
Few-Round Learning for Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In federated learning (FL), a number of distributed clients targeting the same task
2 collaborate to train a single global model without sharing their data. The learning
3 process typically starts from a randomly initialized or some pretrained model.
4 In this paper, we aim at *designing an initial model* based on which an arbitrary
5 group of clients can obtain a global model for its own purpose, within only a few
6 rounds of FL. The key challenge here is that the task of the group conducting FL
7 are generally unknown when the initial model is prepared. Our idea is to take
8 a meta-learning approach to construct the initial model so that any group with a
9 possibly unseen task can obtain a high-accuracy global model within only R rounds
10 of FL. Our meta-learning itself could be done via federated learning among willing
11 participants and is based on an episodic arrangement to mimic the R rounds of
12 FL followed by inference in each episode. Extensive experimental results show
13 that our method generalizes well for arbitrary groups of clients and provides large
14 performance improvements given the same overall communication/computation
15 resources, compared to other baselines relying on known pretraining methods (e.g.,
16 federated meta-learning ideas targeting personalization).

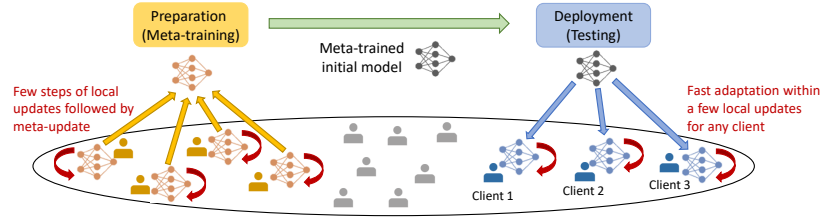
17 1 Introduction

18 Rapidly growing amounts of valuable data are being collected at distributed edge nodes such as
19 mobile phones, wearable client devices and smart vehicles/drones. Directly sending these local data
20 to the central server for model training raises significant privacy concerns. To address this issue, an
21 emerging trend known as federated learning (FL) [13, 9, 1, 11, 20, 16, 15], where server uploading
22 of local data is not necessary, has been actively researched. In FL, a large group of distributed clients
23 interested in solving the same task (e.g., classification on given categories of images) collaborate in
24 training a single global model without sharing their data. While standard supervised learning uses
25 some dataset D to find the model ϕ that would minimize a loss function $f(\phi, D)$, FL in comparison
26 seeks the model ϕ that minimizes the averaged version of the local losses $f(\phi, D_k)$, computed at
27 each node k using local data D_k . The learning process typically starts from a randomly initialized or
28 some pretrained model and is carried out through iterative aggregation of the local model updates.

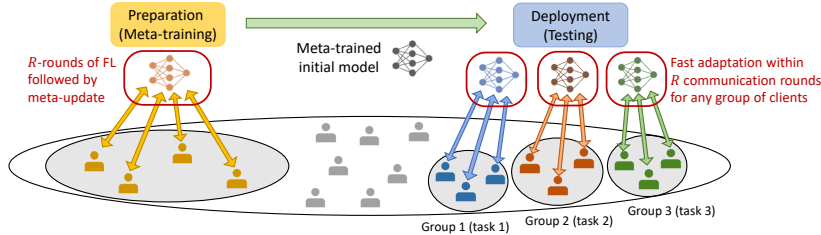
29 1.1 Backgrounds and Main Contributions

30 **Motivation.** Unfortunately, FL generally requires a large number of communication rounds between
31 the server and the clients for model exchange, to achieve a desired level of performance. This makes
32 the implementation of FL a significant challenge in bandwidth-limited or time-sensitive applications.
33 Especially in real-time applications (e.g., connected vehicles or drones) where the model should
34 quickly adapt to dynamically evolving environments, the requirement on many communication rounds
35 becomes a major bottleneck.

36 **Goal and challenge.** To tackle this problem from the service provider’s perspective, we aim to
37 *prepare an initial model* that can quickly adapt to any group (focusing on its own task) of clients
38 within only a few rounds of FL. The key challenge here is that the task of the group conducting



(a) **Personalized FL:** meta-training geared to few steps of local updates at any client.



(b) **Proposed few-round FL:** meta-training geared to few rounds of FL at any group.

Figure 1: Basic concept of personalized FL and proposed few-round learning method. In the meta-training phase (or preparation stage) of our scheme, the service provider prepares the initial model with willing participants. Once this preparation is over, the service provider offers this initial model to any group of clients (possibly including the meta-training participants) who hope to classify an arbitrary task within few-round FL. In order to facilitate R -round FL at deployment, we take an episodic training strategy that mimics actual inference preceded by R FL rounds in the meta-training phase with a new set of participants chosen in each episode.

39 FL is generally not known when the service provider prepares the initial model. In the context of
 40 classification, different tasks mean classification involving different sets of classes. For example,
 41 classifying diseases A, B, C (task 1) is a different task compared to classifying diseases D, E, F (task
 42 2). Since the group conducting FL can include classes that are unseen during preparation, *existing FL*
 43 *approaches cannot tackle this problem.*

44 **Key idea.** Our key idea is to adopt meta-learning (which enables reliable prediction even when the
 45 task at inference is unseen when the model was meta-trained) to prepare the initial model that enables
 46 *few-round FL*. Once meta-training is over, the service provider would offer the trained model to
 47 some clients who want to solve a common task after collaborating through a quick few rounds of
 48 FL. These clients may or may not be the participants of the earlier meta-training phase, and their
 49 classification task is generally considered unseen during meta-training. A high-level description of
 50 our idea is depicted in Fig. 1(b). Given a small target value R , we take an episodic training approach
 51 to enable R -round FL for any group of clients. In essence, we find the initial model ϕ that would
 52 minimize the average of local losses $f(\theta^R(\phi), D_k)$, where $\theta^R(\phi)$ is the model to be updated from ϕ
 53 through R rounds of FL among future clients in the deployment stage. Despite the high practical
 54 significance of this problem formulation, to the best of our knowledge, this is the first work to propose
 55 a meta-learning strategy geared to *few-round FL*. It is also worth mentioning that model preparation
 56 is not a real-time requirement and can often be done when bandwidth demands are sparse.

57 **Comparison with personalized FL.** We stress that our idea has a different purpose and approach
 58 relative to the recent line of works on federated meta-learning [12, 4], which initiate a model for
 59 personalized optimizations at local clients (see Fig. 1(a)). The goal of these approaches is to obtain a
 60 personalized local model at each client within a few steps of gradient descents, in the deployment
 61 stage. To achieve this goal, in the preparation stage, a few steps of local updates and meta-update are
 62 first performed at each participant independently (with its own local data), and FL (or aggregation) is
 63 adopted just to take advantage of data of various participants: these approaches seek ϕ that minimizes
 64 the average of local losses $f(\theta_k(\phi), D_k)$, where $\theta_k(\phi)$ is the local model updated from ϕ through a
 65 number of gradient steps using local data D_k . In contrast to personalized FL that focuses on local
 66 client models in the deployment stage, our *few-round learning* inherit the ability of FL at deployment
 67 to obtain a *global model*. Hence, for our scheme, it is inevitable to adopt FL in the preparation stage
 68 to mimic the R -round FL scenario at deployment; in the preparation stage, meta-update is performed
 69 at each participant after they collaboratively perform R FL rounds. To sum up, our approach aims to
 70 prepare an initial model that leads to a *global model* within “few rounds of FL”, while personalized
 71 FL aims for an initial model leading to *personalized model* within “few steps of local updates” only
 72 using its own data. These are obviously two completely different problems with distinct solutions.

73 **Main contributions.** Technically, we utilize a model-agnostic meta-learning (MAML) approach to
74 prepare the initial model via episodic training strategy. While directly applying MAML independently
75 to each local model leads to existing solutions on personalized FL [12, 4], in our approach, R rounds
76 of local updates and aggregations are first performed in each episode before the meta-update process.
77 This unique episode construction compared to personalized FL methods mimics the deployment
78 stage where actual inference is preceded by an R -round FL procedure. Another key ingredient in our
79 solution is to adopt *prototype aggregation* in each FL round to construct global prototypes that serve
80 as better class representatives compared to the locally computed prototypes, in learning embedding
81 space. This strategy is especially effective when a non-IID (independent, identically distributed) data
82 distribution across clients tends to induce a significantly biased model after performing local updates.
83 The global prototypes serve as prior knowledge, a form of regularization, and prevent local models
84 from overfitting to the local data. Moreover, the global prototypes (reflecting all classes across clients)
85 can assist the local model to learn a more general embedding space. We call this approach a global
86 prototype-assisted learning (GPAL) strategy. Our main contributions are summarized as follows:

- 87 • We formulate a **new problem of high practical significance**, namely, **few-round learning**,
88 where the goal is to prepare an initial model that can quickly adapt to any group of clients within
89 only a few rounds of FL.
- 90 • We propose a **meta-training algorithm specifically geared to R rounds of FL** followed by
91 inference, to be performed by a group of clients on a possibly unseen task.
- 92 • We **guarantee convergence** of our meta-training algorithm via theoretical analysis.
- 93 • We show via experiments that our scheme **outperforms existing pretraining approaches**
94 including fine-tuning via FedAvg and personalized FL in both IID and non-IID scenarios.

95 1.2 Related Works

96 **Few-shot learning.** Few-shot learning is an instantiation of meta-learning. In the context of image
97 classification, few-shot learning typically involves episodic training where each episode of training
98 data is arranged into a few training (support) sample images and validation (query) samples to mimic
99 inference that uses only a few examples [19]. Through a repetitive exposure to a series of varying
100 episodes with different sets of image classes, the model learns to handle new tasks (classification
101 against unseen classes) each time. Two widely-known few-shot learning methods with different
102 philosophical twists, which are also conceptually relevant to the present work, are MAML [5] and
103 Prototypical Networks [18]. MAML attempts to generate an initial model from which different
104 models targeting different tasks can be obtained quickly via just a few gradient updates. The idea is
105 that the initial model is learned via meta-training to develop an internal representation that is close in
106 some sense to a variety of unseen tasks. Prototypical Networks, on the other hand, learn embedding
107 space such that model outputs cluster around class prototypes, the class-specific centroids of the
108 embedder outputs. With episodic training, simple Prototypical Networks are surprisingly effective in
109 learning inductive bias for successful generalization to new tasks.

110 We stress that our few-round learning scheme (that targets few global rounds of FL) has different
111 purpose and technical approach compared to the existing works on few-shot learning (that targets few
112 shots of data sample). Nevertheless, we take advantage of both concepts on MAML and Prototypical
113 Networks to achieve our own goal: we adopt MAML in updating the initial model specifically geared
114 to R -round FL, and adopt both *prototype aggregation* and *prototype-assisted learning* strategies to
115 learn a general embedding space and successfully handle the non-IID issue in FL.

116 **Federated meta-learning.** Recent research activity has focused on improving model personalization
117 via federated meta-learning [12, 3, 4, 7]. The common goal of these works is to generate an initial
118 model based on which each new client can find its own optimized model via a few local gradient steps
119 and using only its own data. In these works, meta-learning employed during federated learning intends
120 to enable each client to handle previously unseen tasks, in the spirit of MAML of [5]. User-specific
121 next-word prediction at individual smartphones, for example, is a possible application. Compared to
122 this line of work, we focus on creating an initial model that leads to a high-accuracy *global model*,
123 rather than personalized models. In this way, we seek to take advantage of a higher variety of data
124 as well as the larger data volume that would be made available through collaborative learning of a
125 group of distributed nodes. A clear example is the diagnosis of a broader class of diseases that would
126 be possible through collaborative training across more examples contributed by a larger group of
127 individuals. Personalized FL methods (e.g., [12, 4]) especially have disadvantage in non-IID settings
128 where each client necessarily lacks a sufficient variety of data. The results are reported in Section 4.

129 **One-shot FL.** Another line of work recently focused on one-shot FL, where the goal is to train a
 130 global model with just one communication round between the server and the clients. The authors of
 131 [6] proposed an ensemble method to choose reliable client-specific models from given clients. In
 132 the work of [17], local clients send XOR-encoded MNIST image data to the server, and the server
 133 decodes it to train the global model. While the server would need certain data in advance to decode
 134 the received results, XOR operation can serve as data augmentation while preserving privacy. In
 135 the fusion learning of [8], each local client uploads both the model parameters and the distribution
 136 parameters to the server. The server generates artificial data samples from the distribution parameters
 137 to train a global model. When the data gets complex, however, it is not clear whether conversion into
 138 a simple distribution would be reliable. Compared to the existing works on one-shot FL that employ
 139 some randomly initialized model, the key difference of our method is the use of *meta-learning* to
 140 prepare an initial model which can adapt to unseen tasks of individual groups’ of clients within R
 141 rounds of FL. The advantage of our scheme compared to these methods is shown in Section 4.

142 2 Proposed Few-Round Learning Algorithm

143 2.1 Problem Setup

144 **Federated learning.** Let N be the number of clients in the system. FL allows each distributed node
 145 k with a dataset D_k to participate in iterative learning of a global model θ without having to reveal its
 146 data to anyone else including the central server. As a given round r starts, each of K participating
 147 nodes (generally chosen anew every round) downloads a global model θ^r from the server and updates
 148 it using its own local data D_k . The updated local models θ_k^{r+1} get all uploaded to the server to be
 149 aggregated to a new model $\theta^{r+1} = \sum_{k=1}^K \mu_k \theta_k^{r+1}$, according to the relative dataset sizes $\mu_k = \frac{|D_k|}{\sum_{j=1}^K |D_j|}$.
 150 The same process gets repeated. FL generally requires a significant number of such global rounds to
 151 achieve the desired accuracy, with each round taking up substantial communication resources.

152 **Problem formulation.** In preparing an initial model ϕ for any group of clients to pursue a few FL
 153 rounds, we use meta-learning based on episodic training, where each episode is constructed to mimic
 154 R FL rounds followed by inference. Once meta-training is over, in the deployment phase, the service
 155 provider offers the trained initial model ϕ to any group of clients wishing to pursue inference on
 156 some common task (possibly unseen during meta-training) after collaborating for R rounds of FL.

157 2.2 Meta-Training (Preparation Stage)

158 More precisely stated, our meta-training phase is to find ϕ that minimizes the objective function

$$F(\phi) = \mathbb{E}_{A_t \sim p(\mathcal{A})} \left[\frac{1}{K} \sum_{k \in A_t} f(\theta^R(\phi), D_k) \right] \quad (1)$$

159 where A_t is a specific group with K participants drawn from $p(\mathcal{A})$, the distribution over all possible
 160 groups, each with K participants; $\theta^R(\phi)$ is the model after R rounds of FL in group A_t , starting from
 161 ϕ ; and D_k is the local dataset of participant k in group A_t . In comparison, the objective function for
 162 personalized FL methods (e.g., Per-FedAvg of [4]) is $F(\phi) = \frac{1}{N} \sum_{k=1}^N f(\theta_k(\phi), D_k)$ where N is the
 163 number of clients in the system and $\theta_k(\phi)$ is the model after a few gradient steps at client k starting
 164 from ϕ . We also reiterate that conventional FL aims at minimizing $F(\phi) = \frac{1}{N} \sum_{k=1}^N f(\phi, D_k)$.

165 To create a training environment matching the actual R -rounds of FL followed by inference at
 166 deployment, in each episode of our meta-training phase, we update the model over R federated
 167 rounds using the support set and then makes a final adjustment (meta-update) using the query set.
 168 This process is repeated as the model is exposed to a series of episodes.

169 The detailed procedure of our meta-training is given in Algorithm 1. For a quick summary, as each
 170 episode t begins, the server selects a new set of K participants. The model ϕ^t , carried over from the
 171 last episodic stage, becomes the initial model θ^0 for the current episode. After R rounds of FL with
 172 each round consisting of local updates via local support sets and a global aggregation, θ^0 evolves
 173 to θ^R . Before moving to the next episode, local meta-updates are done based on θ^R using the local
 174 query sets to adjust the initial model θ^0 , in the spirit of MAML. As these meta-updated models get
 175 aggregated to ϕ^{t+1} at the server, the new episode can begin.

176 2.2.1 R Rounds of Local Updates and Aggregations

177 In defining the loss function, we utilize the class prototypes and associated distance metric of [18],
 178 a proven method of simple yet effective learning of embedding space. For each communication

Algorithm 1 Proposed Meta-Training Algorithm for Few-Round Learning

Input: Initialized model ϕ^0 **Output:** Model ϕ^T after T training episodes

- 1: **for** each training episode $t = 0, 1, \dots, T - 1$ **do**
- 2: The server constructs a group $A_t \sim p(A)$ of K participants chosen out of N users.
- 3: Each participant $k \in A_t$ splits D_k into support set S_k and query set Q_k .
- 4: $\theta^0 \leftarrow \phi^t$
- 5: **for** each communication round $r = 0, 1, \dots, R - 1$ **do**
- 6: **for** each participant k **in parallel do**
- 7: Download θ^r and Γ^{r-1} from the server (download only θ^r when $r = 0$)
- 8: **for** each class $c \in C_k$ **do**
- 9: $\Gamma_k^r(c) = \frac{1}{|S_k(c)|} \sum_{x \in S_k(c)} g_{\theta^r}(x)$ // Local prototype calculation with support set S_k
- 10: **end for**
- 11: $\theta_k^{r+1} \leftarrow \theta^r - \alpha \nabla_{\theta^r} f(\theta^r, S_k)$ // Local update of θ with support set S_k and GPAL
- 12: **end for**
- 13: $\theta^{r+1} = \sum_{k=1}^K \lambda_k \theta_k^{r+1}$ // Model aggregation; λ_k is relative support set size
- 14: $\Gamma^r = \left\{ \sum_{k=1}^K \lambda_k \Gamma_k^r(c) \mid c = 1, 2, \dots, N_c \right\}$ // Prototype aggregation
- 15: **end for**
- 16: **for** each participant k **in parallel do**
- 17: Download θ^R, Γ^{R-1} from the server.
- 18: Compute local prototypes based on Q_k .
- 19: $\theta_k^0 \leftarrow \theta^0 - \beta \nabla_{\theta^R} f(\theta^R, Q_k)$ // Local meta-update of θ^0 with query set Q_k and GPAL
- 20: **end for**
- 21: $\phi^{t+1} = \sum_{k=1}^K \mu_k \theta_k^0$ // Aggregation of meta-updated models; μ_k is relative data size
- 22: **end for**

179 round r , we not only aggregate the global model θ^{r+1} but also the global prototypes $\Gamma^r = \{\Gamma^r(c) \mid c =$
 180 $1, 2, \dots, N_c\}$ for all classes, where N_c is the number of classes over all clients.

181 **Model and global prototype download.** In the beginning of round $r \geq 1$, the server has the global
 182 model θ^r and the global prototypes $\Gamma^{r-1} = \{\Gamma^{r-1}(c) \mid c = 1, 2, \dots, N_c\}$ from the previous round $r - 1$.
 183 Each participant k first downloads θ^r and Γ^{r-1} from the server. Since there is no global prototype in
 184 the first round, the participants only download the model θ^0 when $r = 0$.

185 **Local prototype calculation.** The local prototype of $\Gamma_k^r(c)$ for participant k is computed as in Line
 186 9 using the downloaded model θ^r , the associated embedder outputs g_θ corresponding to the local
 187 support samples $S_k(c)$ labeled c . This local prototype serves as a representative of class c calculated
 188 based on the local data (support set) of client k .

189 **Loss calculation from local prototypes.** Let Γ_k^r be the set of all classes of prototypes at participant
 190 k : $\Gamma_k^r = \{\Gamma_k^r(c) \mid c \in C_k\}$, where C_k is a set of all classes at participant k . Now using S_k, θ^r and Γ_k^r ,
 191 each participant k computes the local loss according to

$$L_{\text{local}}^{S_k}(\theta, \Gamma_k^r(c)) = \frac{1}{\sum_{c \in C_k} |S_k(c)|} \sum_{c \in C_k} \sum_{x \in S_k(c)} \left\{ d(g_\theta(x), \Gamma_k^r(c)) + \log \sum_{c' \neq c} \exp(-d(g_\theta(x), \Gamma_k^r(c'))) \right\}, \quad (2)$$

192 based on Euclidean distance $d(\cdot)$ between $\Gamma_k^r(c)$ and $g_\theta(x)$ for $x \in S_k(c)$.

193 **Auxiliary loss from global prototypes.** Relying only on the loss function of (2) based on the local
 194 prototype tends to bias the model, especially when data distributions across different clients are
 195 non-IID. This generally leads to a performance degradation of the global model. To get around, we
 196 propose a global prototype-assisted learning (GPAL) strategy, where the global prototypes serve
 197 as prior knowledge in a form of regularization to prevent local models from overfitting to their
 198 local data. Moreover, the global prototypes, reflecting classes not limited to the local dataset,
 199 can assist the local model to learn a more general embedding space. Given the global prototypes
 200 $\Gamma^{r-1} = \{\Gamma^{r-1}(c) \mid c = 1, 2, \dots, N_c\}$ and $\{g_\theta(x) \mid x \in S_k\}$, the auxiliary loss $L_{\text{aux}}^{S_k}(\theta^r, \Gamma^{r-1})$ can be
 201 computed by replacing local prototypes Γ_k^r with global prototypes Γ^{r-1} in (2).

202 **Local update based on GPAL.** Based on the local loss $L_{\text{local}}^{S_k}(\theta, \Gamma_k^r)$ computed using local prototypes
 203 and the auxiliary loss $L_{\text{aux}}^{S_k}(\theta, \Gamma^{r-1})$ based on global prototypes, the objective function becomes

$$f(\theta^r, S_k) = \gamma L_{\text{local}}^{S_k}(\theta^r, \Gamma_k^r) + (1 - \gamma) L_{\text{aux}}^{S_k}(\theta^r, \Gamma^{r-1}) \quad (3)$$

204 where γ is a balancing coefficient. For $r = 0$, we have $f(\theta^r, S_k) = L_{\text{local}}^{S_k}(\theta^r, \Gamma_k^r)$ since the global
 205 prototype is not defined in the first global round. Line 11 of Algorithm 1 performs local update
 206 accordingly, where α is the learning rate. We call this strategy GPAL.

207 In FL, the clients can perform multiple local updates, say E times. Hence, the process of local
 208 prototype computation in Line 9 of Algorithm 1, loss computation in (3) and local update of Line 11
 209 can be repeated E times to obtain θ_k^{r+1} .

210 **Model and prototype aggregations.** After performing local updates, each participant k sends its
 211 updated local model θ_k^{r+1} and the computed local prototypes Γ_k^r to the server. Then the server
 212 aggregates them according to Lines 13 and 14 in Algorithm 1, where the weighting factor $\lambda_k =$
 213 $\frac{|S_k|}{\sum_{j=1}^K |S_j|}$ reflects the relative support set sizes.

214 The above local update and global aggregation processes are repeated for R global rounds ($r =$
 215 $0, 1, \dots, R-1$), and the server obtains θ^R and Γ^{R-1} in a given episode.

216 2.2.2 One-Round Local Meta-Update and Aggregation

217 Towards the end of each episode processing stage, the participants download θ^R and Γ^{R-1} from
 218 the server. Each participant k uses its query set Q_k to compute the local prototypes Γ_k^R as in as in
 219 Line 9. The query loss $f(\theta^R, Q_k)$ is calculated similar to (3) based on Q_k , θ^R , Γ^{R-1} and Γ_k^R . The
 220 meta-update would call for taking the derivative of this loss with respect to θ^0 : $\nabla_{\theta^0} f(\theta^R, Q_k) =$
 221 $\nabla_{\theta^R} f(\theta^R, Q_k) \times \frac{\partial \theta^R}{\partial \theta^0} = \nabla_{\theta^R} f(\theta^R, Q_k) \times \left(\prod_{r=0}^{R-1} \sum_{j=1}^K \lambda_j \frac{\partial}{\partial \theta^r} (\theta^r - \alpha \nabla_{\theta^r} f(\theta^r, S_j)) \right)$. But one would
 222 need the double derivatives from other user locations as well, which is highly inconvenient. Ignoring
 223 the double derivative terms, we simply replace $\nabla_{\theta^0} f(\theta^R, Q_k)$ with $\nabla_{\theta^R} f(\theta^R, Q_k)$, as in Line 19.
 224 This is the same as making a first-order approximation to the MAML-like meta-update, as often
 225 done in the implementation of MAML variants including the original work of [5]. All our reported
 226 experimental results as well as convergence analysis in the present paper reflect this choice. The
 227 server finally aggregates the meta-updated models from all participants. The next episode begins as
 228 the server selects a new set of K participants.

229 2.3 Testing (Deployment Stage)

230 In the actual deployment or test phase, given a group of clients, the server sets $\theta^0 = \phi^T$ and then
 231 leads R rounds of FL to obtain θ^R and Γ^{R-1} . Now given a test sample, we make prediction based on
 232 θ^R and Γ^{R-1} : the model output is first computed using θ^R and then comparison is made with the
 233 distances from all global prototypes in Γ^{R-1} to reach a decision.

234 3 Convergence Analysis

235 We provide theoretical analysis to guarantee a certain convergence behavior for our meta-training
 236 algorithm for nonconvex loss functions $f_k(\phi) := f(\phi, D_k)$. We need the following assumptions
 237 commonly made in convergence analyses of FL involving meta-learning, e.g., [12, 4].

238 **Assumption 1.** For all i , f_i is L -smooth, i.e., $\|\nabla f_i(\phi_1) - \nabla f_i(\phi_2)\| \leq L \|\phi_1 - \phi_2\|$ for any ϕ_1, ϕ_2 .

239 **Assumption 2.** Let $l_i(\phi; x)$ be the loss function for a single data point $x \in D_i$ of participant i . For
 240 all $i = 1, 2, \dots, N$, the variance of the loss gradients across data samples at a given participant is
 241 bounded, i.e., $\mathbb{E}_{x \in D_i} [\|\nabla l_i(\phi; x) - \nabla f_i(\phi)\|^2] \leq V_d$ for any ϕ .

242 **Assumption 3.** Let $f(\phi) = \frac{1}{N} \sum_{i=1}^N f_i(\phi)$ be the average local loss of all participants in
 243 the system. The variance of the gradient of loss f_i across participants is bounded, i.e.,
 244 $\frac{1}{N} \sum_{i=1}^N \|\nabla f_i(\phi) - \nabla f(\phi)\|^2 \leq V_p$ for any ϕ .

245 Two key lemmas and a theorem below establish the convergence of our method. All proofs are in
 246 Supplementary Material.

247 **Lemma 1.** Assume that the learning rate α is in the range $(0, 1/L]$. Then, the global loss function
 248 $F(\phi)$ in (1) is L_F -smooth, where $L_F = L2^R$.

249 **Lemma 2.** Define the local loss of our scheme $F_k(\phi) := f_k(\theta^R(\phi))$ at participant k . For a group
 250 A with K clients, define the loss averaged within that group $\mathcal{F}_A(\phi) := \frac{1}{K} \sum_{k \in A} F_k(\phi)$. Assume
 251 $\alpha \in (0, 1/L]$. Then, the variance of the gradient of $\mathcal{F}_A(\phi)$ across groups is bounded as

$$|\mathcal{A}|^{-1} \sum_{A \in \mathcal{A}} \|\nabla \mathcal{F}_A(\phi) - \nabla F(\phi)\|^2 \leq V_p K^{-1} \quad (4)$$

252 where \mathcal{A} is the set of all possible groupings of K participants drawn from a pool of N individuals.
 253 **Theorem 1.** Suppose Assumptions 1, 2, 3 hold and $\alpha \in (0, 1/L]$. Let $|D|$ be the mini-batch size at
 254 the meta-update processes of all participants. Then, Algorithm 1 guarantees the following upper
 255 bound on the loss gradient associated with our learned model ϕ^T :

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla F(\phi^t)\|^2] \leq \frac{4(F(\phi^0) - F(\phi^*))}{\beta T} + \epsilon(\beta, R, |D|, K) \quad (5)$$

256 where ϕ^* is the optimal solution of (1) and $\epsilon(\beta, R, |D|, K) = \beta L 2^{R+2} (V_d |D|^{-1} + V_p K^{-1})$.
 257 As the number of episodes T increases, the upper bound of (5) settles to ϵ . For a given smoothness
 258 L , assumed loss gradient variance bounds (V_d, V_p) and a targeted number of FL rounds R , the error
 259 term ϵ is controlled by the meta-update learning rate β , the mini-batch size $|D|$ and the per-episode
 260 number of participants K . For any reasonable value R , practical choices of β , $|D|$ and K can make ϵ
 261 sufficiently small, as discussed in in Supplementary Material using representative parameter values.

262 4 Experiments

263 We validate our algorithm on CIFAR-100 [10], *miniImageNet* [19], FEMNIST[2]. Following the data
 264 splits in [14], for CIFAR-100 and *miniImageNet*, 100 classes are divided into 64 train, 16 validation
 265 and 20 test classes. For FEMNIST, we divide 62 classes into 52 alphabet (uppercase, lowercase) and
 266 10 digit classes. For each class of FEMNIST, we sort the images by its name and choose first 600
 267 samples. After all, we have 600 samples for each class in every dataset. 52 alphabet classes are set
 268 to train classes, while 10 digit classes are set to test classes. The train classes are utilized for the
 269 preparation stage, and the test classes are utilized at deployment to model the unseen tasks.

270 **Comparison schemes.** First, as a simplest baseline, we consider FedAvg [13], where a randomly
 271 initialized model is trained for R FL rounds at deployment. The preparation stage is not considered
 272 for this scheme. Thus, direct performance comparison would be obviously unfair for FedAvg, but we
 273 just want to show what kind of performance improvement is possible by meta-learned initialization
 274 versus random initialization. Second, we consider a FedAvg-based fine-tuning, where the model is
 275 first pretrained by conducting FedAvg in each episode during preparation, and then fine-tuned with
 276 new clients for R FL rounds via FedAvg at deployment. For example, in *miniImageNet*, a 64-way
 277 classifier model is first pretrained in the preparation stage. Next, the last linear layer is replaced by a
 278 Xavier-initialized layer, and then the overall model is fine-tuned to the group at deployment. We also
 279 consider fine-tuning based on one-shot FL [6], where the local models are sampled and aggregated by
 280 ensemble cross-validation. We allow a larger number of available clients (in the deployment stage)
 281 for this scheme to accommodate user sampling. The model is first pretrained via FedAvg during
 282 preparation, and then fine-tuned based on the scheme of [6] for R rounds at deployment. Finally,
 283 although comparison with personalized FL [4] is tricky as the goal is different, a global model can still
 284 be trained by repeating local updates and aggregations for R FL rounds starting from the initialized
 285 model geared to client personalization. We also provide comparison with this “forced” global model.
 286 For our few-round learning (FRL), we try both a linear classifier and a distance-based classifier [18]
 287 for comparison. For the linear classifier, we connect an additional linear layer behind CNN layers, as
 288 in other baselines. The distance-based classifier utilizes prototypes instead of using the linear layer.
 289 For the distance-based classifier that utilizes prototypes, we observe the effect of our GPAL strategy.

290 **Preparation stage.** We assume $N = 64$ participants in the system in the preparation stage for
 291 CIFAR-100 and *miniImageNet*. We assume $N = 52$ for FEMNIST. For every dataset, following
 292 [13], training data samples are prepared into $2N$ shards of 300 samples each, such that each shard
 293 corresponds to one image class. Each participant is given two shards, and these two shards may
 294 belong to either a common class or two distinct classes. This models non-IID data distributions
 295 across participants. To construct each episode, the server then randomly selects $K = 10$ out of N
 296 participants. Each participant uses one half of its local data from each class as support samples, and
 297 the remaining half as query samples. We typically set the target number of global rounds to $R = 3$.
 298 Each episode of our scheme requires 4 global rounds in the meta-training phase: 3 rounds of local
 299 updates and aggregation, and 1 round of local meta-update and aggregation. For a fair comparison,
 300 we let all baselines to consume the same amount of communication resources in the preparation
 301 stage: up to 40,000 communication rounds between the server and participants (other than FedAvg
 302 that employs no preparation rounds). Hence, our scheme is meta-trained over up to 10,000 episodes,
 303 taking 4 rounds in each episode. We also reiterate that model preparation at the service provider is
 304 not a real-time requirement and can be done when bandwidth demands are sparse; this offers an even
 305 more favorable performance/complexity tradeoff options for the proposed scheme.

Table 1: Performance with only unseen classes at deployment in an IID setup.

Methods	CIFAR-100	miniImageNet	FEMNIST
FedAvg	51.55 ± 0.38%	38.80 ± 0.26%	74.76 ± 0.35%
Fine-tuning via FedAvg	63.18 ± 0.41%	61.58 ± 0.47%	91.95 ± 0.28%
Fine-tuning via one-shot FL [6]	64.71 ± 0.37%	65.23 ± 0.43%	93.62 ± 0.26%
FRL: Linear classifier (Ours)	67.32 ± 0.37%	67.75 ± 0.35%	94.86 ± 0.13%
FRL: Distance-based classifier (Ours)	69.74 ± 0.31%	68.05 ± 0.34%	95.07 ± 0.10%
FRL: Distance-based classifier + GPAL (Ours)	72.93 ± 0.32%	69.31 ± 0.33%	96.61 ± 0.09%
Personalized FL: Linear classifier [4]	60.87 ± 0.31%	61.88 ± 0.32%	93.19 ± 0.12%
Personalized FL: Distance-based classifier	61.88 ± 0.32%	67.61 ± 0.31%	94.11 ± 0.12%

Table 2: Performance with only unseen classes at deployment in a non-IID setup.

Methods	CIFAR-100	miniImageNet	FEMNIST
FedAvg	34.85 ± 0.27%	29.74 ± 0.22%	59.22 ± 0.18%
Fine-tuning via FedAvg	44.33 ± 0.37%	33.39 ± 0.41%	58.23 ± 0.65%
Fine-tuning via one-shot FL [6]	35.11 ± 0.46%	27.16 ± 0.42%	57.88 ± 0.67%
FRL: Linear classifier (Ours)	52.98 ± 0.42%	53.51 ± 0.43%	85.14 ± 0.44%
FRL: Distance-based classifier (Ours)	63.85 ± 0.43%	61.07 ± 0.41%	88.60 ± 0.42%
FRL: Distance-based classifier + GPAL (Ours)	66.87 ± 0.40%	63.41 ± 0.39%	92.42 ± 0.32%
Personalized FL: Linear classifier [4]	51.54 ± 0.38%	52.42 ± 0.42%	80.59 ± 0.51%
Personalized FL: Distance-based classifier	58.11 ± 0.39%	55.83 ± 0.35%	88.07 ± 0.37%

306 **Deployment stage.** At deployment, we distribute the initial model obtained in the preparation stage
307 to a new group of clients. To measure the performance, we obtain the average test accuracy with
308 a 95% confidence interval over 1000 different groups (with $K = 10$ clients in each group) after R
309 rounds of FL. For testing, in one case we randomly sample τ classes from test classes that have not
310 been seen during preparation and distribute across $K = 10$ clients. In the other case, we randomly
311 sample τ_u classes from the unseen test classes and $\tau - \tau_u$ classes from the train classes seen during
312 meta-learning. We consider both IID and non-IID distributions. In the IID setup, the data samples
313 from each class are equally distributed to $K = 10$ clients. In the non-IID setup, we distribute data as
314 in the preparation stage. The support set is utilized for R FL rounds and the server calculates test
315 accuracy with the global model and the gathered query sets of all clients. For the one-shot FL scheme,
316 we allow 20 clients and the server samples $K = 10$ of them to aggregate. We focus on a 5-way setup
317 (i.e., $\tau = 5$) in the main paper with the $\tau = 10$ case reported in Supplementary Material.

318 **Implementation details.** The structure of the model follows the setting of [5] and [18], containing 4
319 consecutive 3×3 convolutional layers with 64 filters. Successively, each CNN output goes through
320 batch normalization, ReLU, and 2×2 max pooling. In the case of CIFAR-100 where the size of image
321 is 32×32 , the last two max pooling layers are omitted to up-scale the feature map. We adopt the SGD
322 optimizer with a learning rate of $\beta = 0.001$ for the meta-learner and a learning rate of $\alpha = 0.0001$ for
323 the learner. We set the mini-batch size to 60 and the number of local epochs at each client to $E = 1$.

324 **Results with unseen classes at deployment.** Tables 1 and 2 show test accuracies averaged over
325 1000 different groups after $R = 3$ global rounds at deployment, where the goal of each group is to
326 classify $\tau = 5$ classes that were unseen during preparation. First, it can be seen that FedAvg yields
327 significantly lower accuracy compared to others, as expected, since it uses a randomly initialized
328 model for training. By pretraining the model, FedAvg-based fine-tuning gives significant performance
329 gains compared to naive application of FedAvg, underlying the importance of initialization efforts.
330 The fine-tuning scheme based on one-shot FL shows further performance improvements in the IID
331 setup. However, since $K = 10$ clients are sampled from 20 clients for this method, there possibly
332 exist some unseen classes when building the global model in the non-IID setup, which lowers the

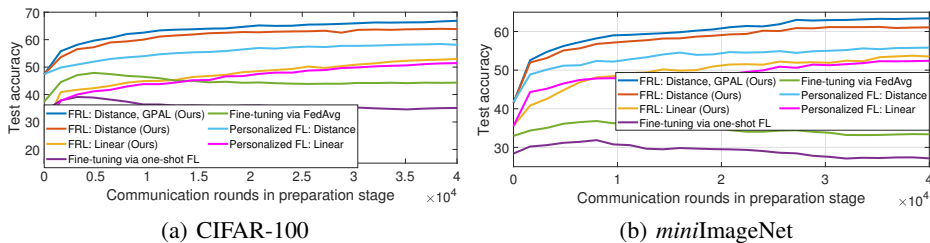


Figure 2: Final test accuracy at deployment, with varying numbers of communication rounds (proportional to the number of episodes for FRL and to the number of pretraining rounds for fine-tuning) in the preparation stage.

Table 3: Performance with both unseen/seen classes at deployment.

Methods	CIFAR-100		miniImageNet	
	IID	Non-IID	IID	Non-IID
FedAvg	50.03 ± 0.42%	34.82 ± 0.31%	42.17 ± 0.36%	30.37 ± 0.26%
Fine-tuning via FedAvg	66.73 ± 0.36%	44.46 ± 0.36%	63.82 ± 0.49%	36.18 ± 0.42%
Fine-tuning via one-shot FL [6]	69.84 ± 0.39%	35.33 ± 0.46%	67.05 ± 0.44%	29.12 ± 0.43%
FRL: Linear classifier (Ours)	68.22 ± 0.38%	53.62 ± 0.45%	69.02 ± 0.39%	55.18 ± 0.46%
FRL: Distance-based classifier (Ours)	70.49 ± 0.36%	65.13 ± 0.43%	70.39 ± 0.38%	62.42 ± 0.43%
FRL: Distance-based classifier + GPAL (Ours)	73.68 ± 0.37%	67.31 ± 0.44%	71.81 ± 0.34%	65.33 ± 0.42%
Personalized FL: Linear classifier [4]	65.09 ± 0.32%	52.08 ± 0.44%	62.05 ± 0.38%	53.53 ± 0.49%
Personalized FL: Distance-based classifier	68.70 ± 0.34%	57.62 ± 0.41%	63.63 ± 0.35%	58.08 ± 0.41%

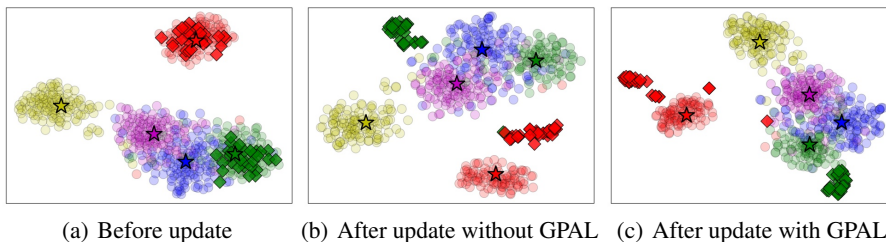


Figure 3: t-SNE visualization of embedding space at a client. The local data samples of the client is illustrated with diamond \diamond , and the data points of all other clients are represented by circle \circ . The global prototypes of each class are shown with \star . GPAL prevents the model from being biased to its local data and enables to learn more general embedding space. This leads to performance improvements as seen in Tables 1, 2, 3 and Fig. 2.

333 performance compared to fine-tuned FedAvg. Interestingly, the performance of personalized FL is
 334 better than the fine-tuning methods, although lagging well behind our methods. This latter observation
 335 is expected given the different design objectives. Our FRL algorithm performs the best, with the
 336 distance-based classifier showing better accuracy compared to the linear classifier. The relative gains
 337 of our methods for non-IID are particularly strong. It can be also seen that the performance of the
 338 global model can be further improved by our GPAL strategy. Fig. 2 shows how the final test accuracy
 339 (after 3 fixed FL rounds at deployment) improves with the number of communication rounds in the
 340 preparation stage. The overall results in Tables 1, 2 and Fig. 2 confirm the advantage of exploiting
 341 meta-learning and global prototype-assisted learning to facilitate few-round FL.

342 **Results with both unseen/seen classes at deployment.** In Table 3, we report test accuracies with
 343 both unseen/seen classes at deployment; the goal of each group is to classify $\tau = 5$ classes, 2 from the
 344 unseen classes and 3 from the seen classes. Since the tasks also handle classes already seen during
 345 preparation, the accuracies are generally higher than in Tables 1, 2. The trend is consistent with the
 346 results in Tables 1 and 2, confirming the advantage of the proposed algorithm.

347 **Effect of global prototype-assisted learning.** To understand the effect of our GPAL method further,
 348 we visualized t-SNE of the embedding space at a client in Fig. 3. CIFAR-100 is considered with
 349 each client having two classes in its local data in a non-IID setup. When only local prototypes are
 350 used for training as in Fig. 3(b), it can be seen that the two classes of the client form clusters without
 351 considering the data samples of other clients (but still well-separated). By considering the global
 352 prototypes (reflecting classes of all participants), in Fig. 3(c), the data points in the local client form
 353 clusters while staying away from all other global prototypes, a clearly desirable feat. This prevents
 354 the local model from being biased to its local data and enables the local model to learn a more general
 355 embedding space compared to the case in Fig. 3(b) considering only the local prototypes.

356 **Other experimental results.** Additional results on other settings including higher-way classification,
 357 larger group size and mismatched R are reported in Supplementary Material.

358 5 Conclusion

359 We proposed a meta-learning strategy to prepare an initial model geared to few-round federated
 360 learning. Given a group of clients with a new task, our meta-trained model generalizes well within
 361 only a few FL rounds. Convergence of our meta-training is guaranteed through theoretical analysis.
 362 Extensive experimental results confirm significant advantages of our idea over different baselines such
 363 as FedAvg-based fine-tuning and personalized FL in various setups. Our solution offers a promising
 364 direction for FL in practice, where minimizing training time and communication resources required
 365 in real-time is among key challenges.

References

- 366
- 367 [1] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir
368 Ivanov, Chloe Kiddon, Jakub Konecny, Stefano Mazzocchi, H Brendan McMahan, et al. Towards
369 federated learning at scale: System design. *arXiv preprint arXiv:1902.01046*, 2019.
- 370 [2] Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H Brendan
371 McMahan, Virginia Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings.
372 *arXiv preprint arXiv:1812.01097*, 2018.
- 373 [3] Fei Chen, Mi Luo, Zhenhua Dong, Zhenguo Li, and Xiuqiang He. Federated meta-learning
374 with fast convergence and efficient communication. *arXiv preprint arXiv:1802.07876*, 2018.
- 375 [4] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning
376 with theoretical guarantees: A model-agnostic meta-learning approach. *Advances in Neural
377 Information Processing Systems*, 33, 2020.
- 378 [5] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adapta-
379 tion of deep networks. In *ICML*, 2017.
- 380 [6] Neel Guha, Ameet Talwalkar, and Virginia Smith. One-shot federated learning. *arXiv preprint
381 arXiv:1902.11175*, 2019.
- 382 [7] Yihan Jiang, Jakub Konečný, Keith Rush, and Sreeram Kannan. Improving federated learning
383 personalization via model agnostic meta learning. *arXiv preprint arXiv:1909.12488*, 2019.
- 384 [8] Anirudh Kasturi, Anish Reddy Ellore, and Chittaranjan Hota. Fusion learning: A one shot
385 federated learning. In *International Conference on Computational Science*, pages 424–436.
386 Springer, 2020.
- 387 [9] Jakub Konecny, H. Brendan McMahan, Felix X. Yu, Ananda Theertha Suresh, and Dave Bacon.
388 Federated learning: strategies for improving communication efficiency. In *NIPS Workshop on
389 Private Multi-Party Machine Learning*, 2016.
- 390 [10] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
391 2009.
- 392 [11] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Chal-
393 lenges, methods, and future directions. *arXiv preprint arXiv:1908.07873*, 2019.
- 394 [12] Sen Lin, Guang Yang, and Junshan Zhang. A collaborative learning framework via federated
395 meta-learning. *arXiv preprint arXiv:2001.03229*, 2020.
- 396 [13] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
397 Communication-efficient learning of deep networks from decentralized data. In *Artificial
398 Intelligence and Statistics*, pages 1273–1282, 2017.
- 399 [14] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. 2017.
- 400 [15] Amirhossein Reisizadeh, Aryan Mokhtari, Hamed Hassani, Ali Jadbabaie, and Ramtin Pedarsani.
401 Fedpaq: A communication-efficient federated learning method with periodic averaging and
402 quantization. *arXiv preprint arXiv:1909.13014*, 2019.
- 403 [16] Felix Sattler, Simon Wiedemann, Klaus-Robert Müller, and Wojciech Samek. Robust and
404 communication-efficient federated learning from non-iid data. *IEEE transactions on neural
405 networks and learning systems*, 2019.
- 406 [17] MyungJae Shin, Chihoon Hwang, Joongheon Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun
407 Kim. Xor mixup: Privacy-preserving data augmentation for one-shot federated learning. *arXiv
408 preprint arXiv:2006.05148*, 2020.
- 409 [18] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In
410 *Advances in neural information processing systems*, pages 4077–4087, 2017.

- 411 [19] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks
412 for one shot learning. In *Advances in neural information processing systems*, pages 3630–3638,
413 2016.
- 414 [20] Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated
415 learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.

416 Checklist

- 417 1. For all authors...
- 418 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
419 contributions and scope? [Yes] We clearly stated our contributions and scope in both
420 abstract and introduction.
- 421 (b) Did you describe the limitations of your work? [No]
- 422 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 423 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them?
424 [Yes]
- 425 2. If you are including theoretical results...
- 426 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Assumptions
427 1, 2 and 3.
- 428 (b) Did you include complete proofs of all theoretical results? [Yes] All proofs are included in
429 Supplementary Material.
- 430 3. If you ran experiments...
- 431 (a) Did you include the code, data, and instructions needed to reproduce the main experimental
432 results (either in the supplemental material or as a URL)? [Yes] Our code and required
433 datasets are provided in Supplementry Material.
- 434 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were
435 chosen)? [Yes] See Section 4 for the details.
- 436 (c) Did you report error bars (e.g., with respect to the random seed after running experiments
437 multiple times)? [Yes] See our results in Tables 1, 2 and 3.
- 438 (d) Did you include the total amount of compute and the type of resources used (e.g., type
439 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4 and Supplementry
440 Material.
- 441 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 442 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 4.
- 443 (b) Did you mention the license of the assets? [No] Experiments were conducted on public
444 and uncommercial datasets.
- 445 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 446 (d) Did you discuss whether and how consent was obtained from people whose data you’re
447 using/curating? [No]
- 448 (e) Did you discuss whether the data you are using/curating contains personally identifiable
449 information or offensive content? [No]
- 450 5. If you used crowdsourcing or conducted research with human subjects...
- 451 (a) Did you include the full text of instructions given to participants and screenshots, if
452 applicable? [No]
- 453 (b) Did you describe any potential participant risks, with links to Institutional Review Board
454 (IRB) approvals, if applicable? [No]
- 455 (c) Did you include the estimated hourly wage paid to participants and the total amount spent
456 on participant compensation? [No]