matrix.

138 different situations. For example, $s(\mathbf{w}^n, \hat{x})$ can be deemed to be a soft prediction of the client n ,
 139 which are calculated with the softmax of logits z^n , i.e., $s(\mathbf{w}^n, \hat{x}) = \frac{\exp(z_c^n/T)}{\sum_{c=1}^C \exp(z_c^n/T)}$, logits z^n is the
 140 output of the last fully connected layer on client n ’s model, T is the temperature hyperparameter of
 141 the softmax function.

142 **Definition 2.** We define $\mathbf{c} \in \mathbb{R}^{N \times N}$ as the *knowledge coefficient matrix*:

$$\mathbf{c} = \begin{Bmatrix} c_{11} & c_{12} & \cdots & c_{1N} \\ c_{21} & c_{22} & \cdots & c_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ c_{N1} & c_{N2} & \cdots & c_{NN} \end{Bmatrix}. \quad (3)$$

143 Our objective is to minimize

$$\min_{\mathbf{w}, \mathbf{c}} \mathcal{L}(\mathbf{w}, \mathbf{c}) := \sum_{n=1}^N \frac{D_n}{D} \mathcal{L}_{per,n}(\mathbf{w}^n) + \rho \|\mathbf{c} - \frac{\mathbf{1}}{N}\|^2, \quad (4)$$

144 where $\mathbf{w} = [\mathbf{w}^1, \dots, \mathbf{w}^N] \in \mathbb{R}^{\sum_{n=1}^N d_n}$ is the concatenated vector of all weights, d_n represents the
 145 dimensions of model parameter \mathbf{w}^n . $\mathbf{1} \in \mathbb{R}^{n^2}$ is the identity matrix whose elements are all equal to
 146 1. The second term in (4) is a regularization term that ensures generalization ability of the whole
 147 learning system. ρ is a regularization parameter that larger than 0.

148 4 KT-pFL Algorithm

149 In this section, we introduce the proposed KT-pFL algorithm, where the local model parameters and
 150 knowledge coefficient matrix are updated alternatively. To enable personalized knowledge transfer in
 151 FL, we train personalized models locally according to the related *collaborative knowledge*. Insert
 152 (2) into (4), and we can design an alternating optimization approach to solve (4), that in each round
 153 we fix either \mathbf{w} or \mathbf{c} by turns, and optimize the unfixed one following an alternating way until a
 154 convergence point is reached.

155 **Update \mathbf{w} :** In each communication round, we first fix \mathbf{c} and optimize (train) \mathbf{w} for several epochs
 156 locally. In this case, updating \mathbf{w} depends on both the private data (i.e., \mathcal{L}_{CE} on $\mathbb{D}_n, n \in [1, \dots, N]$),
 157 that can only be accessed by the corresponding client, and the public data (i.e., \mathcal{L}_{KL} on \mathbb{D}_r), which is
 158 accessible for all clients. We propose a two-stage updating framework for \mathbf{w} :

159 • *Local Training:* Train \mathbf{w} on each client’s private data by applying a gradient descent step:

$$\mathbf{w}^n \leftarrow \mathbf{w}^n - \eta_1 \nabla_{\mathbf{w}^n} \mathcal{L}_n(\mathbf{w}^n; \xi_n), \quad (5)$$

160 where ξ_n denotes the mini-batch of data \mathbb{D}_n used in local training, η_1 is the learning rate.

Algorithm 1 KT-pFL Algorithm

Input: Private dataset \mathbb{D}_n , public dataset \mathbb{D}_r , personalized model \mathbf{w}^n , $n = 1, \dots, N$, regularization parameter ρ and learning rate η_1, η_2, η_3 . Total communication rounds T .

Output: Trained personalized models $\mathbf{w} = [\mathbf{w}^1, \dots, \mathbf{w}^N]$.

```

1: Initialize the model parameters  $\mathbf{w}_0$  and knowledge coefficient matrix  $\mathbf{c}_0$ .
2: procedure SERVER-SIDE OPTIMIZATION
3:   Distribute  $\mathbf{w}_0$  and  $\mathbf{c}_0$  to each client
4:   for each communication round  $t \in \{1, 2, \dots, T\}$  do
5:     for each client  $n$  in parallel do
6:        $\mathbf{w}_{t+1}^n \leftarrow \text{ClientLocalUpdate}(n, \mathbf{w}_t^n, \mathbf{c}_{t,n})$ 
7:       Update knowledge coefficient matrix  $\mathbf{c}$  via (7)
8:       Distribute  $\mathbf{c}_{t+1}$  to all clients
9:   procedure CLIENTLOCALUPDATE( $n, \mathbf{w}_t^n, \mathbf{c}_{t,n}$ )
10:    Client  $n$  receives  $\mathbf{w}_t^n$  and  $\mathbf{c}_n$  from the server
11:    for each local epoch  $i$  from 1 to  $E$  do
12:      for mini-batch  $\xi_t \subseteq \mathbb{D}_n$  do
13:        Local Training: update model parameters on private data via (5)
14:      for each distillation step  $j$  from 1 to  $R$  do
15:        for mini-batch  $\xi_{r,t} \subseteq \mathbb{D}_r$  do
16:          Distillation: update model parameters on public data via (6)
17:    return local parameters  $\mathbf{w}_{t+1}^n$ 

```

161 • *Distillation:* Transfer knowledge from personalized soft prediction to each local client based on
 162 public dataset:

$$\mathbf{w}^n \leftarrow \mathbf{w}^n - \eta_2 \nabla_{\mathbf{w}^n} \mathcal{L}_{KL} \left(\sum_{m=1}^N \mathbf{c}_m^{*,T} \cdot s(\mathbf{w}^m, \xi_r), s(\mathbf{w}^n, \xi_r) \right), \quad (6)$$

163 where ξ_r denotes the mini-batch of public data \mathbb{D}_r , and η_2 is the learning rate. $\mathbf{c}_m^* =$
 164 $[c_{m1}, c_{m2}, \dots, c_{mN}]$ is the *knowledge coefficient* vector for client m , which can be found in
 165 m -th row of \mathbf{c} . Note that all *collaborative knowledge* and *knowledge coefficient matrix* are required
 166 to obtain the personalized soft prediction in this stage, which can be collected in the server.

167 **Update \mathbf{c} :** After updating \mathbf{w} locally for several epochs, we turn to fix \mathbf{w} and update \mathbf{c} in the server.

$$\mathbf{c} \leftarrow \mathbf{c} - \eta_3 \lambda \sum_{n=1}^N \frac{D_n}{D} \nabla_{\mathbf{c}} \mathcal{L}_{KL} \left(\sum_{m=1}^N \mathbf{c}_m \cdot s(\mathbf{w}^{m,*}, \xi_r), s(\mathbf{w}^{n,*}, \xi_r) \right) - 2\eta_3 \rho \left(\mathbf{c} - \frac{\mathbf{1}}{N} \right), \quad (7)$$

168 where η_3 is the learning rate for updating \mathbf{c} .

169 Algorithm 1 demonstrates the proposed KT-pFL algorithm and the idea behind it is shown in Figure 1.
 170 In every communication round of training, the clients use local SGD to train several epochs based
 171 on the private data and then send the *collaborative knowledge* (e.g., soft predictions on public data)
 172 to the server. When the server receives the *collaborative knowledge* from each client, it aggregates
 173 them to form the personalized soft predictions according to the knowledge coefficient matrix. The
 174 server then sends back the personalized soft prediction to each client to perform local distillation.
 175 The clients then iterate for multiple steps over public dataset¹. After that, the knowledge coefficient
 176 matrix is updated in the server while fixing the model parameters \mathbf{w} .

177 **Performance Guarantee.** Theorem 1 provides the performance analysis of the personalized model
 178 when each client owns Non-IID data. Detailed description and derivations are deferred to Appendix.

179 **Theorem 1.** Denote the n -th local distribution and its empirical distribution by \mathcal{D}_n and $\hat{\mathcal{D}}_n$ respec-
 180 tively, and the hypothesis $h \in \mathcal{H}$ trained on $\hat{\mathcal{D}}_n$ by $h_{\hat{\mathcal{D}}_n}$. There always exist $c_{m,n}^*$, $m = 1, \dots, N$,
 181 such that the expected loss of the personalized ensemble model for the data distribution \mathcal{D}_n of client
 182 n is not larger than that of the single model only trained with local data: $\mathcal{L}_{\mathcal{D}_n}(\sum_{m=1}^N c_{m,n}^* h_{\hat{\mathcal{D}}_m}) \leq$
 183 $\mathcal{L}_{\mathcal{D}_n}(h_{\hat{\mathcal{D}}_n})$. Besides, there exist some problems where the personalized ensemble model is strictly
 184 better, i.e., $\mathcal{L}_{\mathcal{D}_n}(\sum_{m=1}^N c_{m,n}^* h_{\hat{\mathcal{D}}_m}) < \mathcal{L}_{\mathcal{D}_n}(h_{\hat{\mathcal{D}}_n})$.

¹Note that our method can work over both labeled and unlabeled public datasets.

185 Theorem 1 indicates that the performance of the personalized ensemble model under some suitable
 186 coefficient matrices is better than that of the model only trained on its local private data, which
 187 theoretically demonstrates the necessity of the parameterized personalization. However, it is challeng-
 188 ing to find such matrices due to the complexity and diversity of machine learning models and data
 189 distributions. In this paper, our designed algorithm KT-pFL can find the desired coefficient matrix
 190 in a gradient descent manner, hence achieving the performance boost. Besides, we have the similar
 191 claim for the relationship between the personalized ensemble model and average ensemble model.

192 **Remark 1.** *There always exist $c_{m,n}^*$, $m = 1, \dots, N$, such that the expected loss of the personal-*
 193 *ized ensemble model for the data distribution \mathcal{D}_n of client n is not larger than that of the aver-*
 194 *age ensemble model: $\mathcal{L}_{\mathcal{D}_n}(\sum_{m=1}^N c_{m,n}^* h_{\hat{\mathcal{D}}_m}) \leq \mathcal{L}_{\mathcal{D}_n}(\frac{1}{N} \sum_{m=1}^N h_{\hat{\mathcal{D}}_m})$. Besides, there exist some*
 195 *problems where the personalized ensemble model is strictly better, i.e., $\mathcal{L}_{\mathcal{D}_n}(\sum_{m=1}^N c_{m,n}^* h_{\hat{\mathcal{D}}_m}) <$*
 196 *$\mathcal{L}_{\mathcal{D}_n}(\frac{1}{N} \sum_{m=1}^N h_{\hat{\mathcal{D}}_m})$.*

197 5 Evaluations

198 5.1 Experimental Setup

199 **Task and Datasets** We evaluate our proposed training framework on three different image classifi-
 200 cation tasks: EMNIST [34], Fashion_MNIST [35] and CIFAR-10 [36]. For each dataset, we apply
 201 two different Non-IID data settings: 1) each client only contains two classes of samples; 2) each
 202 client contains all classes of samples, while the number of samples for each class is different from
 203 that of a different client. All datasets are split randomly with 75% and 25% for training and testing,
 204 respectively. The testing data on each client has the same distribution with its training data. For all
 205 methods, we record the average test accuracy of all local models for evaluation.

206 **Model Structure:** Four different lightweight model structures including LeNet [37], AlexNet [38],
 207 ResNet-18 [39], and ShuffleNetV2 [40] are adopted in our experiments. Our pFL system has 20
 208 clients, who are assigned with four different model structures, i.e., five clients per model.

209 **Baselines:** Although KT-pFL is designed for effective personalized federated learning, it can be
 210 applied to both heterogeneous and homogeneous model cases. We first conduct experiments on
 211 heterogeneous systems where the neural architecture of the local models are different among clients.
 212 We compare the performance of KT-pFL to the non-personalized distillation-based method named
 213 FedDF [16]² and other simple versions of KT-pFL including Sim-pFL and TopK-pFL. Sim-pFL
 214 calculates the knowledge coefficient by using cosine similarity between two local soft predictions.
 215 Instead of combining the knowledge from all clients, TopK-pFL obtains the personalized soft
 216 predictions only from K clients³ who have higher value of cosine similarity.

217 To demonstrate the generalization and effectiveness of our proposed training framework, we further
 218 compare KT-pFL to FedAvg [1] and state-of-the-art pFL methods including Per-Fedavg [9], pFedMe
 219 [6] and FedAMP [8] under the homogeneous model setting. Note that Per-FedAvg is a MAML-based
 220 method which aims to optimize the one-step gradient update for its personalized model. pFedMe
 221 and FedAMP are two regularized-based methods. The former one obtains personalized models by
 222 controlling the distance between local model and global model, while the later one facilitates the
 223 collaboration of clients with similar data distribution.

224 **Implementation** The experiments are implemented in PyTorch. We simulate a set of clients and a cen-
 225 tralized server on one deep learning workstation (i.e., Intel(R) Core(TM) i9-9900KF CPU@3.6GHz
 226 with one NVIDIA GeForce RTX 2080Ti GPU).

227 5.2 Results

228 Subject to space constraints, we only report the most important experimental results in this section.
 229 Please refer to Appendix for the details of different Non-IID data settings on EMNIST, Fash-

²FedDF exchanges model parameters with the server and performs knowledge transfer after averaging the local model parameters at the server side. In our heterogeneous setting, the local clients exchange soft predictions with the server, so we omit the parameters aggregation phase of FedDF in our experiments.

³In our experiments, we set K as 5.

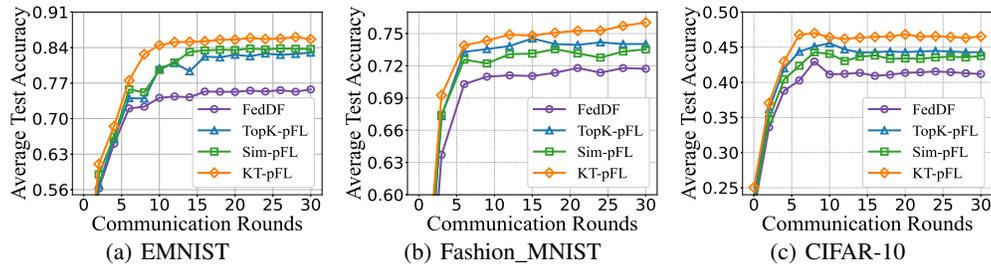


Figure 2: Performance comparison of FedDF, TopK-pFD, Sim-pFD, and KT-pFL in average test accuracy on three datasets (Non-IID case 1: each client contains all labels).

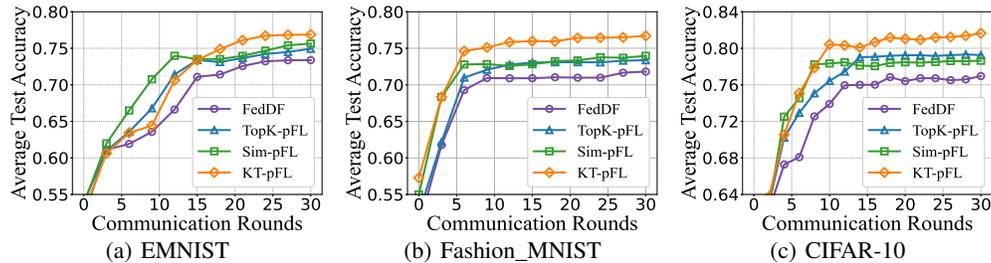


Figure 3: Performance comparison of FedDF, TopK-pFD, Sim-pFD, and KT-pFL in average test accuracy on three datasets (Non-IID case 2: each client contains only two labels).

230 ion_MNIST and CIFAR10, the implementation details and the hyperparameter settings of all the
 231 methods, and also extra results about the convergence and robustness analysis.

232 5.2.1 Performance Comparison

233 **Heterogeneous FL.** For all the methods and all the data settings, the batch size on private data and
 234 public data are 128 and 256, respectively, the number of local epochs is 20 and the distillation steps
 235 is 1 in each communication round of pFL training. Unless mentioned, otherwise the size of public
 236 data used in each communication round is 3000, the learning rate is set to 0.01 for EMNIST and
 237 Fashion_MNIST, and 0.02 for CIFAR-10.

Table 1: The comparison of final test accuracy on different datasets with heterogeneous models (i.e., Lenet, AlexNet, ResNet-18 and ShuffleNetV2).

Method	EMNIST (%)		FashionMNIST (%)		CIFAR-10 (%)	
	Non-IID_1	Non-IID_2	Non-IID_1	Non-IID_2	Non-IID_1	Non-IID_2
FedDF [16]	72.75	75.41	71.89	72.76	42.03	39.50
Sim-pFD	82.72	74.42	72.82	72.71	42.50	39.90
TopK-pFD	83.00	73.91	72.24	73.39	42.58	39.34
KT-pFL	84.65	76.30	75.48	75.60	43.82	41.04

238 Figure 2 and 3 show the curves of the average test accuracy during the training process on four
 239 different models with three datasets, which include the results of FedDF, Sim-pFL, TopK-pFL and
 240 pFL. We also summarize the final average test accuracy in Table 1. In two cases of Non-IID settings,
 241 KT-pFL obtains comparable or even better accuracy performance than others. The performance of
 242 FedDF is worse than the other three methods. The reason is that taking the global aggregation of all
 243 local soft predictions trained on the Non-IID data from different clients produces only one global
 244 soft prediction, which cannot be well adapted to each client. The other two personalized methods,
 245 i.e., Sim-pFL and TopK-pFL, achieve comparably performance on most of cases with FedDF. Our
 246 proposed method KT-pFL has the best performance, because each client can adaptively aggregate all
 247 local soft predictions to form a personalized one instead of being restricted to a global soft prediction.

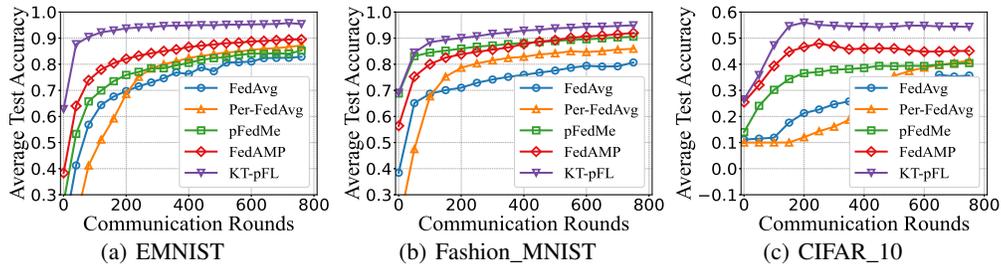


Figure 4: Performance comparison of FedAvg, Per-Fedavg, pFedMe, FedAMP and KT-pFL in average test accuracy on three datasets. The Non-IID data setting: each client contains all labels. 20 clients with homogeneous models: CNN [1]. Learning rate: 0.005.

Table 2: The comparison of final test accuracy on different datasets with homogeneous models (i.e., the same CNN architecture as [1]).

Method	EMNIST (%)	Fashion_MNIST (%)	CIFAR-10 (%)
Fedavg [1]	85.86	80.69	35.69
Per-Fedavg [9]	86.16	85.91	41.37
pFedMe [6]	84.26	90.97	40.44
FedAMP [8]	88.51	90.79	40.11
KT-pFL	94.40	93.93	52.35

248 **Homogeneous FL.** Apart from heterogeneous system, we further extend KT-pFL to homogeneous
 249 model setting. To conduct fair comparison with baselines, we exchange model parameters with
 250 the server to perform knowledge transfer. Specifically, each client maintains a personalized global
 251 model at the server side and each personalized global model is aggregated by a linear combination of
 252 local model parameters and knowledge coefficient matrix. In this case, we compare the performance
 253 of homogeneous-version of KT-pFL with FedAvg, Per-Fedavg, pFedMe and FedAMP. Figure 4
 254 and Table 2 show the curve of the average test accuracy during training and the final average test
 255 accuracy on three different datasets, respectively. For all datasets, KT-pFL dominates the other four
 256 methods on average test accuracy with less variance. Besides, the concept of FedAMP is similar
 257 to our proposed method KT-pFL. The difference is that the weights for each local model during
 258 personalized aggregation phase are updated in a gradient descent manner in KT-pFL.

259 5.2.2 Effect of hyperparameters

260 To understand how different hyperparameters such as E , R , ρ , and $|\xi_r|$ can affect the training
 261 performance of KT-pFL in different settings, we conduct various experiments on three dataset with
 262 $\eta_1 = \eta_2 = \eta_3 = 0.01$.

263 **Effect of Local Epochs E .** To reduce communication overhead, the server tends to allow clients to
 264 have more local computation steps, which can lead to less global updates and thus faster convergence.
 265 Therefore, we monitor the behavior of KT-pFL using a number of values of E , whose results are given
 266 in Table 3. The results show that larger values of E can benefit the convergence of the personalized
 267 models. There is, nevertheless, a trade-off between the computations and communications, i.e.,
 268 while larger E requires more computations at local clients, smaller E needs more global commu-
 269 nication rounds to converge. To do such trade-off, we fix $E = 20$ and evaluate the effect of other
 270 hyperparameters accordingly.

271 **Effect of Distillation Steps R .** As the performance of the knowledge transfer is directly related to
 272 the number of distillation steps R , we compare the performance of KT-pFL under different setting of
 273 distillation steps (e.g., $R = 1, 2, 3, 5$). The number of local epochs is set to 20. The results show that
 274 larger value of R cannot always lead to better performance, which means a moderate number of the
 275 distillation steps is necessary to approach to the optimal performance.

276 **Effect of ρ .** As mentioned in (4), ρ is the regularization term to do the trade-off between personal-
 277 ization ability and generalization ability. For example, larger value of ρ means that the knowledge

Table 3: The comparison of final test accuracy on different datasets with homogeneous models.

Dataset	# Local epochs (E)				# Distillation steps (R)			
	5	10	15	20	1	2	3	5
EMNIST (%)	90.15	92.76	91.03	90.15	91.76	91.40	91.18	91.54
Fashion_MNIST (%)	88.73	89.14	88.58	89.42	89.14	89.90	89.07	88.08
CIFAR-10 (%)	58.30	59.38	59.34	59.24	59.24	57.99	59.22	58.91

Table 4: The comparison of final test accuracy on different datasets with homogeneous models

Dataset	# Regularization parameter ρ				# Batch size of public data ($ \xi_r $)			
	0.1	0.3	0.5	0.7	500	1000	3000	5000
EMNIST (%)	90.96	90.66	90.88	91.76	91.25	91.47	91.66	91.32
Fashion_MNIST (%)	89.49	89.30	88.56	89.14	88.08	89.40	89.50	89.14
CIFAR-10 (%)	59.08	58.52	59.56	58.40	58.93	58.79	58.91	58.40

278 coefficient of each local soft prediction should approach to $\frac{1}{N}$, and thus more generalization ability
 279 should be guaranteed. Table 4 shows the results of KT-pFL with different value of ρ . In most settings,
 280 a significantly large value of ρ will hurt the performance of KT-pFL.

281 **Effect of $|\xi_r|$.** Table 4 demonstrates the effect of batch size of public data $|\xi_r|$ used in local distillation
 282 phase. When the size of the mini-batch is increased, KT-pFL has higher average test accuracy.
 283 However, very large $|\xi_r|$ will not only slow down the convergence of KT-pFL but also incur higher
 284 computation time at the clients. During the experiments, the value of $|\xi_r|$ is configured to 3000.

285 5.2.3 Efficiency Evaluation



Figure 5: Communication Efficiency (X-axis units: MBytes) on three datasets. KL-pFL: soft prediction-based personalized federated learning; Conv-pFL: conventional parameter-based personalized federated learning. (20 models, 30 communication rounds for all datasets. In this experiment, the public dataset is stored on the server.)

286 To evaluate the communication overhead, we record three aspects information: Data (batch size used
 287 in distillation phase), Param (model parameters) and Soft Prediction without using data compression
 288 techniques. The results are shown in Figure 5. Compared with conventional parameter-based pFL
 289 methods, the communication overhead on soft prediction-based KT-pFL is far less than Conv-pFL.

290 Broader Impact

291 FL has been emerged as new paradigm to collaboratively train models among multiple clients in a
 292 privacy-preserving manner. Due to the diversity of users (e.g., statistical and systematic heterogeneity,
 293 etc.), applying personalization in FL is essential for future trend. Our method KT-pFL not only breaks
 294 the barriers of homogeneous model constraint, which can significantly reduce the communication
 295 overhead during training, but also improves the training efficiency via a parameterized update
 296 mechanism without additional computation overhead at the client side. This research has the potential
 297 to enable various devices to cooperatively train ML tasks based on customized neural network
 298 architectures.

299 **References**

- 300 [1] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas.
301 Communication-efficient learning of deep networks from decentralized data. In *Proceedings of*
302 *the 20th International Conference on Artificial Intelligence and Statistics, AISTATS*, 2017.
- 303 [2] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia
304 Smith. Federated optimization in heterogeneous networks. In *Proceedings of Machine Learning*
305 *and Systems, MLSys*, 2020.
- 306 [3] Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In
307 *Proceedings of International Conference on Machine Learning, ICML*, pages 4615–4625, 2019.
- 308 [4] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and
309 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In
310 *Proceedings of International Conference on Machine Learning, ICML*, pages 5132–5143, 2020.
- 311 [5] Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. Fair resource allocation in
312 federated learning. In *Proceedings of 8th International Conference on Learning Representations,*
313 *ICLR*, 2020.
- 314 [6] Canh T. Dinh, Nguyen H. Tran, and Tuan Dung Nguyen. Personalized federated learning with
315 moreau envelopes. In *Proceedings of Advances in Neural Information Processing Systems 33:*
316 *Annual Conference on Neural Information Processing Systems, NeurIPS*, 2020.
- 317 [7] Filip Hanzely and Peter Richtárik. Federated learning of a mixture of global and local models.
318 *arXiv preprint arXiv:2002.05516*, 2020.
- 319 [8] Yutao Huang, Lingyang Chu, Zirui Zhou, Lanjun Wang, Jiangchuan Liu, Jian Pei, and Yong
320 Zhang. Personalized cross-silo federated learning on non-iid data. In *Proceedings of the AAAI*
321 *Conference on Artificial Intelligence, AAAI*, 2021.
- 322 [9] Alireza Fallah, Aryan Mokhtari, and Asuman E. Ozdaglar. Personalized federated learning with
323 theoretical guarantees: A model-agnostic meta-learning approach. In *Proceedings of Advances*
324 *in Neural Information Processing Systems 33: Annual Conference on Neural Information*
325 *Processing Systems, NeurIPS*, 2020.
- 326 [10] Avishek Ghosh, Jichan Chung, Dong Yin, and Kannan Ramchandran. An efficient framework
327 for clustered federated learning. In *Proceedings of Advances in Neural Information Processing*
328 *Systems 33: Annual Conference on Neural Information Processing Systems, NeurIPS*, 2020.
- 329 [11] Yishay Mansour, Mehryar Mohri, Jae Ro, and Ananda Theertha Suresh. Three approaches for
330 personalization with applications to federated learning. *arXiv preprint arXiv:2002.10619*, 2020.
- 331 [12] Eunjeong Jeong, Seungeun Oh, Hyesung Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun
332 Kim. Communication-efficient on-device machine learning: Federated distillation and augmen-
333 tation under non-iid private data. *arXiv preprint arXiv:1811.11479*, 2018.
- 334 [13] Daliang Li and Junpu Wang. FedMD: Heterogenous federated learning via model distillation.
335 *arXiv*, oct 2019.
- 336 [14] Hongyan Chang, Virat Shejwalkar, Reza Shokri, and Amir Houmansadr. Cronus: Robust
337 and heterogeneous collaborative learning with black-box knowledge transfer. *arXiv preprint*
338 *arXiv:1912.11279*, 2019.
- 339 [15] Sohei Itahara, Takayuki Nishio, Yusuke Koda, Masahiro Morikura, and Koji Yamamoto.
340 Distillation-based semi-supervised federated learning for communication-efficient collaborative
341 training with non-iid private data. *arXiv preprint arXiv:2008.06180*, 2020.
- 342 [16] Tao Lin, Lingjing Kong, Sebastian U. Stich, and Martin Jaggi. Ensemble distillation for
343 robust model fusion in federated learning. In *Proceedings of Advances in Neural Information*
344 *Processing Systems 33: Annual Conference on Neural Information Processing Systems, NeurIPS*,
345 2020.

- 346 [17] Kangkang Wang, Rajiv Mathews, Chloé Kiddon, Hubert Eichner, Françoise Beaufays,
347 and Daniel Ramage. Federated evaluation of on-device personalization. *arXiv preprint*
348 *arXiv:1910.10252*, 2019.
- 349 [18] Johannes Schneider and Michail Vlachos. Personalization of deep learning. *arXiv preprint*
350 *arXiv:1909.02803*, 2019.
- 351 [19] Manoj Ghuhana Arivazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunav Choudhary.
352 Federated learning with personalization layers. *arXiv preprint arXiv:1912.00818*, 2019.
- 353 [20] Filip Hanzely, Slavomír Hanzely, Samuel Horváth, and Peter Richtárik. Lower bounds and
354 optimal algorithms for personalized federated learning. In *Proceedings of Advances in Neural*
355 *Information Processing Systems 33: Annual Conference on Neural Information Processing*
356 *Systems, NeurIPS*, 2020.
- 357 [21] Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Adaptive personalized
358 federated learning. *arXiv preprint arXiv:2003.13461*, 2020.
- 359 [22] Yihan Jiang, Jakub Konečný, Keith Rush, and Sreeram Kannan. Improving federated learning
360 personalization via model agnostic meta learning. *arXiv preprint arXiv:1909.12488*, 2019.
- 361 [23] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adap-
362 tation of deep networks. In *Proceedings of the 34th International Conference on Machine*
363 *Learning, ICML*, volume 70, pages 1126–1135, 2017.
- 364 [24] Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S. Talwalkar. Federated multi-
365 task learning. In *Proceedings of Advances in Neural Information Processing Systems 30: Annual*
366 *Conference on Neural Information Processing Systems, NeurIPS*, 2017.
- 367 [25] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network.
368 *arXiv preprint arXiv:1503.02531*, 2015.
- 369 [26] Ying Zhang, Tao Xiang, Timothy M Hospedales, and Huchuan Lu. Deep mutual learning. In
370 *Proceedings of 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR*,
371 pages 4320–4328, 2018.
- 372 [27] Jang Hyun Cho and Bharath Hariharan. On the efficacy of knowledge distillation. In *Proceedings*
373 *of 2019 IEEE/CVF International Conference on Computer Vision, ICCV*, pages 4794–4802,
374 2019.
- 375 [28] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student
376 improves imagenet classification. In *Proceedings of 2020 IEEE/CVF Conference on Computer*
377 *Vision and Pattern Recognition, CVPR*, pages 10687–10698, 2020.
- 378 [29] Chenglin Yang, Lingxi Xie, Siyuan Qiao, and Alan L Yuille. Training deep neural networks
379 in generations: A more tolerant teacher educates better students. In *Proceedings of the AAAI*
380 *Conference on Artificial Intelligence, AAAI*, volume 33, pages 5628–5635, 2019.
- 381 [30] Mary Phuong and Christoph H Lampert. Distillation-based training for multi-exit architectures.
382 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1355–
383 1364, 2019.
- 384 [31] Guile Wu and Shaogang Gong. Peer collaborative learning for online knowledge distillation. In
385 *Proceedings of the AAAI Conference on Artificial Intelligence, AAAI*, 2021.
- 386 [32] Qiushan Guo, Xinjiang Wang, Yichao Wu, Zhipeng Yu, Ding Liang, Xiaolin Hu, and Ping
387 Luo. Online knowledge distillation via collaborative learning. In *Proceedings of the IEEE/CVF*
388 *Conference on Computer Vision and Pattern Recognition, CVPR*, pages 11020–11029, 2020.
- 389 [33] Ilai Bistriz, Ariana Mann, and Nicholas Bambos. Distributed distillation for on-device learning.
390 In *Proceedings of Advances in Neural Information Processing Systems 33: Annual Conference*
391 *on Neural Information Processing Systems, NeurIPS*, 2020.

- 392 [34] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. EMNIST: extending
393 MNIST to handwritten letters. In *Proceedings of 2017 International Joint Conference on Neural
394 Networks, IJCNN*, 2017.
- 395 [35] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for
396 benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 397 [36] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
398 2009.
- 399 [37] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning
400 applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 401 [38] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep
402 convolutional neural networks. In *Proceedings of Advances in Neural Information Processing
403 Systems 25: 26th Annual Conference on Neural Information Processing Systems, NeurIPS*,
404 pages 1106–1114, 2012.
- 405 [39] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
406 recognition. In *Proceedings of 2016 IEEE Conference on Computer Vision and Pattern
407 Recognition, CVPR*, pages 770–778, 2016.
- 408 [40] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines
409 for efficient cnn architecture design. In *Proceedings of the European conference on computer
410 vision, ECCV*, pages 116–131, 2018.

411 Checklist

- 412 1. For all authors...
- 413 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
414 contributions and scope? [Yes]
- 415 (b) Did you describe the limitations of your work? [N/A]
- 416 (c) Did you discuss any potential negative societal impacts of your work? [No] We believe
417 our work has the potential to be implemented in practical scenarios.
- 418 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
419 them? [Yes]
- 420 2. If you are including theoretical results...
- 421 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 4
- 422 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix
- 423 3. If you ran experiments...
- 424 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
425 mental results (either in the supplemental material or as a URL)? [Yes] See Appendix
- 426 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
427 were chosen)? [Yes] See Appendix
- 428 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
429 ments multiple times)? [N/A]
- 430 (d) Did you include the total amount of compute and the type of resources used (e.g., type
431 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 5
- 432 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 433 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 434 (b) Did you mention the license of the assets? [N/A]
- 435 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 436
- 437 (d) Did you discuss whether and how consent was obtained from people whose data you’re
438 using/curating? [N/A]

- 439 (e) Did you discuss whether the data you are using/curating contains personally identifiable
440 information or offensive content? [N/A]
- 441 5. If you used crowdsourcing or conducted research with human subjects...
- 442 (a) Did you include the full text of instructions given to participants and screenshots, if
443 applicable? [N/A]
- 444 (b) Did you describe any potential participant risks, with links to Institutional Review
445 Board (IRB) approvals, if applicable? [N/A]
- 446 (c) Did you include the estimated hourly wage paid to participants and the total amount
447 spent on participant compensation? [N/A]