

# Federated Mixture of Experts

Matthias Reisser<sup>1,2</sup> Christos Louizos<sup>2</sup> Efstratios Gavves<sup>1</sup> Max Welling<sup>1,2</sup>

## Abstract

Federated learning (FL) has emerged as the predominant approach for collaborative training of neural network models across multiple users, without the need to gather the data at a central location. One of the important challenges in this setting is data heterogeneity, *i.e.* different users have different data characteristics. For this reason, training and using a single global model might be suboptimal when considering the performance of each of the individual user’s data. In this work, we tackle this problem via Federated Mixture of Experts, `FedMix`, a framework that allows us to train an ensemble of specialized models. `FedMix` adaptively selects and trains a user-specific selection of the ensemble members. We show that users with similar data characteristics select the same members and therefore share statistical strength while mitigating the effect of non-i.i.d data. Empirically, we show through an extensive experimental evaluation that `FedMix` improves performance compared to using a single global model across a variety of different sources of non-i.i.d.-ness.

## 1. Introduction

An ever-increasing amount of devices are being connected to the internet, sensing their environment, and generating vast amounts of data. The term federated learning (FL) has been established to describe the scenario where we aim to learn from the data generated by this “federation” of devices (McMahan et al., 2016). Not only does the number of sensing devices increase, but also their processing power is increasing continuously to the point that it becomes viable to perform inference and training of machine learning models on device. In federated learning, the goal is to learn from these client devices’ data without collecting the data

<sup>1</sup>QUVA Lab, University of Amsterdam <sup>2</sup>Qualcomm AI Research. Correspondence to: Matthias Reisser <mreisser@qti.qualcomm.com>.

Arxiv preprint. Datasets were downloaded and experimented on by the first author at the University of Amsterdam.

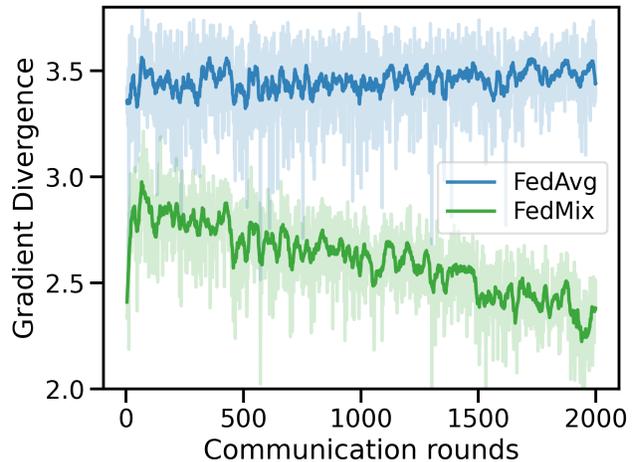


Figure 1. A sliding window of the gradient divergence (defined in Appendix D), on Cifar10 in the setup of Section 4 for `FedAvg` and `FedMix` ( $K = 4$ ).

centrally, which naturally allows for more private exchange of information.

Several challenges arise in the federated scenario. Federated devices are generally resource-constrained, both in their computational capacity as well as in communication bandwidth and latency. In a practical example, a smartphone has limited heat dissipation capacity and must communicate via Wi-Fi. From a global perspective, devices’ processing power and network connection can be highly heterogeneous across geographical regions and socio-economical status of device owners, causing practical issues (Bonawitz et al., 2019) and raising questions of fairness in FL (Li et al., 2019; Mohri et al., 2019). Apart from this cross-device setting, challenges in the so-called cross-silo setting focus on privacy concerns, while computation and communication constrains move to the background (Kairouz et al., 2019).

One of the key challenges in FL that we aim to address in this work is the non-i.i.d. nature of the shards of data that are distributed across devices. In non-federated machine learning, assuming independent and identically distributed data is generally justifiable and not detrimental to model performance. In FL however, each client performs a series of parameter updates on its own data shard to amortize the costs of communication. Over time, the direction of

progress across shards with non-i.i.d. data starts diverging (as shown in Figure 1), which can set back training progress, significantly slow down convergence and decrease model performance (Hsu et al., 2019).

To this end, we propose Federated Mixture of Experts (FedMix), an algorithm for FL that allows for training an ensemble of specialized models instead of a single global model. In FedMix, expert models are learning to specialize in regions of the input space such that, for a given expert, each client’s progress on that expert is aligned. FedMix allows each client to learn which experts are relevant for its shard and we show how it can be extended for inference on a previously unseen client. FedMix shows competitive performance against the established standard in FL, FedAvg (McMahan et al., 2016; Deng et al., 2020) across a range of visual classification tasks.

## 2. Federated Mixture of Experts

Federated learning (McMahan et al., 2016) deals with the problem of learning a server model with parameters  $\mathbf{w}$ , *e.g.*, a neural network, from a dataset of  $N$  datapoints  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$  that is distributed across  $S$  shards, *i.e.*,  $\mathcal{D} = \mathcal{D}_1 \cup \dots \cup \mathcal{D}_S$ , *without* accessing the shard-specific datasets directly. By defining a loss function  $\mathcal{L}_s(\mathcal{D}_s; \mathbf{w})$  per shard, the total risk can be written as

$$\arg \min_{\mathbf{w}} \sum_{s=1}^S \frac{N_s}{N} \mathcal{L}_s(\mathcal{D}_s; \mathbf{w}), \quad (1)$$

$$\mathcal{L}_s(\mathcal{D}_s; \mathbf{w}) := \frac{1}{N_s} \sum_{i=1}^{N_s} L(\mathcal{D}_{si}; \mathbf{w}). \quad (2)$$

It is easy to see that this objective corresponds to empirical risk minimization over the joint dataset  $\mathcal{D}$  with a loss  $L(\cdot)$  for each datapoint. In federated learning one is interested in reducing the communication costs; for this reason (McMahan et al., 2016) propose to do multiple gradient updates for  $\mathbf{w}$  in the inner optimization objective for each shard  $s$ , thus obtaining “local” models with parameters  $\mathbf{w}_s$ . These multiple gradient updates are denoted as “local epochs”, *i.e.*, amount of passes through the entire local dataset, with an abbreviation of  $E$ . Each of the shards then communicates the local model  $\mathbf{w}_s$  to the server and the server updates the global model at “round”  $t$  by averaging the parameters of the local models  $\mathbf{w}^t = \sum_s \frac{N_s}{N} \mathbf{w}_s^t$ . This constitutes federated averaging (FedAvg) (McMahan et al., 2016), the standard in federated learning.

One of the main challenges in federated learning is the fact that usually the data are non-i.i.d. distributed across the shards  $S$ , that is  $p(\mathcal{D}|s_i) \neq p(\mathcal{D}|s_j)$  for  $i \neq j$ . On the one hand, this can make learning a single global model from all of the data with the classical FedAvg problematic. On the other hand, there is one extreme that does not suffer

from this issue; learning  $S$  individual models, *i.e.*, only optimizing  $\mathbf{w}_s$  on  $\mathcal{D}_s$ . Although these individual models by definition do not suffer from non-i.i.d. data, clearly we should aim to do better and exchange meaningful information between clients to learn more robust and expressive models.

With FedMix, we propose to strike a balance between the two aforementioned extremes; learning a single global model and learning  $S$  individual models. For this reason, we revisit an old model formulation, the Mixture of Experts (MoE). The classical formulation of a MoE model (Jacobs et al., 1991; Jordan & Jacobs, 1994) contains a set of  $K$  experts and a gating mechanism that is responsible for choosing an expert for a given data-point. A MoE model for a data point  $(\mathbf{x}, y)$  can generally be described by

$$p_{\mathbf{w}_{1:K}, \theta}(y|\mathbf{x}) = \sum_{z=1}^K p_{\mathbf{w}_z}(y|\mathbf{x}, z) p_{\theta}(z|\mathbf{x}), \quad (3)$$

where  $z$  is a categorical variable that denotes the expert,  $\mathbf{w}_k$  are the parameters of expert  $k$  and  $\theta$  are the parameters of the selection mechanism.

The MoE was proposed as a model for datasets where different subsets of the data exhibit different relationships between input  $\mathbf{x}$  and output  $y$ . Instead of training a single global model to fit this relationship everywhere, each expert performs well on a different subset of the input space. The gating function models the decision boundary between input regions, assigning data-points from subsets of the input region to their respective experts.

In this work, we show that, in the federated scenario, subdividing the input region through a MoE can alleviate the consequences of non-i.i.d. data by aligning gradient updates across experts (Figure 1).

In Federated Mixture of Experts (FedMix) we enrich this model by conditioning the gating mechanism on the shard assignment  $s$ . Whatever characteristics make shard  $s$  different from other shards can manifest in learning a different, localized gating mechanism that does not need to be communicated to the server. In choosing  $K = 1$ , FedMix recovers the standard setting of federated averaging.  $K = S$  in combination with fixing  $p(z = s|\mathbf{x}, s) = 1$  recovers  $S$  independent models. From a global perspective, we are interested in maximizing the following single objective:

$$\sum_{s=1}^S \sum_{i=1}^{N_s} \log p_{\mathbf{w}_{1:K}, \theta_s}(y_{s,i}|\mathbf{x}_{s,i}, s) = \sum_{s=1}^S \sum_{i=1}^{N_s} \log \left[ \sum_{z=1}^K p_{\theta_s}(z|\mathbf{x}_{s,i}, s) p_{\mathbf{w}_z}(y_{s,i}|\mathbf{x}_{s,i}, z) \right] \quad (4)$$

While it is possible to optimize Eq. 4 directly, we have found empirically that it is hard to achieve both: avoiding collapse to a single expert, thus obtaining FedAvg, and specialization of the experts. Instead, we propose to form a variational lower-bound on Eq. 4 with a global variational approximation  $q_\phi(z|\dots)$  to the true posterior  $p(z|\mathbf{x}, y, s)$  with parameters  $\phi$ . At test time,  $p(y|\mathbf{x}^*, s) = \sum_{k=1}^K p(y|\mathbf{x}^*, z)p(z|\mathbf{x}^*, s)$  can be readily evaluated without requiring  $q$ . This allows us to condition  $q_\phi(z|\dots)$  on any available side-information *at training time* that might result in better specialization in the non-i.i.d. federated scenario. In this paper, we discuss several sources of non-i.i.d.-ness: label skew, *i.e.*, different distributions  $p(y|s)$  per shard, input-transformations, *i.e.*, different  $p(x|s)$ , and different mappings  $p(y|x, s)$ , such as label permutations.

In practice, other or additional known sources of misalignment could be included to further improve this approximation, such as a manufacturer-id for a medical device in a medical scenario, a geographic identifier, or general domain-specific information. We show a variety of artificial scenarios in the experimental section. For exposition in this paper, we always use  $q_\phi(z|y)$ , however the formulae and algorithms are equally applicable to other side-information. The lower bound to be maximized in FedMix therefore is as follows:

$$\sum_{s=1}^S \sum_{i=1}^{N_s} \log p_{\mathbf{w}_{1:K}, \theta_s}(y_{s,i}|\mathbf{x}_{s,i}, s) \geq \sum_{s=1}^S \sum_{i=1}^{N_s} \mathbb{E}_{q_\phi(z|y_{s,i})} [\log p_{\mathbf{w}_z}(y_{s,i}|\mathbf{x}_{s,i}, z)p_{\theta_s}(z|\mathbf{x}_{s,i}, s)] + \beta H(q_\phi(z|y_{s,i})), \quad (5)$$

where  $\beta$  is a hyperparameter that controls the entropy of the approximate posterior distribution; a low  $\beta$  will result in  $q_\phi(z|y)$  that are more concentrated around their most probable value whereas higher values will encourage more uncertain distributions. While the parameters in Eq. 5 can be optimized via gradient descent, a closed-form update for the parameters  $\phi$  based on Lagrange multipliers exists if each conditional is parameterized directly, *i.e.*  $\phi_c = q(z|y = c) \in [0, 1]^K$ . The solution for the probabilities of each expert  $k$  conditioned on a given category  $c$  becomes

$$\phi_{c,k} = \frac{\left( \prod_{i=1}^{N_{s,c}} p(y_{s,i}, z = k|\mathbf{x}_{s,i}, s) \right)^{\frac{1}{N_{s,c}\beta}}}{\sum_{z=1}^K \left( \prod_{i=1}^{N_{s,c}} p(y_{s,i}, z = k|\mathbf{x}_{s,i}, s) \right)^{\frac{1}{N_{s,c}\beta}}} \quad (6)$$

$$= \frac{\exp\left(\frac{1}{N_{s,c}\beta} \sum_{i=1}^{N_{s,c}} \log p(y_{s,i}, z = k|\mathbf{x}_{s,i}, s)\right)}{\sum_{z=1}^K \exp\left(\frac{1}{N_{s,c}\beta} \sum_{i=1}^{N_{s,c}} \log p(y_{s,i}, z = k|\mathbf{x}_{s,i}, s)\right)}, \quad (7)$$

*i.e.*, a softmax where the logits correspond to the average of the log-joint probabilities of the datapoints that belong to class  $y = c$  and the prior probabilities of expert  $k$  in shard

$s$ . We have omitted the dependence on parameters  $\mathbf{w}_z, \theta_s$  for clarity. The full derivation is given in Appendix B. It is interesting to see that the entropy hyperparameter  $\beta$  acts as a “temperature” for the softmax, thus directly encouraging low or high entropy solutions for  $q_\phi(z|y)$ . This is in-line with our prior discussion about  $\beta$ . In addition to the closed-form solution, such a parameterization is efficient w.r.t. communication overhead. Conditioning  $q_\phi(z|y)$  on  $s$  is possible and results in localized approximations with parameters  $\phi_s$  that do not need to be communicated, however we found a global approximation to help align the gating mechanisms across shards.

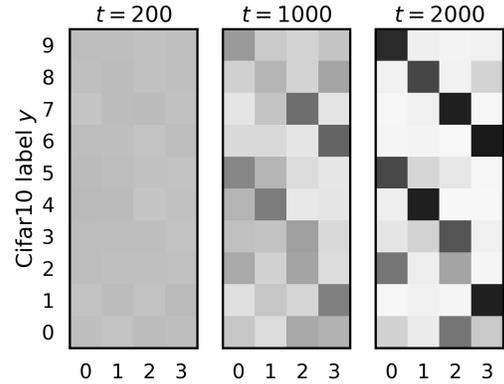


Figure 2. Visualization of  $q_\phi(z|y)$  at different communication rounds  $t$  for FedMix with  $K = 4$  on Cifar10. Greyscale corresponds to probabilities; white corresponds to zero and black corresponds to one. Probabilities sum to one across experts (horizontally).

**Expert Specialization** Specialization of the experts is a key ingredient for FedMix to be successful; with specialization, the gradients for each expert become aligned across shards (see Figure 1) and the (training set) performance in general improves.

Nevertheless, we find a fundamental trade-off in the MoE formulation: highly specialized experts are useful only if the local gating networks make correct selections, which is not always the case for test data. Therefore, a successful application of the MoE formulation requires striking a balance between the specialisation of the experts and their robustness to being wrongly selected by the routing mechanism. During training, experts initially receive gradients from all data-points until  $q_\phi(z|y)$  concentrates and enforces specialisation. We can thus control the speed of specialisation by tuning  $\beta$  and performing dampening on the update of the probabilities of  $q_\phi(z|y)$ . We empirically find that this leads to experts that, while they do not specialize as aggressively, can provide reasonable predictions even on datapoints outside of their “expertise”, thus improving the

overall performance of FedMix.

A possible drawback of specialization is that sometimes FedMix prematurely completely prunes experts, *i.e.*,  $p_{\theta_s}(z = k|\mathbf{x}, s) \approx 0 \forall \mathbf{x}, s$ . This can be undesirable as we lose model capacity that can be used for better modeling the data. As  $q_\phi(z|y)$  is one of the main training signals of  $p_{\theta_s}(z|\mathbf{x}, s)$ , we introduce the marginal entropy term in the server,  $H(E_{p(y)}[q_\phi(z|y)])$ , as a regularizer that encourages using all of the experts. Figure 2 show how  $q_\phi(z|y)$  converges over time towards specialization of experts.

**Personalization** Each of the specialized expert models that FedMix provides can be thought as containing information about the specific subset of the clients that selects each particular expert. As a result, these experts can serve as a better starting point for personalization according to the data on a specific device. Personalization can be achieved by simply finetuning the models obtained from the server, *i.e.*,  $\mathbf{w}_{1:K}$ , on the client specific training set  $\mathcal{D}_{s,\text{train}}$  thus obtaining  $\mathbf{w}_{1:K}^s$ . Under the assumption that  $p(\mathcal{D}_{\text{train}}|s) \approx p(\mathcal{D}_{\text{test}}|s)$ , which is not unreasonable in the federated learning scenario where each device has its own data generating mechanism, such personalized models can then have better prediction capabilities on the test data of that device. Finetuning is performed by optimizing Eq. 5 for a small number of steps (*e.g.*,  $E = 1$ ) with respect to  $\mathbf{w}_{1:K}$ ,  $\phi$  and  $\theta_s$ . It should be mentioned that this finetuning procedure is not limited to FedMix and can be performed on any federated learning method that involves global parameters shared by all of the clients, such as FedAvg.

During our experiments, we measure the performance of FedMix by measuring the average local accuracy of these client specific personalized models on their respective client specific test sets. In order to avoid the extra finetuning step for each round, we use the last version of the model that each client communicated to the server as the personalized model. For fair comparisons, we employ the same finetuning procedure for all of our other baselines as well.

**Server Side Updates** In a general federated learning algorithm, a central server selects a subset  $S' \subset \{1, \dots, S\}$  of clients at time  $t$  and transmits the current estimate of the global parameters  $\mathbf{w}^t$  to them. These clients perform a series of mini-batch gradient updates with data from their shard  $\mathcal{D}_s$  on a local loss function, which can come at the price of each client moving in possibly different directions in parameter space. In generalized FedAvg (Reddi et al., 2020), the server interprets  $\Delta_s^t = \mathbf{w}^t - \mathbf{w}_s^{t+1}$  as a single-step gradient update from client  $s$ , averages those gradients and applies an optimizer such as Adam (Kingma & Ba, 2014) to receive  $\mathbf{w}^{t+1}$ . In light of non-i.i.d. data across clients, this averaging strategy can result in slow progress since averaging updates in a highly non-convex parameter

space can be sub-optimal. In FedMix, this effect is mitigated since for a given expert, the data that is used to update its parameters are aligned better across shards.

FedMix offers a second way to improve convergence speed by modifying the server-side updates. In generalized FedAvg, the individual gradients returned by the subset  $S'$  of clients are averaged according to

$$\Delta^t = \sum_{s=1}^{S'} p(s) \cdot \Delta_s^t, \quad p(s) = \frac{N_s}{N_{S'}}. \quad (8)$$

In FedMix, we can speed up convergence by considering expert-specific updates  $\Delta_{k,s}^t = \mathbf{w}_k^t - \mathbf{w}_{k,s}^{t+1}$ . If a client  $s$  pruned away expert  $k$  from its local gating mechanism,  $\Delta_{k,s}^t$  will be zero. We propose to normalize the effective magnitude of the resulting update  $\Delta_k$  by up-weighting the updates of all other clients that do consider expert  $k$  for their local mixture:

$$\Delta_k^t = \sum_{s=1}^{S'} p(s|z=k) \cdot \Delta_{k,s}^t, \quad p(s|z=k) \propto p(z=k|s)p(s) \quad (9)$$

Computing  $p(z|s) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_s}[p_{\theta_s}(z|\mathbf{x})]$  prior to sending updates to the server involves evaluating potentially large neural network models which might not be desirable, depending on the situation and size of the local dataset. Therefore we approximate  $p(z|s) \approx q_\phi(z|s) = \mathbb{E}_{y \sim \mathcal{D}_s}[q_\phi(z|y)]$ , which involves just a single matrix multiplication.

**Privacy implications** The update rule described in Eq. 9 requires access to the marginal  $q(z|s) = \sum_y p(y|s)q_\phi(z|y)$  at the server. At the same time, the server has access to the parameters  $\phi$  that were used in computing  $p(z|s)$  before being sent to the server. Therefore, in principle, it could solve  $q(z|s) = \sum_y p(y|s)q_\phi(z|y)$  with respect to  $p(y|s)$  and thus obtain the marginal label distribution at the client. In practice this is not as straightforward to do as the probability matrix  $q_\phi(z|y)$  is not always invertible and solutions that use the pseudo-inverse, empirically, are not very accurate in capturing the entire distribution. With the additional constraints that the marginal needs to sum to one, contains only positive elements and that  $N_s \cdot p(y|s) \in \mathbb{Z}$ , in some cases, a reconstruction can become possible. As the number of classes exceeds the number of experts, success becomes more unlikely. We leave a thorough characterization of these properties to future work and discuss empirically in Appendix C the danger of leaking the marginal label distribution in FedMix and FedAvg.

**Communication costs** Reducing the amount of communication via the internet plays a central role, especially in the cross-device setting of FL. FedMix directly increases communication by a factor of  $K$  and is therefore more suited to the cross-silo setting of FL. Nevertheless, FedMix offers

**Algorithm 1** The FedMix algorithm.  $\alpha, \beta$  are the client and server learning rates;  $\gamma$  is a dampening factor and  $Z$  a normalization constant.

---

```

function SERVER SIDE
    Initialize  $\phi$  and  $K$  vectors  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_K]$ 
    for round  $t$  in  $1, \dots, T$  do
         $S' \leftarrow$  random subset of the clients
        Initialize  $\Delta_{\mathbf{W}}^t = \mathbf{0}, \Delta_{\phi}^t = \mathbf{0}$ 
        for  $s$  in  $S'$  do
             $\mathbf{W}_s^t, \phi_s^t, p(z|s) \leftarrow$  CLIENT SIDE( $s, \phi, \mathbf{W}$ )
        end for
         $p(s|z) \leftarrow p(z|s)p(s) / \sum_{s \in S'} p(z|s)p(s)$ 
        for  $s$  in  $S'$  do
             $\Delta_{\mathbf{w}_k}^t + = p(s|z = k)(\mathbf{w}_k^{t-1} - \mathbf{w}_{s,k}^t) \quad \forall k$ 
             $\Delta_{\phi}^t + = \frac{N_s}{N_{S'}}(\phi^{t-1} - \phi_s^t)$ 
        end for
         $\Delta_{\phi}^t - = \nabla_{\phi} H(\sum_c q_{\phi}(z|y=c)p(y=c))$ 
         $\mathbf{w}_{1:K}^{t+1} \leftarrow$  ADAM( $\Delta_{\mathbf{w}_{1:K}}^t, \beta$ )
         $\phi^{t+1} \leftarrow$  ADAM( $\Delta_{\phi}^t, \beta$ )
    end for
end function

function CLIENT SIDE( $s, \phi, \mathbf{W}$ )
    Get local parameters  $\theta_s$ 
    for epoch  $e$  in  $1, \dots, E$  do
        for batch  $b \in B$  do
             $\phi'_c = p_{\mathbf{w}_z, \theta_s}(y_b = c, z|\mathbf{x}_b, s)^{1/(\beta N_{s,c})} / Z$ 
             $\phi_c \leftarrow \gamma \phi'_c + (1 - \gamma) \phi'_c \quad \triangleright$  Dampening
             $L_s \leftarrow \mathbb{E}_{q_{\phi}(z|y_b)}[\log p_{\mathbf{w}_z, \theta_s}(y_b, z|\mathbf{x}_b, s)]$ 
             $\mathbf{W} + = \alpha \nabla_{\mathbf{W}} L_s$ 
             $\theta_s + = \alpha \nabla_{\theta_s} L_s$ 
        end for
    end for
     $q(z|s) \leftarrow \mathbb{E}_{y \sim \mathcal{D}_s}[q_{\phi}(z|y)]$ 
    return  $\mathbf{w}_{1:K}, \phi, q(z|s)$ 
end function
    
```

---

several possible avenues for reducing communication costs. One such way is to “prune away” experts locally from the MoE if  $q_{\phi}(z|s)$  does not surpass a threshold  $\eta/K$ . Alternatively, since each expert is only modelling a subset of the data-space, their required modelling capacity and therefore parameter size can be reduced compared to the FedAvg model. Finally, sharing those parameters between experts that exhibit small gradient divergence or via a parameter efficient construction, such as with rank-1 factors (Dusenberry et al., 2020), are viable options to reduce communication. We explore pruning of experts in Appendix F and leave alternatives to future work.

**Designing robust gates** In the federated scenario,  $N_s$  is often much smaller than  $N$  and especially small in relation to the complexity of the data we try to model. Any local-

ized parameters therefore are prone to overfitting. On the other hand, the global parameters of an expert are trained using all data-points assigned to that expert across all shards, allowing to learn more robust features.

We can make use of the robustness of these experts’ features for the gating mechanism by conditioning on them instead of training an entirely separate model for  $p_{\theta_s}(z|\mathbf{x}, s)$ . Let us define  $h_k(\mathbf{x})$  as intermediary features of expert  $k$ . In order to scale with the number of experts, we introduce the local vector  $\pi_s \in \mathbb{R}_+^K, \sum_k \pi_{k,s} = 1$  with which the intermediate features are averaged before applying a linear transformation to compute the input to the softmax gates:

$$h_s(\mathbf{x}) = \sum_{k=1}^K \pi_{k,s} h_k(\mathbf{x})$$

$$p_{\theta_s}(z|\mathbf{x}, s) = \text{SM}(\mathbf{A}_s^T h_s(\mathbf{x}) + \mathbf{b}_s) \quad (10)$$

where  $\theta_s = (\pi_s, \mathbf{A}_s, \mathbf{b}_s)$  are local learnable parameters and SM represents the softmax function.

**Inference at test time** We consider three variants for test-time evaluation of FedMix. In the first case, a client  $s$  that participated in training is presented with a new data point  $(\mathbf{x}^*, s)$ . Predictions can then be straightforwardly done by selecting the  $y$  that maximizes  $\sum_{z=1}^K p(y|\mathbf{x}^*, z)p(z|\mathbf{x}^*, s)$ . In the second, more challenging, scenario a new client  $s^*$  is introduced together with a new labelled local data set  $\mathcal{D}_{s^*}$ . Here we propose to instantiate and train the local gating mechanism by optimizing the parameters  $\theta_s$  of  $p_{\theta_s}(z|\mathbf{x}, s^*)$  via the local objective. Afterwards, predictions can be made in a manner similar to the first case.

Finally, we consider the case in which a new client  $s^*$  has no labelled dataset available. Without a local gating function, simply ensembling experts exhibits almost random behaviour since experts can be overly confident on out-of-distribution data (Snoek et al., 2019). We therefore propose to ensemble across local gating mechanisms to compute  $p(z|\mathbf{x}^*) = \sum_{s=1}^S p_{\theta_s}(z|\mathbf{x}^*, s)p(s)$ ; a method which works well in practice. In Appendix H we discuss a more principled approach based on a complete graphical model perspective.

### 3. Related Work

FedMix has similarities to many recent works in the topic of federated learning. Two methods closely related to ours are described in (Sattler et al., 2019; Briggs et al., 2020). The authors propose to perform hierarchical clustering on the updates returned from each shard in order to incrementally create separate models for groups of users, with a cluster assignment mechanism based on handcrafted heuristics. FedMix instead takes a different approach; it starts with a fixed set of  $K$  models and then optimizes with gradi-

ent descent at each shard a per-datapoint model assignment mechanism that can better fit the peculiarities of the local dataset.

Another closely related work is presented by (Mansour et al., 2020), where the authors propose to similarly create an ensemble of  $K$ -models and assign to each shard the model that achieves the lowest training loss on the local dataset. This is closer to the assignment that happens in FedMix with one main difference; FedMix takes into account the uncertainty in the selection mechanism as well with  $p(z|x, s)$  instead of selecting the top performing component during training. This is beneficial early in training where the models have not fully specialized yet. Using local and global model parts has also been explored by (Liang et al., 2020). The authors propose to have a local feature extractor at each shard and a global classifier on top of those features as opposed to having  $K$ -separate models and a local selection mechanism as in FedMix. This setup yields improvements upon the vanilla federated averaging algorithm, however there are two potential drawbacks; first, empirically, the authors had to start their procedure from a pre-trained model with FedAvg and secondly, they have to ensemble all of the different feature extractors for predictions in new shards. In our experiments, we omit the pre-training step and show that the ensembling strategy fails.

Federated learning in the non-i.i.d setting can also be improved upon in other ways. (Li et al., 2018) propose to employ a proximal regularizer at the shard level in order to prevent the local models from drifting too far from the global model, thus making federated learning more robust. (Jiang et al., 2019) notice that FedAvg and Reptile (Nichol et al., 2018), a meta-learning algorithm, are essentially the same algorithm and thus propose fine tuning with Reptile in order to improve the personalized performance of the global model. In a similar vein, there are promising new works that explore the meta-learning view of federated learning (Chen et al., 2018; Khodak et al., 2019; Fallah et al., 2020). Improving the personalized performance of the global model has also been done without meta-learning in works such as by (Deng et al., 2020; Mansour et al., 2020). In general, such improvements are complementary to FedMix and can be used to further enhance its performance. We refer the interested readers to the recent surveys by (Kairouz et al., 2019; Kulkarni et al., 2020).

## 4. Experiments

We evaluate FedMix with  $K = 4$  experts across several datasets and non-i.i.d. settings. Results with different number of experts  $K$  can be found in Appendix E. For label-skew we use Cifar10, Cifar100 (Krizhevsky et al., 2009) and feminist (Caldas et al., 2018b), a 62-way image classification problem on hand-written digits and letters that is also nat-

urally non-i.i.d due to the different writing styles of 3500 users. In Appendix A we detail the experimental setup and provide additional ablation studies in Appendix E. For non-i.i.d.-ness in the input, we perform experiments on rotated MNIST (LeCun et al., 2010). For differences in the classification mechanism  $p(y|x, s)$  we explore label-permutation on Cifar10, following (Sattler et al., 2019), and show the strength of FedMix as a federated clustering algorithm.

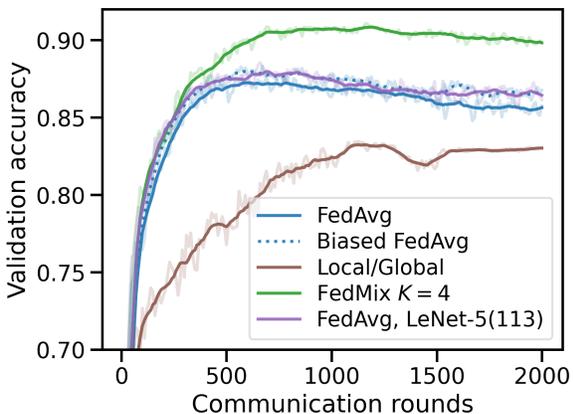


Figure 3. Average accuracy across all clients (y-axis) as a function of communication rounds for Cifar 10. Best viewed in color.

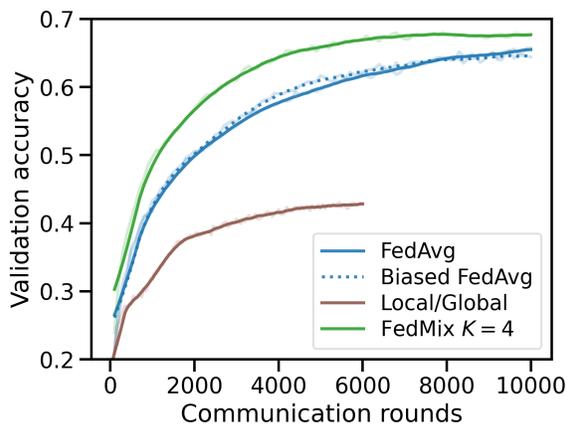


Figure 4. Average accuracy across all clients (y-axis) as a function of communication rounds for Cifar 100. Best viewed in color.

### 4.1. Label Skew

We compare FedMix along several dimensions to baselines such as (generalized) FedAvg (Reddi et al., 2020), *biased* FedAvg, and the Local/Global approach of (Liang et al., 2020). In *biased* FedAvg, we allow each client to learn a personalized bias vector  $b_s$  of its output layer. This will allow *biased* FedAvg to model the label skew at each client

Table 1. Average test-set accuracies across clients and communication costs (rounds and GB) for Cifar10, Cifar100 and Femnist at the end of training. We report both, the local accuracy after fine-tuning for one local epoch as well as the accuracy when evaluating at the server across the union of all local test-sets.

Method	Cifar 10			Cifar 100			Femnist		
	Local	Global	Comm.	Local	Global	Comm.	Local	Global	Comm.
FedAvg	85.98%	69.91%	2k, 137.3GB	65.67%	46.97%	10k, 205.53GB	90.93%	86.26%	5.2k, 235.71GB
biased FedAvg	86.82%	70.62%	2k, 130.92GB	64.41%	43.33%	10k, 205.52GB	90.72%	85.83%	5.2k, 235.71GB
Local/Global	83.13%	10.00%	2k, 129.93GB	42.82%	6.60%	6k, 116.73GB	85.54%	2.17%	5.2k, 235.61GB
FedMix K=4	<b>88.54%</b>	<b>76.42%</b>	2k, 522.5GB	<b>67.54%</b>	<b>47.48%</b>	10k, 822.12GB	<b>91.47%</b>	<b>87.00%</b>	5.2k, 942.84GB

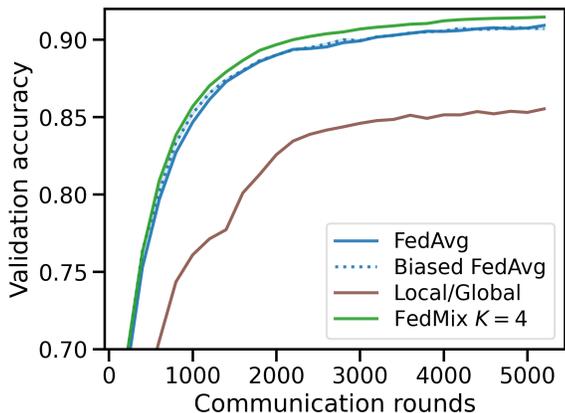


Figure 5. Average accuracy across all clients (y-axis) as a function of communication rounds for Femnist. Best viewed in color.

but, fundamentally, cannot model any other form of non-i.i.d.-ness. Similarly, (Liang et al., 2020) propose to split the model into local and global components by having local feature extractors and learning the upper layers of the neural network via FedAvg. We experimented with splitting LeNet-5 at every intermediate layer and report results with the best performing split: keeping the input layer local. For ResNet-20, splitting after the first block performed best. We show that training with FedMix achieves higher personalized model accuracy as well as global model accuracy on the server.

Figures 3, 4 and 5 shows learning curves for these different settings. These curves were obtained by averaging the local clients’ model accuracy on the validation set using the model that they last communicated to the server, thereby serving as a proxy for the local models’ performance. Table 1 shows that FedMix consistently outperforms the other approaches at the cost of more GBs communicated. Figure 3 shows, however, that for the same communication costs per communication round, FedMix offers better performance than FedAvg, as can be seen in comparing FedMix to FedAvg LeNet-5 with 113 channels. The channel count has been chosen such that the model size approximates 4

separate standard LeNet experts respectively. The advantage of *biased* FedAvg and Local/Global shows if we do not allow a local client to fine-tune the global model parameters. This setting becomes relevant if a clients’ local training data does not generalize to the test data, either because of a distribution shift or because it is too small. With finetuning we find that FedMix outperforms these approaches given that the local training data is sufficient. We show learning curves with the server-models in Appendix G.

For Femnist we observe a less dramatic improvement in performance compared to FedAvg since its source of non-i.i.d.-ness is not only expressed through label-skew.

#### 4.2. Rotated MNIST

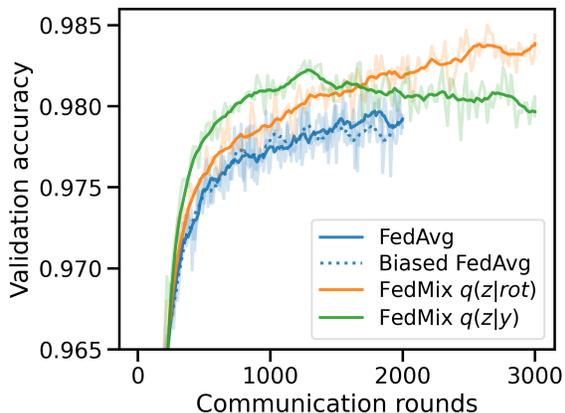


Figure 6. Only rotation

To show that FedMix is not limited to label-skew, we create a federated rotated MNIST dataset with 100 clients. Instead of label skew, each client randomly chooses a multiple of 45 degrees from a different probability distribution over 8 possible rotation angles to rotate a digit with. Each client’s distribution is drawn from Dir( $\alpha = 1.0$ ). At test time, each data-point is randomly rotated according to the client’s distribution. Additionally, we create a dataset where instead of uniform sampling of labels, we replicate the non-i.i.d.

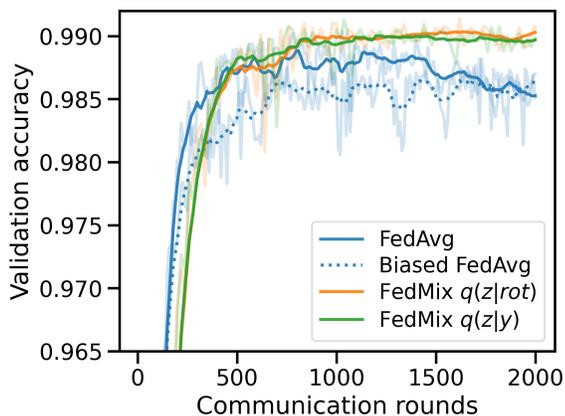


Figure 7. Rotation and labels

label skew described for Cifar10 above and combine it with the rotation non-i.i.d.-ness.

We compare FedMix where  $q$  is conditioned on  $y$  or on the degrees of rotation for a data-point against baselines. Figure 6 shows that FedMix benefits from being conditioned on the correct side-information for this task. Conditioning on  $y$  initially improves performance, however degrades in the end. In the presence of both sources of non-i.i.d.-ness, Figure 7 shows that FedMix improves performance regardless of which conditioning is chosen.

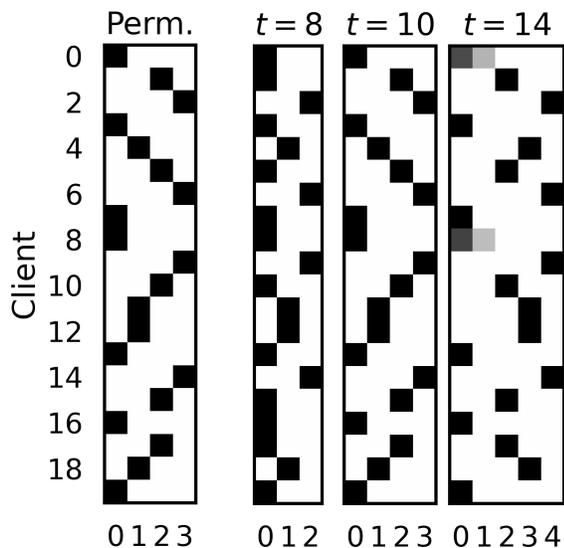


Figure 8. Ground truth and  $q(z|s)$  for  $K \in \{3, 4, 5\}$  after different number of communication rounds  $t$ . The ordering of columns is arbitrary. Greyscale represents probabilities (white: 0; black: 1)

### 4.3. Label Permutations

Apart from non-i.i.d.-ness in  $p(y)$  and  $p(x)$ , we can expect the mapping  $p(y|x)$  itself to be different between clients. We replicate the experimental setup of (Sattler et al., 2019) with 20 clients for Cifar10,  $C = 1.0$  and  $E = 3$ , albeit with LeNet-5. Each client is randomly assigned one of 4 different label permutations, determining the cluster assignment  $q(z|s)$  that FedMix has to learn. Although FedMix is designed to distinguish different regions of the input space, we show that it can perform user clustering. The gating function  $p(z|x, s)$  learns to correctly identify, for each data point, the expert corresponding to the permutation of  $s$ , thus recovering the original clustering; Figure 8 illustrates this effect (please note that the ordering of columns is arbitrary) for different number of experts  $K$ , using a dampening value of  $\gamma = 0.75$ .

When  $K$  matches the number of clusters in the data generating process, we see that FedMix correctly identifies the user clustering after just 10 communication rounds, which is much faster than the results presented in (Sattler et al., 2019). In the cases where  $K$  differs from the ground truth number of clusters we observe two phenomena, depending on whether we have fewer or more experts. When  $K$  is smaller, *i.e.*,  $K = 3$ , we observe that after 8 rounds, one expert takes responsibility for two of the four permutations whereas the other two experts correctly identify the remaining clusters. In the case of more experts than clusters, *i.e.*,  $K = 5$ , we observe that after 14 rounds FedMix splits the responsibility of one permutation across two experts and uses the other three experts to model each of the remaining permutations. Overall, we observe that FedMix has intuitive behaviour and, given enough capacity, can correctly identify the label permutations of each client.

## 5. Discussion

With FedMix we have introduced a federated learning algorithm that explicitly takes the non-i.i.d. characteristics of a federated dataset into account. We showed FedMix’ strength across a variety of non-i.i.d. datasets ranging from label-skew to input transformations and label permutations. Clients can learn to align specialized experts on sub-regions of the data space and achieve higher performance compared to FedAvg, in situations where the source of the non-i.i.d. nature is known. This assumption is very strong in real-world federated scenarios and we expect a more flexible alignment process than a global  $q$  to be the most interesting avenue for future research. In the future, we will explore ways to perform automatic selection of  $K$  as well as automatic selection of architecture elements to share between experts, trading-off gradient alignment and communication budgets.

## References

- Bonawitz, K., Eichner, H., Grieskamp, W., Huba, D., Ingerman, A., Ivanov, V., Kiddon, C., Konecny, J., Mazzocchi, S., McMahan, H. B., et al. Towards federated learning at scale: System design. *arXiv preprint arXiv:1902.01046*, 2019.
- Briggs, C., Fan, Z., and Andras, P. Federated learning with hierarchical clustering of local updates to improve training on non-iid data. *arXiv preprint arXiv:2004.11791*, 2020.
- Caldas, S., Konečný, J., McMahan, H. B., and Talwalkar, A. Expanding the reach of federated learning by reducing client resource requirements. *arXiv preprint arXiv:1812.07210*, 2018a.
- Caldas, S., Wu, P., Li, T., Konečný, J., McMahan, H. B., Smith, V., and Talwalkar, A. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097*, 2018b.
- Chen, F., Luo, M., Dong, Z., Li, Z., and He, X. Federated meta-learning with fast convergence and efficient communication. *arXiv preprint arXiv:1802.07876*, 2018.
- Deng, Y., Kamani, M. M., and Mahdavi, M. Adaptive personalized federated learning. *arXiv preprint arXiv:2003.13461*, 2020.
- Dusenberry, M., Jerfel, G., Wen, Y., Ma, Y., Snoek, J., Heller, K., Lakshminarayanan, B., and Tran, D. Efficient and scalable bayesian neural nets with rank-1 factors. In *International conference on machine learning*, pp. 2782–2792. PMLR, 2020.
- Fallah, A., Mokhtari, A., and Ozdaglar, A. Personalized federated learning: A meta-learning approach. *arXiv preprint arXiv:2002.07948*, 2020.
- Hsu, T.-M. H., Qi, H., and Brown, M. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., and Hinton, G. E. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.
- Jiang, Y., Konečný, J., Rush, K., and Kannan, S. Improving federated learning personalization via model agnostic meta learning. *arXiv preprint arXiv:1909.12488*, 2019.
- Jordan, M. I. and Jacobs, R. A. Hierarchical mixtures of experts and the em algorithm. *Neural computation*, 6(2): 181–214, 1994.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K., Charles, Z., Cormode, G., Cummings, R., et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- Khodak, M., Balcan, M.-F. F., and Talwalkar, A. S. Adaptive gradient-based meta-learning methods. In *Advances in Neural Information Processing Systems*, pp. 5915–5926, 2019.
- Kim, H., Kim, T., and Youn, C.-H. On federated learning of deep networks from non- $\{iid\}$  data: Parameter divergence and the effects of hyperparametric methods, 2020. URL <https://openreview.net/forum?id=SJeOAJStwB>.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Krizhevsky, A. et al. Learning multiple layers of features from tiny images. 2009.
- Kulkarni, V., Kulkarni, M., and Pant, A. Survey of personalization techniques for federated learning. *arXiv preprint arXiv:2003.08673*, 2020.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- LeCun, Y., Cortes, C., and Burges, C. Mnist handwritten digit database. 2010.
- Li, T., Sahu, A. K., Zaheer, M., Sanjabi, M., Talwalkar, A., and Smith, V. Federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*, 2018.
- Li, T., Sanjabi, M., and Smith, V. Fair resource allocation in federated learning. *arXiv preprint arXiv:1905.10497*, 2019.
- Liang, P. P., Liu, T., Ziyin, L., Salakhutdinov, R., and Morency, L.-P. Think locally, act globally: Federated learning with local and global representations. *arXiv preprint arXiv:2001.01523*, 2020.
- Mansour, Y., Mohri, M., Ro, J., and Suresh, A. T. Three approaches for personalization with applications to federated learning. *arXiv preprint arXiv:2002.10619*, 2020.
- McMahan, H. B., Moore, E., Ramage, D., Hampson, S., et al. Communication-efficient learning of deep networks from decentralized data. *arXiv preprint arXiv:1602.05629*, 2016.
- Mohri, M., Sivek, G., and Suresh, A. T. Agnostic federated learning. *arXiv preprint arXiv:1902.00146*, 2019.

- Nalisnick, E., Matsukawa, A., Teh, Y. W., Gorur, D., and Lakshminarayanan, B. Do deep generative models know what they don't know? *arXiv preprint arXiv:1810.09136*, 2018.
- Nichol, A., Achiam, J., and Schulman, J. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018.
- Reddi, S., Charles, Z., Zaheer, M., Garrett, Z., Rush, K., Konečný, J., Kumar, S., and McMahan, H. B. Adaptive federated optimization. *arXiv preprint arXiv:2003.00295*, 2020.
- Sattler, F., Müller, K.-R., and Samek, W. Clustered federated learning: Model-agnostic distributed multi-task optimization under privacy constraints. *arXiv preprint arXiv:1910.01991*, 2019.
- Snoek, J., Ovadia, Y., Fertig, E., Lakshminarayanan, B., Nowozin, S., Sculley, D., Dillon, J., Ren, J., and Nado, Z. Can you trust your model's uncertainty? evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems*, pp. 13969–13980, 2019.
- Wu, Y. and He, K. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 3–19, 2018.

## A. Experimental Setup

For all experiments in the experimental section, we use a SGD optimizer with a learning rate of 0.05 locally and the Adam (Kingma & Ba, 2014) optimizer with its default hyperparameters at the server by interpreting the difference of the local from the global model as a gradient (Reddi et al., 2020). For `FedMix`, the features  $h_k(\mathbf{x})$  are defined as the input to a expert  $k$ 's output layer. Unless otherwise mentioned, we choose  $\gamma = 0.99$  for the dampening value of the  $\phi_{c,s}$  updates locally.

### A.1. Label Skew

**Cifar10** For Cifar10, we replicate the federated data split of (Hsu et al., 2019). The dataset is split across 100 clients, whose data-points are drawn according to their label from a  $\text{Dir}(\alpha = 1.0)$  distribution without replacement. For the base model, we use a LeNet-5 architecture (LeCun et al., 1998). We use a batch size of  $B = 64$  locally. We sample 10 clients without replacement on each round (but with replacement across rounds) and train for  $E = 1$  local epochs. We found dropout to be necessary to avoid overfitting to the local training data-sets. For all experiments with the LeNet-5 architecture on Cifar10 we found a dropout rate of 0.3 for the second convolutional layer and a rate of 0.1 for the first fully connected layer to be helpful, except for the Local/Global (Liang et al., 2020) experiments, where we consequently omit dropout. We found a value of  $\beta = 0.8$  to be best in balancing between robust experts and fast specialization.

**Cifar100** For Cifar100 we replicate the data split of (Reddi et al., 2020). The dataset is split into 500 clients by using a hierarchical model over the coarse and fine labels, with the same hyperparameters as the ones provided by (Reddi et al., 2020). The other hyperparameters are the same as Cifar10 with the exception of the batch size, where we use  $B = 20$ , as well as the architecture, where we use a ResNet-20 with group normalization (Wu & He, 2018) layers instead of batch normalization (Ioffe & Szegedy, 2015). We augment the data by random cropping from a 4 pixel padded image and horizontal flipping. We found  $\beta = 1.0$  to perform well for Cifar100, resulting in reliable specialisation.

**Femnist** Finally, for the Femnist dataset, we similarly followed the setup of (Reddi et al., 2020) with the same LeNet-5 architecture and hyperparameters of Cifar10 with the exception of the batch size where we used  $B = 20$ . We found no danger of overfitting, so we omit dropout here. Similarly to Cifar10, we found a value of  $\beta = 0.8$  to serve as a good tradeoff.

### A.2. Rotated Mnist

For the experiments on rotated Mnist, we use the LeNet-5 architecture (LeCun et al., 1998) and keep the same hyperparameters as for Cifar10, *i.e.* a local batch size of  $B = 64$ , 10 clients per round and a single local epoch  $E = 1$  but do not use dropout. When conditioning on the label information *i.e.* using  $q(z|y)$ , we choose  $\beta = 0.8$  and  $\gamma = 0.99$ . When using  $q(z|rot)$ , we found such a high dampening factor does not lead to any specialization. Instead, lowering it to  $\gamma = 0.5$  achieved specialized experts.

### A.3. Label Permutations

For the label permutation experiments, we again make use of LeNet-5 on the Cifar10 dataset. Replicating the setup of (Sattler et al., 2019), we deviate from our previous hyperparameters and chose  $S = 20$  clients, no client subsampling (*i.e.*  $C = 1.0$ ) and three local epochs  $E = 3$  per client. We omit dropout since we do not train these models to convergence. We found `FedMix` to be very robust in finding the correct cluster assignment in a small amount of steps. A value of  $\gamma = 0.75$  quickly causes specialization across all values of  $K$  we experiment with.

## B. Solving for $\phi_c$

Here we present the derivation for  $\phi$  discussed in the main text. We aim to maximize Eq. 5 on a specific client  $s$  for the parameters  $\phi$  of  $q_\phi(z|y)$ , *i.e.* solve  $\arg \max_\phi L_s(\phi)$  with

$$L_s(\phi) = \sum_{i=1}^{N_s} \mathbb{E}_{q_\phi(z|y_{s,i})} [\log p_{\mathbf{w}_z}(y_{s,i}|\mathbf{x}_{s,i}, z) p_{\theta_s}(z|\mathbf{x}_{s,i}, s)] + \beta H(q_\phi(z|y_{s,i})). \quad (11)$$

We assume a parameterization in the form of  $q(z = k|y = c) = \phi_{c,k}$ ,  $\sum_{z=1}^K \phi_{c,k} = 1$ . We can extend out maximization target to

$$L_s(\phi) = \sum_{c=1}^C \sum_{i=1}^{N_{s,c}} \mathbb{E}_{q_\phi(z|y_{s,i}=c)} [\log p_{\mathbf{w}_z}(y_{s,i} = c|\mathbf{x}_{s,i}, z) p_{\theta_s}(z|\mathbf{x}_{s,i}, s)] + \beta H(q_\phi(z|y_{s,i})) \quad (12)$$

$$= \sum_{c=1}^C \sum_{i=1}^{N_{s,c}} \left( \sum_{k=1}^K \phi_{c,k} [\log p_{\mathbf{w}_z}(y_{s,i} = c|\mathbf{x}_{s,i}, z) p_{\theta_s}(z|\mathbf{x}_{s,i}, s)] - \beta \sum_{k=1}^K \phi_{c,k} \log \phi_{c,k} \right), \quad (13)$$

where  $C$  corresponds to the number of classes,  $N_{s,c}$  the number of data-points on shard  $s$  belonging to class  $c$  and in the second line we have substituted  $q$  for its direct parameterization.

The requirement of  $\sum_{z=1}^K \phi_{c,k} = 1, \forall c$  transforms the optimization problem into a constrained optimization problem that we can approach by introducing lagrangian multipliers  $\lambda_c$  and optimizing the Lagrangian, *i.e.*  $\arg \max_{\phi, \lambda} \tilde{L}_s(\phi, \lambda)$ :

$$\tilde{L}_s(\phi, \lambda) = \sum_{c=1}^C \sum_{i=1}^{N_{s,c}} \left( \sum_{k=1}^K \phi_{c,k} [\log p_{\mathbf{w}_z}(y_{s,i} = c|\mathbf{x}_{s,i}, z) p_{\theta_s}(z|\mathbf{x}_{s,i}, s)] - \beta \sum_{k=1}^K \phi_{c,k} \log \phi_{c,k} \right) + \sum_{c=1}^C \lambda_c \sum_{k=1}^K (\phi_{c,k} - 1). \quad (14)$$

The problem can be decomposed into  $C$  independent problems  $\arg \max_{\phi_c, \lambda_c} \tilde{L}_{s,c}(\phi_c, \lambda_c)$ , solving for  $\phi_c$  and  $\lambda_c$  each. Setting the gradient with respect to the lagrangian multiplier to zero recovers our constraint:

$$\frac{\partial}{\partial \lambda_c} \tilde{L}_{s,c}(\phi_c, \lambda_c) = 0 \quad (15)$$

$$\sum_{z=1}^K \phi_{c,k} = 1 \quad (16)$$

We compute the gradient  $\frac{\partial}{\partial \phi_{c,k}} \tilde{L}_{s,c}(\phi_c, \lambda_c)$  as follows, set it to zero and solve for  $\phi_{c,k}$ :

$$\frac{\partial}{\partial \phi_{c,k}} \tilde{L}_{s,c}(\phi_c, \lambda_c) = \sum_{i=1}^{N_{s,c}} [\log p_{\mathbf{w}_z}(y_{s,i} = c|\mathbf{x}_{s,i}, z) p_{\theta_s}(z|\mathbf{x}_{s,i}, s) - \beta (\log \phi_{c,k} + 1)] + \lambda_c = 0 \quad (17)$$

$$= \sum_{i=1}^{N_{s,c}} [\log p_{\mathbf{w}_z}(y_{s,i} = c|\mathbf{x}_{s,i}, z) p_{\theta_s}(z|\mathbf{x}_{s,i}, s)] - \beta N_{s,c} \log \phi_{c,k} - \beta N_{s,c} + \lambda_c = 0 \quad (18)$$

$$\log \phi_{c,k} = \frac{1}{\beta N_{s,c}} \left( \sum_{i=1}^{N_{s,c}} [\log p_{\mathbf{w}_z}(y_{s,i} = c|\mathbf{x}_{s,i}, z) p_{\theta_s}(z|\mathbf{x}_{s,i}, s)] - \beta N_{s,c} + \lambda_c \right) \quad (19)$$

$$\phi_{c,k} = \exp \left( \frac{\lambda_c}{\beta N_{s,c}} - 1 \right) \left( \prod_{i=1}^{N_{s,c}} p_{\mathbf{w}_z}(y_{s,i} = c|\mathbf{x}_{s,i}, z) p_{\theta_s}(z|\mathbf{x}_{s,i}, s) \right)^{\frac{1}{\beta N_{s,c}}} \quad (20)$$

Combining the last result for  $\phi_{c,k}$  with Eq. 16, we can proceed to solve for  $\lambda_c$ :

$$\sum_{z=1}^K \phi_{c,k} = 1 = \exp\left(\frac{\lambda_c}{\beta N_{s,c}} - 1\right) \sum_{z=1}^K \left( \prod_{i=1}^{N_{s,c}} p_{\mathbf{w}_z}(y_{s,i} = c | \mathbf{x}_{s,i}, z) p_{\theta_s}(z | \mathbf{x}_{s,i}, s) \right)^{\frac{1}{\beta N_{s,c}}} \quad (21)$$

$$\exp\left(1 - \frac{\lambda_c}{\beta N_{s,c}}\right) = \sum_{z=1}^K \left( \prod_{i=1}^{N_{s,c}} p_{\mathbf{w}_z}(y_{s,i} = c | \mathbf{x}_{s,i}, z) p_{\theta_s}(z | \mathbf{x}_{s,i}, s) \right)^{\frac{1}{\beta N_{s,c}}} \quad (22)$$

$$1 - \frac{\lambda_c}{\beta N_{s,c}} = \log \left( \sum_{z=1}^K \left( \prod_{i=1}^{N_{s,c}} p_{\mathbf{w}_z}(y_{s,i} = c | \mathbf{x}_{s,i}, z) p_{\theta_s}(z | \mathbf{x}_{s,i}, s) \right)^{\frac{1}{\beta N_{s,c}}} \right) \quad (23)$$

$$\lambda_c = \beta N_{s,c} \left( 1 - \log \left( \sum_{z=1}^K \left( \prod_{i=1}^{N_{s,c}} p_{\mathbf{w}_z}(y_{s,i} = c | \mathbf{x}_{s,i}, z) p_{\theta_s}(z | \mathbf{x}_{s,i}, s) \right)^{\frac{1}{\beta N_{s,c}}} \right) \right) \quad (24)$$

Finally, we can integrate our solution for  $\lambda_c$  into Eq. 20. Note that the exponential expression containing  $\lambda_c$  simplifies by cancelling out  $\beta N_{s,c}$  and the factor 1:

$$\phi_{c,k} = \exp \left( - \log \left( \sum_{z=1}^K \left( \prod_{i=1}^{N_{s,c}} p_{\mathbf{w}_z}(y_{s,i} = c | \mathbf{x}_{s,i}, z) p_{\theta_s}(z | \mathbf{x}_{s,i}, s) \right)^{\frac{1}{\beta N_{s,c}}} \right) \right) \cdot \left( \prod_{i=1}^{N_{s,c}} p_{\mathbf{w}_z}(y_{s,i} = c | \mathbf{x}_{s,i}, z) p_{\theta_s}(z | \mathbf{x}_{s,i}, s) \right)^{\frac{1}{\beta N_{s,c}}} \quad (25)$$

$$= \frac{\left( \prod_{i=1}^{N_{s,c}} p_{\mathbf{w}_z}(y_{s,i} = c | \mathbf{x}_{s,i}, z) p_{\theta_s}(z | \mathbf{x}_{s,i}, s) \right)^{\frac{1}{\beta N_{s,c}}}}{\left( \sum_{z=1}^K \left( \prod_{i=1}^{N_{s,c}} p_{\mathbf{w}_z}(y_{s,i} = c | \mathbf{x}_{s,i}, z) p_{\theta_s}(z | \mathbf{x}_{s,i}, s) \right)^{\frac{1}{\beta N_{s,c}}} \right)} \quad (26)$$

$$= \frac{\exp\left(\frac{1}{\beta N_{s,c}} \sum_{i=1}^{N_{s,c}} \log p_{\mathbf{w}_z}(y_{s,i} = c | \mathbf{x}_{s,i}, z) p_{\theta_s}(z | \mathbf{x}_{s,i}, s)\right)}{\sum_{z=1}^K \exp\left(\frac{1}{\beta N_{s,c}} \sum_{i=1}^{N_{s,c}} \log p_{\mathbf{w}_z}(y_{s,i} = c | \mathbf{x}_{s,i}, z) p_{\theta_s}(z | \mathbf{x}_{s,i}, s)\right)}, \quad (27)$$

where in the second step we have taken the exponent of the log of the product over  $N_{s,c}$  data points, both, in the numerator and denominator. It is now easy to see how the solution for  $\phi_{c,k}$  corresponds to a temperature  $\beta$  modulated softmax over the log-likelihood of datapoints belonging to class  $c$ . In practice, a client  $s$  performs mini-batch stochastic gradient descent on the parameters  $\mathbf{w}_z, \theta_s$ . We therefore want to avoid evaluating all  $N_s$  data-points in order to compute  $\phi_{c,k}$ . Even though we only require the forward pass through our model, therefore not causing any memory issues, we want to avoid the additional computational overhead. Similarly to mini-batch stochastic gradient descent, we therefore approximate the costly average over  $N_{s,c}$  datapoints in Eq. 27 by considering only a mini-batch of  $M_s$  data-points, where  $M_{s,c}$  denotes all data-points in the mini-batch  $M_s$  that belong to class  $c$ :

$$\phi'_{s,c} = \frac{\exp\left(\frac{1}{\beta M_{s,c}} \sum_{i=1}^{M_{s,c}} \log p_{\mathbf{w}_z}(y_{s,i} = c | \mathbf{x}_{s,i}, z) p_{\theta_s}(z | \mathbf{x}_{s,i}, s)\right)}{\sum_{z=1}^K \exp\left(\frac{1}{\beta M_{s,c}} \sum_{i=1}^{M_{s,c}} \log p_{\mathbf{w}_z}(y_{s,i} = c | \mathbf{x}_{s,i}, z) p_{\theta_s}(z | \mathbf{x}_{s,i}, s)\right)}. \quad (28)$$

In order to stabilize the update of  $\phi_{s,c}$  from round  $t$  to round  $t + 1$  during the local client's training epoch, we introduce a dampening factor  $\gamma$ ,

$$\phi_{s,c}^{t+1} = \gamma \phi_{s,c}^t + (1 - \gamma) \phi_{s,c}^{t+1}, \quad (29)$$

thereby completing the update rule for  $\phi_{s,c}$ .

## C. Privacy Implications

Privacy is one of the key motivations for research and deployment of Federated Learning. Even though privacy is not a focus of this paper, we briefly discuss some implications of making explicit use of  $q(y|s)$  in `FedMix` in the main text. Here, we want to shine light on the practical possibility to reconstruct  $p(y|s)$  at the server-side in standard `FedAvg`, arguing that `FedMix` does not reveal information that is not already accessible by the server.

Assume a randomly initialized model being sent to a client  $s$ , where the client performs a single full batch update step on the output layer’s bias vector  $b_s$ . Assuming a softmax cross-entropy loss  $L_s$ , the average gradient with respect to a the  $k$ -th entry  $b_k$  takes the form of

$$\frac{\partial L_s}{\partial b_k} = \frac{1}{N_s} \sum_{i=1}^{N_s} \mathbb{1}[y_{i,k} = y_{i,\text{true}}] - \pi_{i,k}, \quad (30)$$

where  $\mathbb{1}$  is the indicator function and  $\pi_{i,k}$  is the softmax probability of class  $k$  of datapoint  $i$ . With a randomly initialized model, these softmax probabilities can be assumed to be uniform, leading to an average gradient of

$$\frac{\partial L_s}{\partial b_k} = \frac{N_{s,k}}{N_s} - \frac{1}{N_c} = p(y|s) - \frac{1}{N_c}, \quad (31)$$

where  $N_c$  is the number of classes. Upon sending the updated bias vector  $b_s = b - \alpha \frac{\partial L_s}{\partial b_k}$  to the server, it can easily reconstruct the marginal label distribution.

Figure 9 shows, for every client in our Cifar10 setup (Appendix A), in red the true marginal  $p(y|s)$  and in blue the reconstructed marginal based on (the same) randomly initialized model being sent to each client. Clearly, we have high congruence. In Figure 10, we investigate the more realistic setup of  $E = 1$  with mini-batch stochastic gradient descent at the client level. In order to avoid reconstructing the multi-step update, we simply normalize the difference  $(b - b_s)$  and interpret it as marginal label distribution. We see that multiple updates (on average: 8) reduce the congruence between the true and reconstructed marginal, however the information leakage is still remarkable.

## D. Gradient Divergence

We aim to track the divergence of updates for a subset  $S'$  of shards at the server at time step  $t$  for `FedMix` and `FedAvg`. Therefore we define a metric inspired by (Kim et al., 2020; Sattler et al., 2019) to define divergence of gradients  $\Delta_{k,i}^t = \omega_k^t - \omega_{k,i}^{t+1}$  for some subset  $\omega_k$  of the parameters  $\mathbf{w}_k$  of expert  $k$  as

$$\text{GD}(\Delta_k^t) = \sum_{i=1}^{S'} \sum_{j=1}^{S'} p(s=i|z=k)p(s=j|z=k) \cdot 0.5 \cdot \left( 1 - \frac{\Delta_{k,i}^t \cdot \Delta_{k,j}^t}{\|\Delta_{k,j}^t\| \cdot \|\Delta_{k,i}^t\|} \right). \quad (32)$$

For `FedAvg`, the above metric collapses to

$$\text{GD}(\Delta^t) = \sum_{i=1}^{S'} \sum_{j=1}^{S'} p(s=i)p(s=j) \cdot 0.5 \cdot \left( 1 - \frac{\Delta_i^t \cdot \Delta_j^t}{\|\Delta_j^t\| \cdot \|\Delta_i^t\|} \right). \quad (33)$$

In Figure 1 in the main text, we plot the sum of  $\text{GD}(\Delta_k^t)$  across all parameters  $\omega_k$  (*i.e.*, convolutional kernels, weights and biases) of the LeNet-5 experts in comparison to  $\omega$  for `FedAvg`.

## E. Ablation studies

We investigate the characteristics of `FedMix` by varying the number of experts  $K$  for Cifar10 and Cifar100. Figure 11 shows learning curves of several values of  $K$  for these two datasets. We can see that a higher number of experts consistently improves accuracy however with a diminishing rate of return. With increasing  $K$ , the modeling task for a single expert becomes progressively easier and data across clients becomes more aligned at the cost of higher communication and computation.



Figure 9. Histogram representation of  $p(y|s)$  as well as its server-side reconstruction for every client  $s$  for the Cifar10 setup described in Appendix A. The x-axis per sub-plot enumerates the 10 classes. For every class  $c$ , the left bar in red represents  $p(y = c|s)$  and the right bar in blue represents its reconstruction. Each client performed a single full data-set update step.

## F. Reducing Communication Costs

We can ignore the communication of experts between server and client if, for this given client, that particular expert will not be updated during training. If  $q_\phi(z = k|s) = E_{p(y|s)}[q_\phi(z = k|y)] = 0$  for an expert  $k$  on client  $s$ , then parameters of that expert observe no gradient and the expert is effectively pruned from the local library of available experts. Since the parameterization of  $q_\phi(z|s)$  does not allow for exact values of zero, we introduce  $\eta$ , such that if  $q_\phi(z = k|s) < \eta/K$ , we consider expert  $k$  to be pruned. Since the value of  $q_\phi(z = k|s)$  during local optimization at the client is subject to change, it is possible that  $q_\phi(z = k|s)$  briefly lies above the threshold even though it had just been pruned away. Therefore the server considers expert  $k$  pruned for a given client  $s'$  only if  $q_\phi(z = k|s = s') < 0.9 \cdot \eta/K$ . Algorithm 2 details this approach. Note that we write  $p(z|s) = \mathbb{E}_{\mathbf{x} \sim D_s}[p_{\theta_s}(z|\mathbf{x}, s)]$ , i.e. the true marginal, however in practice we make use of the cheap to compute marginal approximation  $q_\phi(z|s) = \mathbb{E}_{y \sim D_s}[q_\phi(z|y)]$ . Figure 12 shows experiments with different values of the threshold  $\eta$ . All models were trained for the same number of communication rounds ( $2k$  for Cifar10 and  $10k$  for Cifar100). As training progresses and clients begin to drop experts, the amount of GBs communicated is reduced, leading to less GB communicated in total. At the same time, pruning does impact performance significantly, leaving room for improvement in future work. Instead of dropping experts, it might be beneficial to prune weights and features of the individual experts themselves. Inspired by Federated Dropout (Caldas et al., 2018a), it is interesting to study the effect of dropping experts with a dropout probability related to  $q(z|s)$  between communication rounds.



Figure 10. Histogram representation of  $p(y|s)$  as well as its server-side reconstruction for every client  $s$  for the Cifar10 setup described in Appendix A. The x-axis per sub-plot enumerates the 10 classes. For every class  $c$ , the left bar in red represents  $p(y = c|s)$  and the right bar in blue represents its reconstruction. Each client performed multiple mini-batch update steps (on average 8).

### G. Non-Personalized Evaluation

In the main text we discussed how FedMix creates experts that serve as better initialization points for local finetuning. In Figure 13 we show learning curves using the server-side model parameters without any local finetuning. We see how FedMix not only serves as a better model for finetuning, but also how the combination of server-side model with the local gating mechanism positions FedMix above FedAvg. The server-side model in combination with the locally trained output bias of *biased* FedAvg shows the importance of the output bias for combating label skew, however comparing to the results from the main text it seems that standard FedAvg can recover the same advantage through finetuning. Local/Global equally shows the strength of combining local feature extractors with the server-side model weights, however cannot recover additional gains through fine-tuning the server-side models.

### H. New Shard Inference

In the main text, we discussed performance of the individual algorithms when evaluating new data on a shard that took part in training and therefore has access to trained local model parameters. This is the case for all considered methods except FedAvg, for which there exist no local parameters. Consequently, FedAvg can easily be applied to new clients without any local training data. Methods with localized parameters, i.e. *biased* FedAvg, Local/Global and FedMix require special consideration.

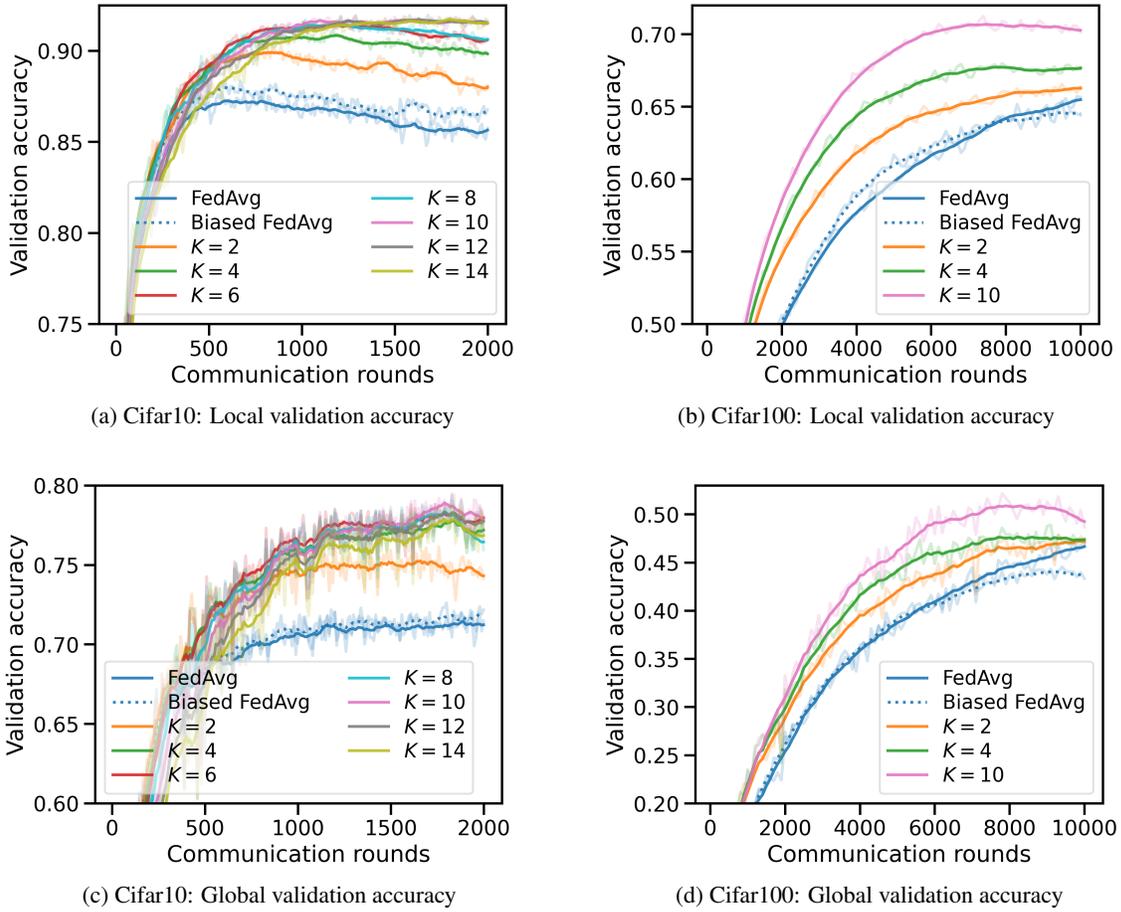


Figure 11. Ablation studies on the effect of  $K$  on the average local accuracy (a, b) and new shard accuracy (c, d) on Cifar10 and Cifar100.

(Liang et al., 2020) propose to ensemble the representations of the local feature extractors across clients before evaluating the global part of the network. We find this approach to work quite poorly in practice. In *biased FedAvg*, we ensemble the individual local biases across clients to receive a single global bias. For *FedMix*, we marginalize the local gating predictions to achieve a global gating prediction  $p(z|\mathbf{x}^*) = \sum_{s=1}^S p_{\theta_s}(z|\mathbf{x}^*, s)p(s)$ . Predictions can then be made by marginalizing across experts using this global gating function:  $p(y|\mathbf{x}^*) = \sum_{z=1}^K p_{w_k}(y|\mathbf{x}^*, z)p(z|\mathbf{x}^*)$ . We realize that this is not a scalable solution to the cross-device FL setting, since evaluating  $S$  gating functions is costly, however it is feasible in the cross-silo setting. We leave techniques for model distillation for future work.

In Figure 14 we use the sum of all local validation sets as a proxy for inference on a new client. Local/Global is left out of the Cifar10 plot since it has random performance throughout. *FedMix* displays very good performance compared to (*biased*) *FedAvg*, however *FedAvg* eventually catches up with *FedMix* on Cifar100. In Figure 11 we see that the global validation accuracy scales with the number of experts  $K$ .

### H.1. Using a Generative Model

As an alternative to marginalizing the local gate predictions using  $p(s)$ , investigating the graphical model in Figure 15 reveals the possibility for marginalization with  $p(s|\mathbf{x}^*)$ :

$$p(z|\mathbf{x}^*) = \sum_{s=1}^S p_{\theta_s}(z|\mathbf{x}^*, s)p(s|\mathbf{x}^*), \quad p(s|\mathbf{x}^*) \propto p(\mathbf{x}^*|s)p(s). \quad (34)$$

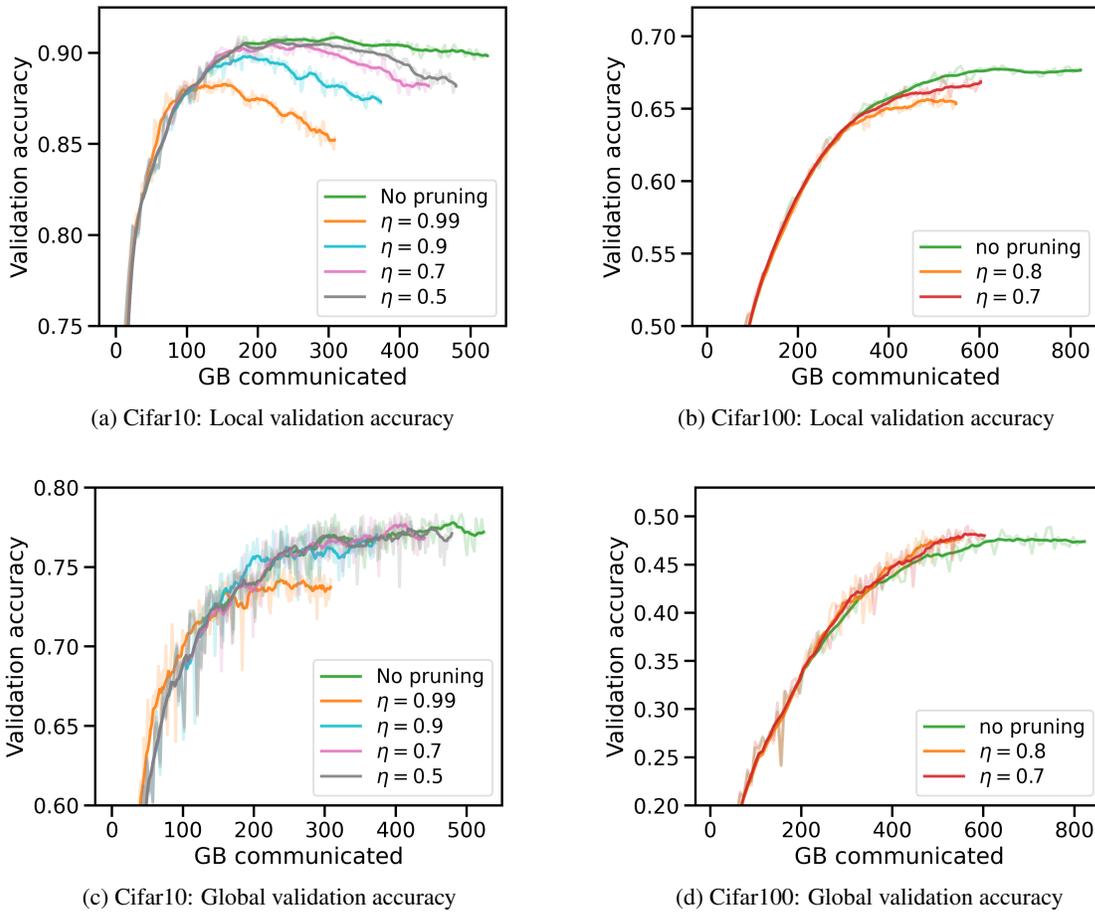


Figure 12. The effect of pruning experts on the average local accuracy (a, b) and new shard accuracy (c,d) on Cifar10 and Cifar100.

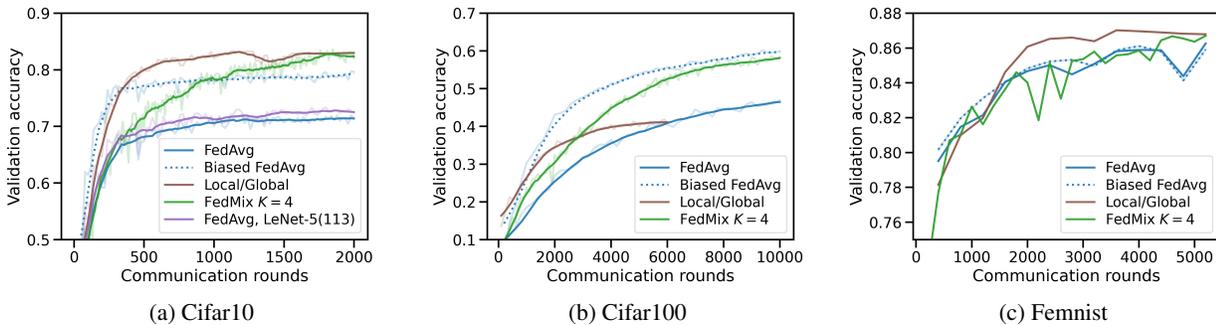


Figure 13. Accuracy using the most recent server weights (y-axis) as a function of communication rounds. Cifar10 models are trained on the standard 45k training split. Best viewed in color.

It is for this third evaluation case that training local generative models  $p(\mathbf{x}|s)$  for each client becomes interesting, as they allow to compute the responsibilities  $p(s|\mathbf{x}^*)$  at test time. In practice, however, the success of this approach depends heavily on the correctness of  $p(s|\mathbf{x}^*)$ , which in turn depends on the ability of the local generative models  $p(\mathbf{x}|s)$  to assign high probability to data that resembles  $\mathcal{D}_s$  and low probability to out-of-distribution data. Training and calibrating generative models for this task is in itself an active area of research (Nalisnick et al., 2018) and investigating how clients in a federated

**Algorithm 2** The FedMix algorithm.  $\alpha, \beta$  are the client and server learning rates;  $\gamma$  is a dampening factor and  $Z$  a normalization constant.  $\eta \in [0, 1]$  is the pruning threshold.

**function** SERVER SIDE

Initialize  $\phi$  and  $K$  vectors  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_K]$   
 Initialize  $p(z|s) = 1/K \forall s$   
**for** round  $t$  in  $1, \dots, T$  **do**  
      $S' \leftarrow$  random subset of the clients  
     Initialize  $\Delta_{\mathbf{W}}^t = \mathbf{0}, \Delta_{\phi}^t = \mathbf{0}$   
     **for**  $s$  in  $S'$  **do**  
          $\mathbf{W}' \leftarrow [\mathbf{w}_k \mid p(z = k|s) \geq 0.9 \cdot \eta/K]$   
          $\mathbf{W}_s^t, \phi_s^t, p(z|s) \leftarrow$  CLIENT SIDE( $s, \phi, \mathbf{W}'$ )  
         Store  $p(z|s)$   
     **end for**  
      $p(s|z) \leftarrow p(z|s)p(s) / \sum_{s \in S'} p(z|s)p(s)$   
     **for**  $s$  in  $S'$  **do**  
          $\Delta_{\mathbf{w}_k}^t + = p(s|z = k)(\mathbf{w}_k^{t-1} - \mathbf{w}_{s,k}^t) \quad \forall k$   
          $\Delta_{\phi}^t + = \frac{N_s}{N_{S'}}(\phi^{t-1} - \phi_s^t)$   
     **end for**  
      $\Delta_{\phi}^t - = \nabla_{\phi} H(\sum_c q_{\phi}(z|y=c)p(y=c))$   
      $\mathbf{w}_{1:K}^{t+1} \leftarrow$  ADAM( $\Delta_{\mathbf{w}_{1:K}}^t, \beta$ )  
      $\phi^{t+1} \leftarrow$  ADAM( $\Delta_{\phi}^t, \beta$ )  
**end for**  
**end function**

**function** CLIENT SIDE( $s, \phi, \mathbf{W}$ )

Get local parameters  $\theta_s$   
**for** epoch  $e$  in  $1, \dots, E$  **do**  
     **for** batch  $b \in B$  **do**  
          $q(z|s) \leftarrow \mathbb{E}_{y \sim \mathcal{D}_s} [q_{\phi}(z|y)]$   
          $K' \leftarrow [k \mid q(z = k|s) \geq \eta/K]$  ▷ Indices of remaining experts  
          $\bar{\phi}_c = p_{\mathbf{w}_z, \theta_s}(y_b = c, z | \mathbf{x}_b, s)^{1/(\beta N_{s,c})} / Z$  ▷ Update for remaining experts  
          $\phi'_c \leftarrow \bar{\phi}_c \cdot \sum_{k \in K'} \phi_{c,k}$  ▷ Reweigh to account for pruned experts  
          $\phi_c \leftarrow \gamma \phi_c + (1 - \gamma) \phi'_c$  ▷ Dampening and update of  $\phi_c$  for the remaining experts  
          $\phi''_{c,k} \leftarrow \frac{\phi_{c,k}}{\sum_{k \in K'} \phi_{c,k}}, k \in K'$  ▷ Renormalize updated  $\phi$  for remaining experts  
          $L_s \leftarrow \mathbb{E}_{q_{\phi''}(z|y_b)} [\log p_{\mathbf{w}_z}(y_b | \mathbf{x}_b, z)] + \mathbb{E}_{q_{\phi}(z|y_b)} [\log p_{\theta_s}(z | \mathbf{x}_b, s)]$   
          $\mathbf{W} + = \alpha \nabla_{\mathbf{W}} L_s$   
          $\theta_s + = \alpha \nabla_{\theta_s} L_s$   
     **end for**  
     **end for**  
      $q(z|s) \leftarrow \mathbb{E}_{y \sim \mathcal{D}_s} [q_{\phi}(z|y)]$   
      $\mathbf{W}' \leftarrow [\mathbf{w}_k \mid q(z = k|s) \geq \eta/K]$   
     **return**  $\mathbf{W}'; \phi; q(z|s)$   
**end function**

setting might exchange information to facilitate this process is not yet explored. We therefore leave a thorough evaluation of this case to future work and only present here a MNIST (LeCun et al., 2010) experiment.

We train FedMix with  $K = 4$  and experts of two hidden layer ReLU MLP with 200 hidden units on MNIST. We split the dataset into  $S = 100$  clients according to the procedure described in (Liang et al., 2020). FedMix achieves 97.7% average validation accuracy compared to 97.0% with FedAvg after 600 communication steps. Independently for each client, we train a small variational autoencoder with a 32-dimensional latent space using a two-layer MLP with 512 and 256 hidden units respectively as encoder and mirrored decoder structure. We optimize the VAEs using Adam with standard

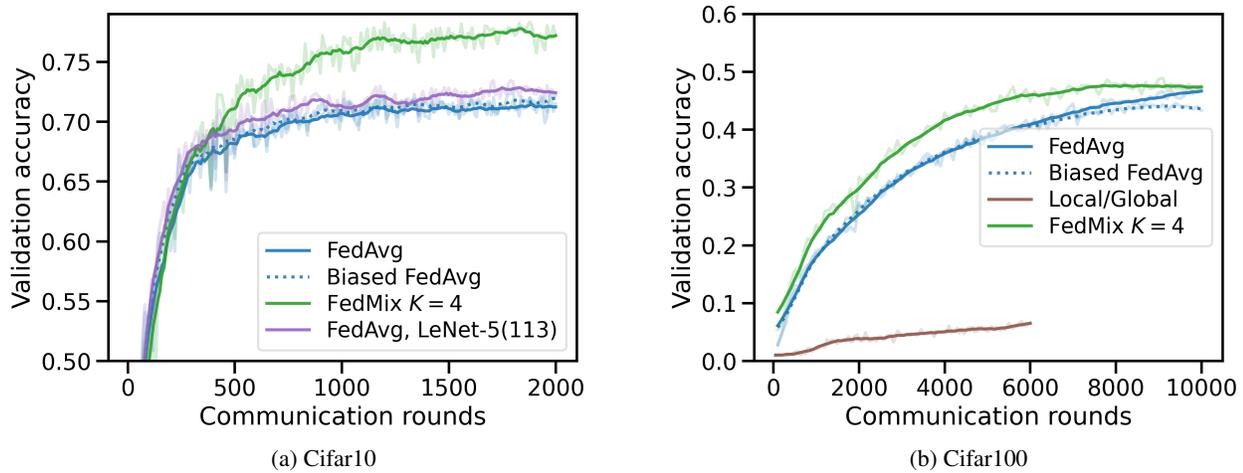


Figure 14. Accuracy on a new client (y-axis) as a function of the amount of communication rounds. Cifar 10 models are trained on the standard 45k training split. Best viewed in color.

hyper parameters,  $B = 10$  and perform early stopping on the local validation sets after no improvement for three epochs. After training, each client communicates their VAE to the server, where we evaluate  $p(s|\mathbf{x}^*)$  to marginalize over the local gating functions according to  $p(z|\mathbf{x}^*) = \sum_{s=1}^S p(z|\mathbf{x}^*, s)p(s|\mathbf{x}^*)$ . With this procedure, FedMix achieves 95.9% test set accuracy, compared to FedAvg with 96.9%. Marginalizing with  $p(s)$  instead of  $p(s|\mathbf{x}^*)$  achieves 96.63%, showing the limitation of the approach in that any error in  $p(s|\mathbf{x}^*)$  propagates into the expert assignment. Reliable out-of-distribution detection capabilities in the individual estimators for  $p(\mathbf{x}|s)$  are therefore necessary.

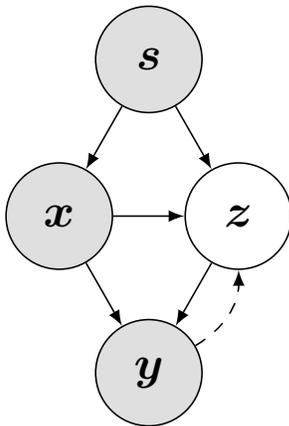


Figure 15. FedMix graphical model. The generative model is depicted with solid lines and the inference model with dashed lines.