

FEDERATED LEARNING WITH MATCHED AVERAGING

Anonymous authors

Paper under double-blind review

ABSTRACT

Federated learning allows edge devices to collaboratively learn a shared model while keeping the training data on device, decoupling the ability to do model training from the need to store the data in the cloud. We propose Federated matched averaging (FedMA) algorithm designed for federated learning of modern neural network architectures *e.g.* convolutional neural networks (CNNs) and LSTMs. FedMA constructs the shared global model in a layer-wise manner by *matching and averaging* hidden elements (*i.e.* channels for convolution layers; hidden states for LSTM; neurons for fully connected layers) with similar feature extraction signatures. Our experiments indicate that FedMA outperforms popular state-of-the-art federated learning algorithms on deep CNN and LSTM architectures trained on real world datasets, while improving the communication efficiency.

1 INTRODUCTION

Edge devices such as mobile phones, sensor networks or vehicles have access to a wealth of data. However, due to concerns raised by data privacy, network bandwidth limitation, and device availability, it's unpractical to gather all local data to the data center and conduct centralized training. To address these concerns, federated learning is emerging (McMahan et al., 2017; Li et al., 2019; Smith et al., 2017; Caldas et al., 2018; Bonawitz et al., 2019) to allow local clients to collaboratively train a shared global model.

The typical federated learning paradigm involves two stages: (i) clients train models over their datasets independently (ii) the data center uploads their locally trained models. The data center then aggregates the received models into a shared global model. One of the standard aggregation methods is FedAvg (McMahan et al., 2017) where parameters of local models are averaged element-wise with weights proportional to sizes of client datasets. FedProx (Sahu et al., 2018) adds a proximal term for client local cost functions, which limits the impact of local updates by restricting them to be close to the global model. Agnostic Federated Learning (AFL) (Mohri et al., 2019), as another variant of FedAvg, optimizes a centralized distribution that is formed by a mixture of the client distributions.

One shortcoming of the FedAvg algorithm is that coordinate-wise averaging of weights may have drastic detrimental effect on the performance and hence hinders the communication efficiency. This issue arises due to the permutation invariant nature of the neural network (NN) parameters, *i.e.* for any given NN there are many variations of it that only differ in the ordering of parameters and constitute local optima which are practically equivalent. Probabilistic Federated Neural Matching (PFNM) (Yurochkin et al., 2019) addresses this problem by finding permutation of the parameters of the NNs before averaging them. PFNM further utilizes Bayesian nonparametric machinery to adapt global model size to heterogeneity of the data. As a result, PFNM has better performance and communication efficiency, however it was only developed for fully connected NNs and tested on simple architectures.

Our contribution In this work (i) we demonstrate how PFNM can be applied to CNNs and LSTMs, however we find that it gives very minor improvement over weight averaging when applied to modern deep neural network architectures; (ii) we propose Federated Matched Averaging (FedMA), a new layers-wise federated learning algorithm for modern CNNs and LSTMs utilizing matching and model size adaptation underpinnings of PFNM; (iii) We empirically study FedMA with real datasets under the federated learning constraints.

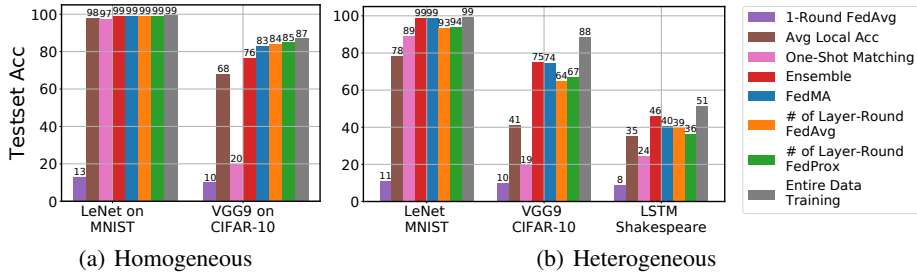


Figure 1: Comparison among various federated learning methods with limited number of communications on LeNet trained on MNIST; VGG-9 trained on CIFAR-10 dataset; LSTM trained on Shakespeare dataset over: (a) homogeneous data partition (b) heterogeneous data partition.

2 FEDERATED MATCHED AVERAGING OF NEURAL NETWORKS

In this section we will discuss permutation invariance classes of prominent neural network architectures and establish the appropriate notion of averaging in the parameter space of NNs. We will begin with the simplest case of a single hidden layer fully connected network, moving on to deep architectures and, finally, convolutional and recurrent architectures.

Permutation invariance of fully connected architectures A basic fully connected (FC) NN can be formulated as $\hat{y} = \sigma(xW^{(1)})W^{(2)}$ (without loss of generality, biases are omitted to simplify notation), where σ is the non-linearity (applied entry-wise). Expanding the preceding expression $\hat{y} = \sum_{i=1}^L W_i^{(2)} \sigma(\langle x, W_i^{(1)} \rangle)$, where $i \cdot$ and $\cdot i$ denote i th row and column correspondingly and L is the number of hidden units. Summation is a permutation invariant operation, hence for any $\{W^{(1)}, W^{(2)}\}$ there are $L!$ practically equivalent parametrizations if this basic NN. It is then more appropriate to write

$$\hat{y} = \sigma(xW^{(1)}\Pi)\Pi W^{(2)}, \text{ where } \Pi \text{ is any } L \times L \text{ permutation matrix.} \quad (1)$$

Recall that permutation matrix acts on rows when applied on the left and on columns when applied on the right. Suppose $\{W^{(1)}, W^{(2)}\}$ are optimal weights, then weights obtained from training on two homogeneous datasets $X_j, X_{j'}$ are $\{W^{(1)}\Pi_j, \Pi_j W^{(2)}\}$ and $\{W^{(1)}\Pi_{j'}, \Pi_{j'} W^{(2)}\}$. It is now easy to see why naive averaging in the parameter space is not appropriate: with high probability $\Pi_j \neq \Pi_{j'}$ and $(W^{(1)}\Pi_j + W^{(1)}\Pi_{j'})/2 \neq \Pi W^{(1)}$ for any Π . To meaningfully average neural networks in the weight space we should first undo the permutation $(W^{(1)}\Pi_j\Pi_j^{-1} + W^{(1)}\Pi_{j'}\Pi_{j'}^{-1})/2 = W^{(1)}$.

2.1 MATCHED AVERAGING FORMULATION

In this section we formulate practical notion of parameter averaging under the permutation invariance. Let w_{jl} be l th neuron learned on dataset j (i.e. l th column of $W^{(1)}\Pi_j$ in the previous example) and $c(\cdot, \cdot)$ be an appropriate similarity function between a pair of neurons. Solution to the following optimization problem are the required inverse permutations:

$$\begin{aligned} \min_{\{\pi_{li}^j\}} \sum_{i=1}^L \sum_{j,l} \min_{\theta_i} \pi_{li}^j c(w_{jl}, \theta_i) \\ \text{s.t. } \sum_i \pi_{li}^j = 1 \forall j, l; \sum_l \pi_{li}^j = 1 \forall i, j. \end{aligned} \quad (2)$$

Then $\Pi_{ji}^{-1} = \pi_{li}^j$ and given weights $\{W_j^{(1)}, W_j^{(2)}\}_{j=1}^J$ provided by J clients, we compute the federated neural network weights $W^{(1)} = \frac{1}{J} \sum_j W_j^{(1)} \Pi_j^{-1}$ and $W^{(2)} = \frac{1}{J} \sum_j \Pi_j^{-1} W_j^{(2)}$. If the client data sizes are imbalanced, similar to Federated Averaging (McMahan et al., 2017), we can use weighted averaging instead of uniform. We refer to this approach as *matched averaging* due to relation of equation 2 to the maximum bipartite matching problem. We note that if $c(\cdot, \cdot)$ is squared

Euclidean distance, we recover objective function similar to k-means clustering, however it has additional constraints on the ‘‘cluster assignments’’ π_{li}^j necessary to ensure that they form permutation matrices.

Solving matched averaging Yurochkin et al. (2019) studied the problem of federated learning of fully connected neural networks. They arrived at a particular form of equation 2 to compute maximum a posteriori estimate (MAP) of their Bayesian nonparametric model based on the Beta-Bernoulli process (BBP) (Thibaux & Jordan, 2007), where similarity $c(w_{jl}, \theta_i)$ is the corresponding posterior probability of j th client neuron l generated from a Gaussian with mean θ_i . Due to the nonparametric aspect, their BBP-MAP inference approach allows number of neurons in the federated model to mildly grow in comparison to the client model sizes. This is an appealing property when client datasets are heterogeneous, which is typical for federated learning. On the downside, their Probabilistic Federated Neural Matching (PFNM) (Yurochkin et al., 2019) is only applicable to fully connected architectures limiting its practicality. On the contrary, our *matched averaging* perspective allows to formulate averaging of widely used architectures such as CNNs and LSTMs as instances of equation 2 and utilize the BBP-MAP as a solver.

2.2 PERMUTATION INVARIANCE OF KEY ARCHITECTURES

Before moving onto the convolutional and recurrent architectures, we discuss permutation invariance in *deep* fully connected networks and corresponding matched averaging approach. We will utilize this as a building block for handling LSTMs and CNN architectures such as VGG (Simonyan & Zisserman, 2014) widely used in practice.

Permutation invariance of deep FCs We extend equation 1 to recursively define deep FC network:

$$x^{(n)} = \sigma(x^{(n-1)}\Pi^{(n-1)}W^{(n)}\Pi^{(n)}), \quad (3)$$

where $n = 1, \dots, N$ is the layer index, $\Pi^{(0)}$ is identity indicating non-ambiguity in the ordering of input features $x = x^{(0)}$ and $\Pi^{(N)}$ is identity for the same in output classes. Conventionally $\sigma(\cdot)$ is any non-linearity except for $\hat{y} = x^{(N)}$ where it is the identity function (or softmax if we want probabilities instead of logits). When $N = 2$, we recover a single hidden layer variant from equation 1. To perform matched averaging of deep FCs obtained from J clients we need to find inverse permutations for every layer of every client. Unfortunately, permutations within any consecutive pair of intermediate layers are coupled leading to a NP-hard combinatorial optimization problem. Instead we consider recursive (in layers) matched averaging formulation. Suppose we have $\{\Pi_j^{(n-1)}\}$, then plugging $\{\Pi_j^{(n-1)}W_j^{(n)}\}$ into equation 2 we find $\{\Pi_j^{(n)}\}$ and move onto next layer. The recursion base for this procedure is $\{\Pi_j^{(0)}\}$, which we know is an identity permutation for any j .

Permutation invariance of CNNs The key observation in understanding permutation invariance of CNNs is that instead of neurons, channels define the invariance. To be more concrete, let $\text{Conv}(x, W)$ define convolutional operation on input x with weights $W \in \mathbb{R}^{C_{in} \times w \times h \times C_{out}}$, where C_{in}, C_{out} are the numbers of input/output channels and w, h are the width and height of the filters. Applying any permutation to the output dimension of the weights and then same permutation to the input channel dimension of the subsequent layer will not change the corresponding CNN’s forward pass. Analogous to equation 3 we can write:

$$x^{(n)} = \sigma(\text{Conv}(x^{(n-1)}, \Pi^{(n-1)}W^{(n)}\Pi^{(n)})). \quad (4)$$

Note that this formulation permits pooling operations as those act within channels. To apply matched averaging for the n th CNN layer we form inputs to equation 2 as $\{w_{jl} \in \mathbb{R}^D\}_{l=1}^{C_{out}^{(n)}}$, $j = 1, \dots, J$, where D is the flattened $C_{in}^{(n)} \times w \times h$ dimension of $\Pi_j^{(n-1)}W_j^{(n)}$. This result can be alternatively derived taking the IM2COL perspective. Similar to FCs, we can recursively perform matched averaging on deep CNNs. The immediate consequence of our CNN permutation invariance discussion is the extension of PFNM (Yurochkin et al., 2019) to CNNs. Empirically (Figure 1) we found that this extension performs well on MNIST with a simpler CNN architecture such as LeNet (LeCun et al., 1998) (4 layers) and significantly outperforms coordinate-wise weight averaging (1 round FedAvg).

However, it breaks down for more complex architecture, e.g. VGG-9 (Simonyan & Zisserman, 2014) (9 layers), needed to obtain good quality prediction on a more challenging CIFAR-10.

Permutation invariance of LSTMs Permutation invariance in the recurrent architectures is associated with the ordering of the hidden states. At a first glance it appears similar to fully connected architecture, however the important difference is associated with the permutation invariance of the hidden-to-hidden weights $H \in \mathbb{R}^{L \times L}$, where L is the number of hidden states. In particular, permutation of the hidden states affects *both* rows and columns of H . Consider a basic RNN $h_t = \sigma(h_{t-1}H + x_tW)$, where W are the input-to-hidden weights. To account for the permutation invariance of the hidden states, we notice that dimensions of h_t should be permuted in the same way for any t , hence

$$h_t = \sigma(h_{t-1}\Pi H \Pi + x_t W \Pi). \quad (5)$$

To match RNNs, the basic sub-problem is to align hidden-to-hidden weights of two clients with Euclidean similarity, which requires minimizing $\|\Pi H_j \Pi - H_{j'}\|_2^2$ over permutations Π . This is a quadratic assignment problem, which is NP-hard. Fortunately, the same permutation appears in an already familiar context of input-to-hidden matching of $W\Pi$. Our matched averaging RNN solution is to utilize equation 2 plugging-in input-to-hidden weights $\{W_j\}$ to find $\{\Pi_j^{-1}\}$. Then federated hidden-to-hidden weights are computed as $H = \frac{1}{j} \sum_j \Pi_j^{-1} H_h \Pi_j^{-1}$ and input-to-hidden weights are computed as before. LSTMs have multiple cell states, each having its individual hidden-to-hidden and input-to-hidden weights. In our matched averaging we stack input-to-hidden weights into $SD \times L$ weight matrix (S is the number of cell states; D is input dimension and L is the number of hidden states) when computing the permutation matrices and then average all weights as described previously. LSTMs also often have an embedding layer, which we handle like a fully connected layer. Finally, we process deep LSTMs in the recursive manner similar to deep FCs.

2.3 FEDERATED MATCHED AVERAGING (FEDMA) ALGORITHM

Defining the permutation invariance classes of CNNs and LSTMs allows us to extend PFNM (Yurochkin et al., 2019) to these architectures, however our empirical study in Figure 1 demonstrates that such extension fails on deep architectures necessary to solve more complex tasks. Our results suggest that recursive handling of layers with matched averaging may entail poor overall solution. To alleviate this problem and utilize the strength of matched averaging on “shallow” architectures, we propose the following layer-wise matching scheme. First, data center gathers *only* the first layer’s weights from the clients and performs one-layer matching described previously to obtain the first layers weights of the federated model. Data center then broadcasts these weights to the clients, which proceed to train all *consecutive* layers on their datasets, keeping the matched federated layers *frozen*. This procedure is then repeated up to the last layer for which we conduct a weighted averaging based on the class proportions of data points per client. We summarize our Federated Matched Averaging (FedMA) in Algorithm 1. The FedMA approach requires communication rounds equal to the number of layers in a network. In Figure 1 we show that with layer-wise matching FedMA performs well on the deeper VGG-9 CNN as well as LSTMs. In the more challenging heterogeneous setting, FedMA outperforms FedAvg, FedProx trained with same number of communication rounds (4 for LeNet and LSTM and 9 for VGG-9) and other baselines, i.e. client individual CNNs and their ensemble.

FedMA with communication We’ve shown that in the heterogeneous data scenario FedMA outperforms other federated learning approaches, however it still lags in performance behind the entire data training. Of course the entire data training is not possible under the federated learning constraints, but it serves as performance upper bound we should strive to achieve. To further improve the performance of our method, we propose *FedMA with communication*, where local clients receive the matched global model at the beginning of a new round and reconstruct their local models with the size equal to the original local models (e.g. size of a VGG-9) based on the matching results of the previous round. This procedure allows to keep the size of the global model small in contrast to a naive strategy of utilizing full matched global model as a starting point across clients on every round.

Algorithm 1: Federated Matched Averaging (FedMA)

Input : local weights of a N -layer architecture $\{W_j^{(1)}, \dots, W_j^{(N)}\}_{j=1}^J$ from J clients
Output: global weights $\{W^{(1)}, \dots, W^{(N)}\}$

```

 $n = 1$ ;
while  $n \leq N$  do
  if  $n < N$  then
     $\{\Pi_j^{-1}\}_{j=1}^J = \text{BBP-MAP}(\{W_j^{(n)}\}_{j=1}^J)$ ; // call BBP-MAP to solve Eq. 2
     $W^{(n)} = \frac{1}{J} \sum_j W_j^{(n)} \Pi_j^{-1}$ ;
  else
     $W^{(n)} = \sum_{k=1}^K \sum_j p_{jk} W_{jl}^{(n)}$  where  $p_k$  is fraction of data points with label  $k$  on worker  $j$ ;
  end
  for  $j \in \{1, \dots, J\}$  do
     $W_j^{(n+1)} \leftarrow \Pi_j^{(n)} W_j^{(n+1)}$ ; // permute the next-layer weights
    Train  $\{W_j^{(n+1)}, \dots, W_j^{(L)}\}$  with  $W^{(n)}$  frozen;
  end
   $n = n + 1$ ;
end

```

3 EXPERIMENTS

We present an empirical study of FedMA with communication and compare it with state-of-the-art methods *i.e.* FedAvg McMahan et al. (2017) and FedProx Sahu et al. (2018); analyze the performance under the growing number of clients and visualize the matching behavior of FedMA to study its interpretability. Our experimental studies are conducted over three real world datasets, detailed information about the datasets and associated models can be found in Table 1.

Experimental Setup We implemented FedMA and the considered baseline methods in PyTorch Paszke et al. (2017). We deploy our empirical study under a simulated federated learning environment where we treat one centralized node in the distributed cluster as the data center and the other nodes as local clients. All nodes in our experiments are deployed on *p3.2xlarge* instances on Amazon EC2. We assume the data center samples all the clients to join the training process for every communication round for simplicity.

Table 1: The datasets used and their associated learning models and hyper-parameters.

Method	MNIST	CIFAR-10	Shakespeare McMahan et al. (2017)
# Data points	60,000	50,000	1,017,981
Model	LeNet	VGG-9	LSTM
# Classes	10	10	80
# Parameters	431k	3,491k	293k
Optimizer	Adam; with AMSGRAD Reddi et al. (2019)		SGD
Hyper-params.	lr: 10^{-3} ; weight decay: 10^{-4}		lr: 0.8(const); momentum: 0.9

For the CIFAR-10 dataset, we use data augmentation (random crops, and flips) and normalize each individual image (details provided in the Supplement). We note that we ignore all batch normalization Ioffe & Szegedy (2015) layers in the VGG architecture and leave it for future work.

For CIFAR-10, we considered two data partition strategies to simulate federated learning scenario: (i) homogeneous partition where each local client has approximately equal proportion of each of the classes; (ii) heterogeneous partition for which number of data points and class proportions are unbalanced. We simulated a heterogeneous partition into J clients by sampling $\mathbf{p}_k \sim \text{Dir}_J(0.5)$ and

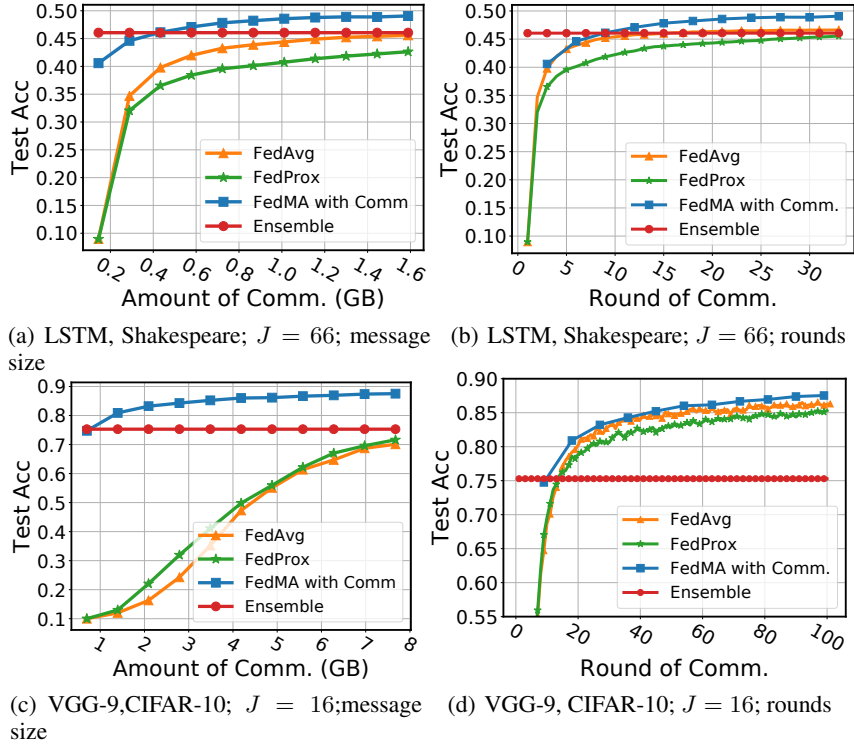


Figure 2: Convergence rates of various methods for VGG-9 trained on CIFAR-10 dataset and LSTM on Shakespeare dataset.

allocating a $\mathbf{p}_{k,j}$ proportion of the training instances of class k to local client j . We use the original test set in CIFAR-10 as our global test set and all test accuracy in our experiments are conducted over that test set. For the Shakespeare dataset, since each speaking role in each play is considered a different client according to Caldas et al. (2018), it’s inherently heterogeneous. We preprocess the Shakespeare dataset by filtering out the clients with datapoints less 10k and get 132 clients in total. We choose 80% of data in training set. We then randomly sample $J = 66$ out of 132 clients in conducting our experiments. We amalgamate all test sets on clients as our global test set.

Communication Efficiency and Convergence Rate In this experiment we study performance of FedMA with communication. Our goal is to compare our method to FedAvg and FedProx in terms of the total message size exchanged between data center and clients (in Gigabytes) and the number of communication rounds (recall that completing one FedMA pass requires number of rounds equal to the number of layers in the local models) needed for the global model to achieve good performance on the test data. We also compare to the performance of an ensemble method. We evaluate all methods under the heterogeneous federated learning scenario on CIFAR-10 with $J = 16$ clients with VGG-9 local models and on Shakespeare dataset with $J = 66$ clients with 1-layer LSTM network. We fix the total rounds of communication allowed for FedMA, FedAvg, and FedProx *i.e.* 11 rounds for FedMA and 99/33 rounds for FedAvg and FedProx for the VGG-9/LSTM experiments respectively. We notice that local training epoch is a common parameter shared by the three considered methods, we thus tune the local training epochs (we denote it by E) (comprehensive results will be discussed in the next section) and report the convergence rate under the best \hat{E} that leads the best end model accuracy over the global test set. We also notice that, there is another hyper-parameter in FedProx *i.e.* the coefficient μ associated with the proxy term, we also tune the parameter using grid search and report the best μ we found *i.e.* 0.001 for both VGG-9 and LSTM experiments. FedMA outperforms FedAvg and FedProx in all scenarios (Figure 2) with its advantage especially pronounced when we evaluate convergence as a function of the message size in Figures 2(a),2(c).

Effect of local training epochs As studied in previous work McMahan et al. (2017); Caldas et al. (2018); Sahu et al. (2018), the number of local training epochs E can affect the perfor-

mance of FedAvg and sometimes lead to divergence. We conduct an experimental study on the effect of E over FedAvg, FedProx, and FedMA on VGG-9 trained on CIFAR-10 under heterogeneous setup. The candidate local epochs we considered are $E \in \{10, 20, 50, 70, 100, 150\}$. For each of the candidate E , we run FedMA for 6 rounds while FedAvg and FedProx for 54 rounds and report the final accuracy that each methods achieves. The result is shown in Figure3. We observed that training longer favors the convergence rate of FedMA, which matches the our assumption that FedMA returns better global model on local models with higher quality. For FedAvg, longer local training leads to deterioration of the final accuracy, which matches the observation in the previous literature McMahan et al. (2017); Caldas et al. (2018); Sahu et al. (2018). FedProx prevents the accuracy deterioration to some extent, however, the accuracy of final model still gets reduced. The result of this experiment suggests that FedMA is the only method that local clients can use to train their model as long as they want.

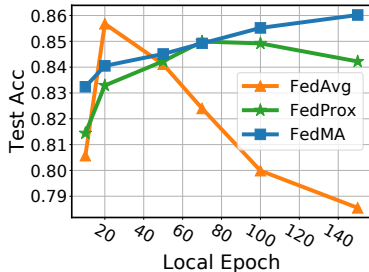


Figure 3: The effect of number of local training epochs on various methods.

Handling data bias Large-scale datasets in the real world are sometimes biased such that for data within class k may consist of different domains *e.g.* geo-diversity. It has been studied that an observable amerocentric and eurocentric bias shown in the widely used ImageNet dataset Shankar et al. (2017); Russakovsky et al. (2015). Classifiers trained on dataset that contains skewed domains (*e.g.* eurocentric training set) in the data space usually perform badly on test set with balanced domains (*e.g.* images that evenly sampled from a broad range of localities) since correlation between domain and classes can dwarf the classifier from the correlation between data feature and classes. We also note that, for different class the dominate domain can be different. We argue that FedMA can handle this type of problem by nature since we split the training tasks that learning the correlation between domain and the class and between data feature to the class to different “clients” and our matching framework can combine those two correlations. We simulate the skewed domain that associated with classes using CIFAR-10 dataset by randomly selecting 5 classes out of 10, and making 95% training images belongs to those classes to be grayscale and 5% training images in the remaining 5 classes to be grayscale. By doing that, we create 5 grayscale images dominated classes and 5 colored images dominated classes. In the test set, there is half grayscale and half colored images for each classes. During the training process, we train the colored images dominated classes and the grayscale images dominated classes on 2 local clients. We then conduct FedMA with communication, FedAvg, and FedProx over the 2 local models the result of this experiment is shown in Figure 4. The entire data training under this skewed domain distributions leads to worse accuracy than normal CIFAR-10 training. FedMA outperforms the entire data training under this scenario.

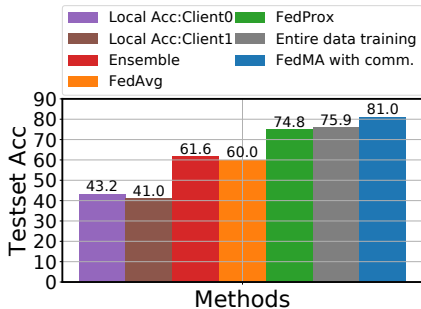


Figure 4: Performance on skewed CIFAR-10 dataset.

Table 2: Additional information from convergence experimental results associated with VGG-9 trained on CIFAR-10 shown in Figure2

Method	FedAvg	FedProx	Ensemble	FedMA
Final Accuracy(%)	86.29	85.32	75.29	87.53
Best local epoch(E)	20	20	N/A	150
Model growth rate	1×	1×	16×	1.11×
Hyper-param.	N/A	$\mu = 0.001$	N/A	$\gamma = 5.0$

Data efficiency It is known that deep learning models perform better when more training data is available. However, under the federated learning constraints, data efficiency has not been stud-

Table 3: Additional information from convergence experimental results associated with LSTM trained on Shakespeare shown in Figure 2

Method	FedAvg	FedProx	Ensemble	FedMA
Final Accuracy(%)	46.63	45.83	46.06	49.07
Best local epoch(E)	2	5	N/A	5
Model growth rate	1 \times	1 \times	66 \times	1.06 \times
Hyper-param.	N/A	$\mu = 0.001$	N/A	$\gamma = 10^{-3}$

ied to the best of our knowledge. The challenge here is that when new clients join the federated system, they each bring their own version of the data distribution, which, if not handled properly, may deteriorate the performance despite the growing data size across the clients. To simulate this scenario we first partition the entire training CIFAR-10 dataset into 5 homogeneous pieces. We then partition each homogeneous data piece further into 5 sub-pieces heterogeneously. Using this strategy, we partition the CIFAR-10 training set into 25 heterogeneous small sub-datasets containing approximately 2k points each. We conduct a 5-step experimental study: starting from a randomly selected homogeneous piece consisting of 5 associated heterogeneous sub-pieces, we simulate a 5-client federated learning heterogeneous problem. For each consecutive step, we add one of the remaining homogeneous data pieces consisting of 5 new clients with heterogeneous sub-datasets. Results of this experiment are presented in Figure 5. Performance of FedMA (with a single pass) improves when new clients are added to the federated learning system, while FedAvg with 9 communication rounds deteriorates.

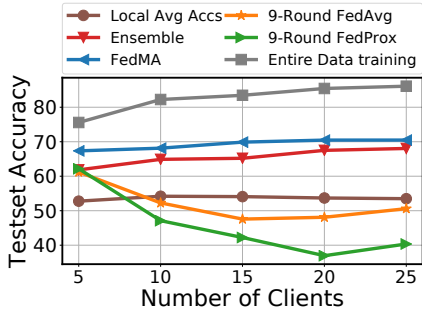


Figure 5: Data efficiency under the increasing number of heterogeneous clients.

Interpretability One of the strengths of FedMA is that it utilizes communication rounds more efficiently than FedAvg. Instead of directly averaging weights element-wise, FedMA first identifies matching groups of convolutional filters and then averages them into the global convolutional filters. It’s natural to ask “*How does the matched filters look like?*”. We visualize the representations generated by a pair of matched local filters, aggregated global filter, and the filter returned by the FedAvg method over the same input image. The result is shown in Figure 6. We observe that the matched filters and the global filter found with FedMA are extracting the same feature of the input image, i.e. filter 0 of client 1 and filter 23 of client 2 are extracting the position of the legs of the horse, and the corresponding matched global filter 0 does the same. For the FedAvg, global filter 0 is the average of filter 0 of client 1 and filter 0 of client 2, which clearly tampers the leg extraction functionality of filter 0 of client 1.

4 CONCLUSION

In this paper, we presented FedMA, a new layer-wise federated learning algorithm designed for modern CNNs and LSTMs architectures utilizing probabilistic matching and model size adaptation. We demonstrate the convergence rate and communication efficiency of FedMA empirically. In the future, we would like to extend FedMA towards finding the optimal averaging strategy. Making FedMA support more building blocks *e.g.* residual structures in CNNs and batch normalization layers is also of interest.

REFERENCES

Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir Ivanov, Chloe Kiddon, Jakub Konecny, Stefano Mazzocchi, H Brendan McMahan, et al. Towards federated learning at scale: System design. *arXiv preprint arXiv:1902.01046*, 2019.

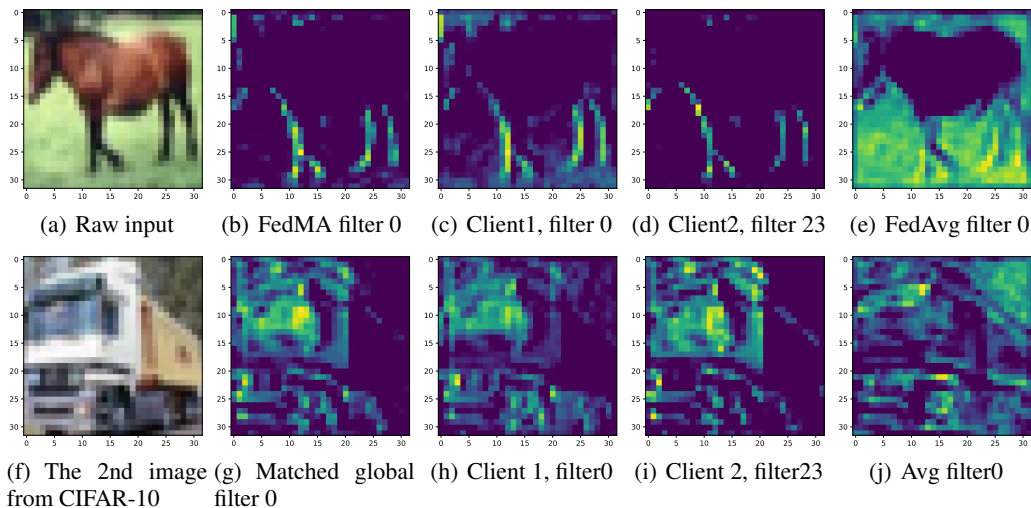


Figure 6: Representations generated by the first convolution layers of locally trained models, FedMA global model and the FedAvg global model.

Sebastian Caldas, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMahan, Virginia Smith, and Ameet Talwalkar. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097*, 2018.

Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.

Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *arXiv preprint arXiv:1908.07873*, 2019.

Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pp. 1273–1282, 2017.

Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In *International Conference on Machine Learning*, pp. 4615–4625, 2019.

Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017.

Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *arXiv preprint arXiv:1904.09237*, 2019.

Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.

Anit Kumar Sahu, Tian Li, Maziar Sanjabi, Manzil Zaheer, Ameet Talwalkar, and Virginia Smith. On the convergence of federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*, 2018.

Shreya Shankar, Yoni Halpern, Eric Breck, James Atwood, Jimbo Wilson, and D Sculley. No classification without representation: Assessing geodiversity issues in open data sets for the developing world. *arXiv preprint arXiv:1711.08536*, 2017.

Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S Talwalkar. Federated multi-task learning. In *Advances in Neural Information Processing Systems*, pp. 4424–4434, 2017.

Romain Thibaux and Michael I Jordan. Hierarchical Beta processes and the Indian buffet process. In *Artificial Intelligence and Statistics*, pp. 564–571, 2007.

Mikhail Yurochkin, Mayank Agarwal, Soumya Ghosh, Kristjan Greenewald, Nghia Hoang, and Yasaman Khazaeni. Bayesian nonparametric federated learning of neural networks. In *International Conference on Machine Learning*, pp. 7252–7261, 2019.

A DETAILS OF CNN ARCHITECTURE AND HYPER-PARAMETERS

Table 4: Detailed information of the VGG-9 architecture used in our experiments, all non-linear activation function in this architecture is ReLU; the shapes for convolution layers follows (C_{in}, C_{out}, c, c)

Parameter	Shape	Layer hyper-parameter
layer1.conv1.weight	$3 \times 32 \times 3 \times 3$	stride:1;padding:1
layer1.conv1.bias	32	N/A
layer2.conv2.weight	$32 \times 64 \times 3 \times 3$	stride:1;padding:1
layer2.conv2.bias	64	N/A
pooling.max	N/A	kernel size:2;stride:2
layer3.conv3.weight	$64 \times 128 \times 3 \times 3$	stride:1;padding:1
layer3.conv3.bias	128	N/A
layer4.conv4.weight	$128 \times 128 \times 3 \times 3$	stride:1;padding:1
layer4.conv4.bias	128	N/A
pooling.max	N/A	kernel size:2;stride:2
dropout	N/A	$p = 5\%$
layer5.conv5.weight	$128 \times 256 \times 3 \times 3$	stride:1;padding:1
layer5.conv5.bias	256	N/A
layer6.conv6.weight	$256 \times 256 \times 3 \times 3$	stride:1;padding:1
layer6.conv6.bias	256	N/A
pooling.max	N/A	kernel size:2;stride:2
dropout	N/A	$p = 10\%$
layer7.fc7.weight	4096×512	N/A
layer7.fc7.bias	512	N/A
layer8.fc8.weight	512×512	N/A
layer8.fc8.bias	512	N/A
dropout	N/A	$p = 10\%$
layer9.fc9.weight	512×10	N/A
layer9.fc9.bias	10	N/A

Table 5: Detailed information of the LSTM architecture in our experiment

Parameter	Shape
encoder.weight	80×8
lstm.weight.ih.10	1024×8
lstm.weight.hh.10	1024×256
lstm.bias.ih.10	1024
lstm.bias.hh.10	1024
decoder.weight	80×256
decoder.bias	80

B DATA AUGMENTATION AND NORMALIZATION DETAILS

In preprocessing the images in CIFAR-10 dataset, we follow the standard data augmentation and normalization process. For data augmentation, random cropping and horizontal random flipping are used. Each color channels are normalized with mean and standard deviation by $\mu_r = 0.491372549, \mu_g = 0.482352941, \mu_b = 0.446666667, \sigma_r = 0.247058824, \sigma_g = 0.243529412, \sigma_b = 0.261568627$. Each channel pixel is normalized by subtracting the mean value in this color channel and then divided by the standard deviation of this color channel.

C EXTRA EXPERIMENTAL DETAILS

Here we report the shapes of final global VGG and LSTM models returned by FRB with communication.

Table 6: Detailed information of the final global VGG-9 model returned by FRB; the shapes for convolution layers follows (C_{in}, C_{out}, c, c)

Parameter	Shape	Growth rate (#global / #original params)
layer1.conv1.weight	$3 \times 47 \times 3 \times 3$	$1.47 \times (1, 269/864)$
layer1.conv1.bias	47	$1.47 \times (47/32)$
layer2.conv2.weight	$47 \times 79 \times 3 \times 3$	$1.81 \times (33, 417/18, 432)$
layer2.conv2.bias	79	$1.23 \times (79/64)$
layer3.conv3.weight	$79 \times 143 \times 3 \times 3$	$1.38 \times (101, 673/73, 728)$
layer3.conv3.bias	143	$1.12 \times (143/128)$
layer4.conv4.weight	$143 \times 143 \times 3 \times 3$	$1.24 \times (184, 041/147, 456)$
layer4.conv4.bias	143	$1.12 \times (143/128)$
layer5.conv5.weight	$143 \times 271 \times 3 \times 3$	$1.18 \times (348, 777/294, 912)$
layer5.conv5.bias	271	$1.06 \times (271/256)$
layer6.conv6.weight	$271 \times 271 \times 3 \times 3$	$1.12 \times (660, 969/589, 824)$
layer6.conv6.bias	271	$1.06 \times (271/256)$
layer7.fc7.weight	4336×527	$1.09 \times (2, 285, 072/2, 097, 152)$
layer7.fc7.bias	527	$1.02 \times (527/512)$
layer8.fc8.weight	527×527	$1.05 \times (277, 729/262, 144)$
layer8.fc8.bias	527	$1.02 \times (527/512)$
layer9.fc9.weight	527×10	$1.02 \times (5, 270/5, 120)$
layer9.fc9.bias	10	$1 \times$
Total Number of Parameters	3, 900, 235	$1.11 \times (3, 900, 235/3, 491, 530)$

Table 7: Detailed information of the LSTM architecture in our experiment

Parameter	Shape	Growth rate (#global / #original params)
encoder.weight	80×21	$2.63 \times (1,680/640)$
lstm.weight.ih.l0	1028×21	$2.64 \times (21,588/8,192)$
lstm.weight.hh.l0	1028×257	$1.01 \times (264,196/262,144)$
lstm.bias.ih.l0	1028	$1.004 \times (1,028/1,024)$
lstm.bias.hh.l0	1028	$1.004 \times (1,028/1,024)$
decoder.weight	80×257	$1.004 \times (20,560/20,480)$
decoder.bias	80	$1 \times$
Total Number of Parameters	310,160	$1.06 \times (310,160/293,584)$