

---

# Building Large Machine Learning Models from Small Distributed Models: A Layer Matching Approach

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Cross-device federated learning (FL) enables massive amount of clients to col-  
2 laborate to train a machine learning model with local data. However, the com-  
3 putational resource of the client devices restricts FL from utilizing large modern  
4 machine learning models that requires sufficient computation. In this paper, we  
5 propose a federated layer matching algorithm that enables the server to build a  
6 deep server machine learning model from relatively shallow client models. The  
7 federated layer matching (FLM) algorithm dynamically averages similar layers in  
8 the client models to the server model, and inserts dissimilar layers as new layers to  
9 the server model. With the proposed algorithm, the clients are able to train small  
10 models based on device capacity, while the server can still obtain a larger and  
11 more powerful server model from the clients with decentralized data. Our numer-  
12 ical experiments show that the proposed FLM algorithm is able to build a server  
13 model 40% larger than the client models, and such a model performs much better  
14 than the model obtained by the classical FedAvg, when using the same amount of  
15 communication resource.

## 16 1 Introduction

17 Machine learning has been widely used in various applications, such as computer vision (CV) and  
18 natural language processing (NLP). To process large amounts of complex data (image, text, time-  
19 series, graph), machine learning (ML) models, especially deep learning models, have become ex-  
20 tremely large in their widths and depths. For example, GPT-3[1] has 175 billion model parameters  
21 with 96 layers; RegNetY [2] has 145 million model parameters and 62 layers; ViT [3] has 632  
22 million model parameters with 32 layers.

23 Even though large models have become state-of-the-art for many application domains, there are  
24 several critical drawbacks to building these models. First of all, training such huge models requires  
25 tremendous amounts of data as well as computation resources. Second, in many applications (such  
26 as those in federated learning), there is a massive amount of data available but distributed among  
27 mobile devices such as laptops, tablets, and mobile phones, which have limited bandwidth, memory,  
28 and computation resources. As a result, it becomes impossible to deploy and train huge ML models  
29 on resource-constrained devices.

30 In this work, we propose a novel approach capable of pooling distributed and heterogeneous re-  
31 sources (both data and computation) to train a large ML model. This is made possible by making  
32 the key observation that modern ML models often have certain *stacked* structures. That is, by repet-  
33 itively stacking certain components on top of each other, one can build larger and larger models. For  
34 example, ResNet-50 repeats the third Bottleneck layer 6 times while ResNet-101 repeats the same  
35 layer 23 times [4]; Bert-small stacks 4 identical attention layers while Bert-large stacks 24 layers [5].  
36 With this observation, we ask and hope to address the following research question:

(Q) Is it possible to progressively and efficiently aggregate small client-side models to construct deep and powerful server-side models, without degrading model performance?

37  
38

39 **Our contributions.** In this project, we address the above question (Q) by considering the situation  
40 that many clients are connected to a server, and the clients possess private data and some computa-  
41 tion power. The goal is to collectively build a large and deep model by utilizing as much local data  
42 and computation power as possible while minimizing the communication overhead in the system.  
43 Towards this end, we propose a *federated layer matching* (FLM) algorithm, which is, to our knowl-  
44 edge, the first method capable of dynamically constructing a large and deep neural network (located  
45 at the server) from shallow ones (located at the clients), by using data distributed over the clients.  
46 The proposed algorithm enjoys many desirable features, such as low communication cost and strong  
47 theoretical guarantees. More specifically, the main contributions of this work are summarized below.

- 48 • **From shallow to deep.** To our knowledge, this is the first work that enables distributed and  
49 systematic construction of a *deep* neural network from *shallow* ones, while the closest existing  
50 works, such as FedMA [6], can only construct *wider* networks with *fixed* depth.
- 51 • **Flexibility.** The proposed method has been carefully designed so that it is *flexible*, in the sense  
52 that it can be used to build a wide range of different neural network structures encountered in  
53 CV and NLP applications. This is made possible by the novel design of certain *layer matching*  
54 mechanism, to be introduced in Section 2.2.
- 55 • **Communication efficiency.** In contrast to the existing (horizontal) FL algorithms such as Fe-  
56 dAvg, which require frequent exchange of the *entire* model between the clients and the server, the  
57 proposed method is communication efficient because it can better utilize the heterogeneous client  
58 models by the proposed *layer matching* technique, rather than simply averaging them.
- 59 • **Superior empirical performance.** Finally, we conduct empirical experiments on ResNet with  
60 Cifar-10 dataset, showing that the proposed algorithm has better performance (measured by pre-  
61 diction accuracy, communication/sample efficiency, etc.) compared with state-of-the-art (SOTA)  
62 algorithms. In certain cases, its performance can even approach that of centralized training algo-  
63 rithms.

## 64 2 Preliminary and Problem Setup

### 65 2.1 Problem Setup

66 Consider the distributed problem with loss function  $f(\cdot)$  and  $K$  clients. Suppose each client has  
67 dataset  $\mathcal{D}_k$ . We aim to solve the following problem:

$$\min_{\Theta} \sum_{k=1}^K f(\Theta; \mathcal{D}_k), \quad (1)$$

68 where  $\Theta$  denotes model parameters. To model the *stacked* structure of the ML models, let us assume  
69 that the client  $k$ 's model parameters, denoted by  $\Theta_k$ , consists of  $M$  blocks  $\Theta_k := [\Theta_{k,1}; \dots; \Theta_{k,M}]$ ,  
70 where  $\Theta_{k,m}$  represents the parameters of the  $m$ th block. Each block  $m$  is stacked with  $L_{k,m}$   
71 layers of the *same* structure. Thus, the parameter of the  $m$ th block can be further written as  
72  $\Theta_{k,m} = [\Theta_{k,m}[1]; \dots; \Theta_{k,m}[L_{k,m}]]$ , where  $\{\Theta_{k,m}[l]\}_{l=1}^{L_{k,m}}$  have the same dimension, and the  
73 model of client  $k$  has  $L_k = \sum_{m=1}^M L_{k,m}$  layers in total. In addition, let us use 0 as the index  
74 of the server, e.g., the server model is denoted as  $\Theta_0$ .

75 Federated Learning (FL) studies the setting of solving (1) in a parameter-server/client system. Fe-  
76 dAvg [7] is a well-studied solution of FL that adopts the computation-then-aggregation strategy. In  
77 the algorithm, the clients  $k = 1, \dots, K$  locally perform a few steps of model updates by optimizing:

$$\min_{\Theta_k} f(\Theta_k; \mathcal{D}_k). \quad (2)$$

78 Then, the server aggregates the updated local models and averages them before sending the updated  
79 global model back to the clients, that is, it performs:

$$\Theta_0 = \frac{1}{K} \sum_{k=1}^K \Theta_k.$$

80 However, due to the limited communication and computation resource on the mobile devices, typical  
 81 FL algorithms can only be used to train small models such as MobileNet or MLPs. These models  
 82 have few parameters than the SOTA models and have sub-optimal performance.

83 A more recent approach is to ensemble the client models to construct a larger server model [8, 6,  
 84 9], which aims to train small models on the clients and ensemble the client models into a large  
 85 server model. In this case, the clients still solve (2) locally with multiple updates. Then, the server  
 86 aggregates the client models  $\{\Theta_k\}_{k=1}^K$  and builds the server model  $\Theta_0$  by a linear or non-linear  
 87 transformation, denoted by

$$\Theta_0 = \text{Agg}(\Theta_1, \dots, \Theta_K).$$

88 In these approaches, the server model  $\Theta_0$  has more parameters than the client models  $\Theta_k$ .

89 However, this line of work still has scalability, client resource requirement, and computation effi-  
 90 ciency limitations. In our work, we focus on a new model ensemble approach: the *layer matching*  
 91 technique. The key idea of our approach is to train shallow models in the clients based on the client  
 92 resource capacity, and the server properly stacks the layers of the client model into a deeper server  
 93 model. Compared with FedAvg, our approach can train a deeper server model with better perfor-  
 94 mance and the ability to fully utilize the massive data and limited computation resources on a large  
 95 number of clients.

## 96 2.2 Related Work

97 **Federated Learning:** The FL problems typically consider the setting that the clients are directly  
 98 connected to a parameter-server and that the communication at the server is the bottleneck of the  
 99 system. The FL algorithms, such as the well-known FedAvg [7], perform multiple local updates  
 100 before one communication step. However, when the data is *heterogeneous* among the agents, it  
 101 is difficult for these algorithms to achieve convergence [10, 11]. Recent algorithms such as the  
 102 FedProx [12], SCAFFOLD [13] and FedPD [14] have developed new techniques to improve upon  
 103 FedAvg. These algorithms require the server and client models to have the same size to perform  
 104 model averaging. This requirement restricts the algorithms from training large ML models when (a  
 105 part of) the clients have limited computation capacity.

106 **Model Ensembling:** Model ensemble method have been widely used in FL [8, 15, 16]. They are  
 107 used to construct server models out of collected client models. In [17, 18], the authors propose  
 108 a straightforward method that directly concatenates local models into a wide model. The server  
 109 model obtained by the above-mentioned *stacking* method scales *linearly* with the client number, and  
 110 unfortunately, it omits the interconnection among the client models. These properties can result in  
 111 a waste of memory and computation resource and have less scalability. A perhaps smarter way of  
 112 performing the model ensemble is the neural matching method [6]. This method tries to identify  
 113 the similarity in hidden elements (e.g., hidden states in LSTM, convolution channels in CNN) and  
 114 match them in different client models and construct a wider server model. Note that, this method  
 115 only matches the parameters in the same layer to extend the width of the aggregated model and  
 116 preserves the depth of the client models. However, many studies have theoretically and empirically  
 117 suggested that the depth of the neural network is critical in strengthening the network representation  
 118 ability and increasing network performance [19, 20].

## 119 3 Federated Layer Matching

### 120 3.1 Algorithm Description

121 In this section, we discuss the proposed layer matching algorithm. The key idea is that, the server  
 122 matches the layers in different client models based on the *similarity* of the layer’s parameters instead  
 123 of the layer’s position. In specific, for each block  $m$  (i.e., a residual block in ResNet), the server  
 124 first aggregates a total of  $\sum_{k=1}^K L_{k,m}$  layers of the same structure. It then clusters these layers by  
 125 the similarity between their parameters  $\Theta_{k,m}[l]$ , and such clustering operation will result in  $L_{0,m}$   
 126 clustered layers, which can be deeper than the client models, i.e.,  $L_{0,m} \geq \max\{L_{k,m}\}_{k=1}^K$ . By  
 127 stacking the clustered layers, the server can construct a deeper server model.

128 Before going into the detailed explanation of the algorithm, let us first define some notations for  
 129 convenience. First, let us define the layer matching pattern for client  $k$  as a matching matrix  $\Pi_k \in$

---

**Algorithm 1** Federated Layer Matching Algorithm
 

---

1: **Inputs:** Data  $\mathcal{D}_k$ , matching frequency  $R$ , maximal iteration  $T$ , initial model  $\Theta_0^0$ .  
 2: **Initialization:** The server broadcasts the initial server model  $\Theta_0^0$  to all clients.  $\Pi_k^0 = I$   
 3: **for**  $t = 1, \dots, T$  **do**  
 4:   **for** Client  $k = 1, \dots, K$  **in parallel do**  
 5:     **Step 1:** Optimize  $\Theta_k^t = \arg \min_{\Theta_k} \mathcal{L}_k(\Theta_k; \Theta_0^{t-1})$ ;  
 6:     Send updated model  $\Theta_k^t$  to server.  
 7:   **end for**  
 8:   **Server:**  
 9:   **if**  $t \bmod R = 0$  **then**  
 10:     **Step 2:** Solves the matching problem (4) and obtains the layer matching pattern  $\Pi_k^t$   
 11:   **else**  
 12:     Uses the matching pattern at the previous round, namely to let  $\Pi_k^t = \Pi_k^{t-1}$   
 13:   **end if**  
 14:   **Step 3:** Construct the server model:  $\Theta_0^t = \frac{\sum_{k=1}^K \Pi_k^t \Theta_k^t}{\sum_{k=1}^K \Pi_k^t \mathbb{1}_{L_k}}$ .  
 15:   **Step 4:** Send the updated server model and matching patterns  $\Pi_k^{t \top} \Theta_0^t$  to the clients  
 16: **end for**  
 17: **Return:** Parameters  $\Theta_0^T$ .

---

130  $\{0, 1\}^{L_0 \times L_k}$ , where  $\Pi_k[l_0, l_k] = 1$  if layer  $l_k$  of the client model matches layer  $l_0$  of the server model  
 131 and  $\Pi_k[l_0, l_k] = 0$  if layer  $l_k$  of the client model does not match layer  $l_0$  of the server model. As  
 132 the layers can only be matched to the layer in the same block, the matching matrix must be a block  
 133 diagonal matrix with  $M$  blocks  $\Pi_k = \text{diag}\{\Pi_{k,1}, \dots, \Pi_{k,M}\}$ , where  $\Pi_{k,m} \in \{0, 1\}^{L_{0,m} \times L_{k,m}}$ .

134 Further, we define the local loss function as

$$\mathcal{L}_k(\Theta_k; \Theta_0) \triangleq f(\Theta_k; \mathcal{D}_k) + r(\Theta_k, \Theta_0),$$

135 where  $r$  is a regularizer defined as:

$$r(\Theta_k, \Theta_0) \triangleq \left\| \Theta_k - \Pi_k^\top \Theta_0 \right\|_2^2.$$

136 Finally, let  $\mathbb{1}_L := [1; 1; \dots; 1] \in \mathbb{R}^L$  denote the all-one vector of size  $L$ .

137 The algorithm at each round of training consists of the following four major steps. Algorithm 1  
 138 provides a detailed description of the algorithm.

- 139 • Step 1: Each client  $k \in [K]$  updates local model  $\Theta_k$  by optimizing the local loss function  $\mathcal{L}_k$ , and  
 140 sends the local model to the server;
- 141 • Step 2: The server cluster the layers in the aggregated client models based on the similarity be-  
 142 tween the layer parameters. By optimizing the layer matching problem (4), the server obtains the  
 143 layer matching patterns  $\Pi_k$  for all clients;
- 144 • Step 3: The server constructs the server model by stacking the clustered layers and reconstruct the  
 145 client models with the server model and the layer matching patterns;
- 146 • Step 4: The server sends the reconstructed client models to each client and starts the next round  
 147 of training.

148 Next, let us provide a detailed explanation to the key step 2 in the proposed algorithm. More detailed  
 149 description of the algorithm are given in Appendix A.

150 **Step 2:** The server-side layer matching is the key step for the proposed FLM algorithm. At the iter-  
 151 ation to perform layer matching  $t \bmod R = 0$ , for each block  $m \in [M]$  in parallel, we sequentially  
 152 match each client layers to the server layers for  $P$  iterations. The matching procedure is illustrated  
 153 in Algorithm 2, which consists four stages, and the major stages 2 and 3 are discussed below.

154 **Stage 2-2:** In this stage, we compute the cost for each client layer  $l_k = 1, \dots, L_{k,m}$  in the block to  
 155 a layer in the server. Assume that we set the maximum number of layers in block  $m$  in the server

---

**Algorithm 2** Layer Matching Step 2
 

---

- 1: **Inputs:** Client models  $\{\Theta_k^t\}_{k=1}^K$ , old server model  $\Theta_0^{t-1}$ , and matching patterns  $\{\Pi_k^{t-1}\}_{k=1}^K$ .
  - 2: **Initialize:**  $\Pi_k^{t,0} = \Pi_k^{t-1}, \forall k \in \{1, \dots, K\}$
  - 3: **for** Block  $m = 1, \dots, M$  in parallel **do**
  - 4:   **for**  $p = 0, \dots, P - 1$  **do**
  - 5:     Randomly select client  $k$  in  $\{1, \dots, K\}$  at random;
  - 6:     **Stage 2-1:** Construct server model  $\Theta_{0,m}^{t,p} = \frac{\sum_{k' \neq k} \Pi_{k',m}^{t,p} \Theta_{k',m}^t}{\sum_{k' \neq k} \Pi_{k',m}^{t,p} \mathbb{1}_{L_{k',m}}}$
  - 7:     **Stage 2-2:** Compute cost matrix  $C_{k,m}^t$  matching  $\Theta_{k,m}^t$  to  $\Theta_{0,m}^{t,p}$  with (3);
  - 8:     **Stage 2-3:** Solve matching pattern  $\Pi_{k,m}^{t,p+1}$  of block  $m$  of client  $k$  by optimizing (4);
  - 9:     **Stage 2-4:** Update all other matching patterns  $\{\Pi_{k',m}\}_{k' \neq k}$ .
  - 10:    **end for**
  - 11: **end for**
  - 12: **Return:** Updated matching patterns  $\{\Pi_k^t = \Pi_k^{t,P}\}_{k=1}^K$ .
- 

156 model as  $\bar{L}_{0,m}$ . Then the costs for matching all layers in block  $m$  for client  $M$  form a cost matrix  
 157  $C_{k,m}^{t,p} \in \mathbb{R}^{\bar{L}_{0,m} \times L_{k,m}}$ , where its entries are defined as:

$$C_{k,m}^{t,p}[l_0, l_k] = c_1(\Theta_{k,m}^t[l_k], \Theta_{0,m}^{t,p}[l_0]) + c_2(\Theta_{k,m}^t[l_k], l_0, \bar{L}_{0,m}), \quad (3)$$

158 where  $c_1(\cdot)$  is the *fidelity score* measures the similarity between the client layer and the server layer,  
 159 and  $c_2(\cdot)$  is the *model complexity penalty* for increasing the number of the server layer. A reasonable  
 160 choice for  $c_1(\cdot)$  is a distance metric (e.g., Euclidean distance or cosine angle) between the server and  
 161 client layers, and the choice for  $c_2(\cdot)$  can be the model selection metric (e.g.,  $K \log(L_{l_0})$  used in  
 162 BIC) on the model size.

163 **Stage 2-3** With the cost matrix  $C_{k,m}^{t,p}$  that describes the cost to match each client layer to the server  
 164 layers, we can optimize the corresponding linear assignment problem:

$$\min_{\Pi_{k,m}^{t,p+1}} \text{Trace} \left( (\Pi_{k,m}^{t,p+1})^\top C_{k,m}^{t,p} \right), \quad \text{s.t. } \Pi_{k,m}^{t,p+1} \in \{0, 1\}^{\bar{L}_{0,m} \times L_{k,m}}, (\Pi_{k,m}^{t,p+1})^\top \mathbb{1}_{\bar{L}_{0,m}} = \mathbb{1}_{L_{k,m}}, \quad (4)$$

165 and obtain the updated matching pattern  $\Pi_{m,k}^{t,p+1}$ . If any of the rows  $L_{0,m} + 1, \dots, \bar{L}_{0,m}$  in the  
 166 matching pattern are not all-zero, then some of the client layers are added to the server model as  
 167 new layers, and the server model becomes deeper. Otherwise, the server model keeps the same  
 168 depth  $L_{0,m}$  or even gets shallower.

169 By repeating Stage 2-1 to 2-4 for  $P$  iterations ( $P \geq K$ ), all the client matching patterns  $\{\Pi_k^t\}_{k=1}^K$   
 170 will be updated. This concludes the entire procedure of Step 2 in Algorithm 1.

171 **Step 3:** In this step, the server constructs the server model with all the updated client matching  
 172 patterns  $\{\Pi_k^t\}_{k=1}^K$  and the client models  $\{\Theta_k^t\}_{k=1}^K$ , by averaging the client layers that match to the  
 173 same server layer, i.e.,

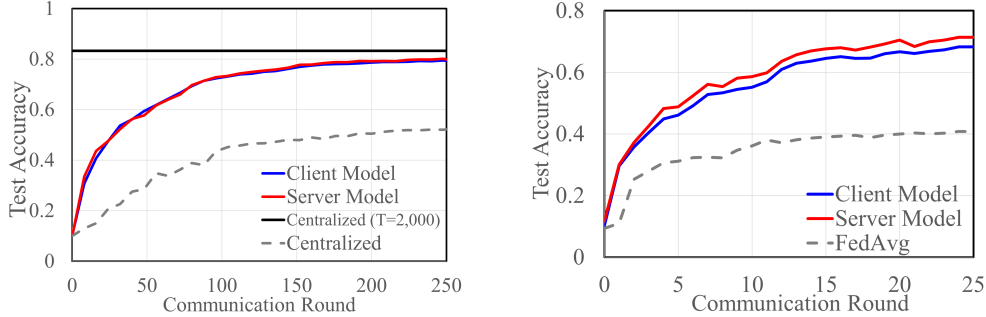
$$\Theta_{0,m}^t[l_0] = \frac{1}{|\{l_k | \Pi_{k,m}^t[l_0, l_k] = 1\}|} \sum_{l_k \in \{l_k | \Pi_{k,m}^t[l_0, l_k] = 1\}} \Theta_{k,m}^t[l_k].$$

174 Now the server model has been constructed, the server then reconstructs the client model with  $\Theta_k^t =$   
 175  $\Pi_k^t \Theta_0^t$ . Such an updated client model will then be used for the local training at the next iteration  
 176  $t + 1$ .

### 177 3.2 Discussion

178 Before we close this section, let us provide some discussions about the algorithm.

179 First, let us discuss the computation/communication resources used at the clients and the servers. In  
 180 each round of training, each client trains a shallow network with parameter  $\Theta_k$ , sends the shallow



(a) Testing accuracies of the client models, the server model, centralized trained model with  $T = 2, 500$  iterations, and the centralized trained model with the same number of updates. (b) Testing accuracies of the client models, the server model trained with FLM, and the model trained with FedAvg with the same number of communication rounds.

Figure 1: Test accuracy of a) IID distributed data, b) Non-IID distributed data on Cifar-10 dataset.

181 network  $\Theta_k^t$  and receives the model  $\Pi_k^t \Theta_0^t$ . The resource cost on the client is the same as Fed-  
 182 Prox [12]. The server needs to solve the layer matching problem described in **Step 3**, which has  
 183 computation complexity of  $\mathcal{O}(K \times \sum_{m=1}^M N_m^3 \times P)$ , and scales linearly with the number of clients.  
 184 Second, let us remark that although the algorithm appears to be rather complicated, in fact its math-  
 185 ematical interpretation is relatively simple. Instead of directly solving (1), we solve the following  
 186 *regularized* problem:

$$\min_{\Theta_0, \{\Theta_k, \Pi_k\}_{k=1}^K} \frac{1}{K} \sum_{k=1}^K \left( f(\Theta_k; \mathcal{D}_k) + \left\| \Theta_k - \Pi_k^\top \Theta_0 \right\|^2 \right). \quad (5)$$

187 Note that not only the client/server parameters  $\Theta_k$ 's are  $\Theta_0$  are optimized, the matching pattern  $\Pi_k$ 's  
 188 are also getting optimized. The client and the server alternately optimize  $\{\Theta_k\}$  and  $\Theta_0, \{\Pi_k\}$  re-  
 189 spectively by using block coordinate descent (BCD)-type algorithm. In the server-side optimization,  
 190 the shape of  $\Pi_k$  and  $\Theta_0$  are dynamically changing based on the heterogeneity of the layers in the  
 191 client models. Following the line of the proof in [21], we can show that the FLM algorithm generates  
 192  $\Theta_0, \{\Theta_k, \Pi_k\}_{k=1}^K$  that converge to the first-order stationary point of (5) with rate  $\mathcal{O}(1/T)$ .

## 193 4 Numerical Experiments

194 We run the experiment based on ResNet model designed for Cifar-10 dataset [22]. In the experiment,  
 195 we use ResNet[ $a, b, c$ ] to denote the number of residual layers in the three residual blocks in the  
 196 ResNet model. The dataset is distributed among the clients following one of the two cases: 1) IID  
 197 case: the samples are uniformly distributed to the clients; 2) Non-IID case: 50% of the samples on  
 198 each client belong to two classes while the other 50% of the samples are uniformly picked from the  
 199 rest of the classes. The detailed descriptions of the experiments are included in Appendix B.

### 200 4.1 Numerical Results

201 In the IID case, we compare the performance of the models obtained with FLM algorithm on  $K = 8$   
 202 clients with centralized training. The result is shown in Figure 1a. In this setting, the client models  
 203 use ResNet[5, 5, 5] and the matched server model has layer number [5, 5, 6]. We can see that the  
 204 performance of the trained client models and that of the matched server model are comparable to the  
 205 centralized model, but the required number of local update is much less than centralized training.

206 In the Non-IID case, we compare the performance of the models obtained with FLM algorithm and  
 207 FedAvg algorithm with  $K = 5$  clients. The result is shown in Figure 1b. The performance of  
 208 the server model, as well as that of the client models trained with FLM under limited number of  
 209 communication round are much better than the model trained with FedAvg. The matched server  
 210 model has layer [6, 7, 8], and the client models are ResNet[5, 5, 5]. In the experiment, FLM can  
 211 better utilize the heterogeneous layers from different clients to build 40% deeper server model to  
 212 deal with data heterogeneity, while FedAvg only averages the client models.

213 **References**

- 214 [1] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan,  
 215 P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances*  
 216 *in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- 217 [2] J. Xu, Y. Pan, X. Pan, S. Hoi, Z. Yi, and Z. Xu, “Regnet: self-regulated network for image  
 218 classification,” *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- 219 [3] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. De-  
 220 hghani, M. Minderer, G. Heigold, S. Gelly *et al.*, “An image is worth 16x16 words: Trans-  
 221 formers for image recognition at scale,” *arXiv preprint arXiv:2010.11929*, 2020.
- 222 [4] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Pro-*  
 223 *ceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–  
 224 778.
- 225 [5] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional  
 226 transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- 227 [6] H. Wang, M. Yurochkin, Y. Sun, D. Papailiopoulos, and Y. Khazaeni, “Federated learning with  
 228 matched averaging,” *arXiv preprint arXiv:2002.06440*, 2020.
- 229 [7] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingerman, V. Ivanov, C. Kiddon,  
 230 J. Konečný, S. Mazzocchi, B. McMahan *et al.*, “Towards federated learning at scale: System  
 231 design,” *Proceedings of Machine Learning and Systems*, vol. 1, pp. 374–388, 2019.
- 232 [8] H.-Y. Chen and W.-L. Chao, “Fedbe: Making bayesian model ensemble applicable to federated  
 233 learning,” *arXiv preprint arXiv:2009.01974*, 2020.
- 234 [9] R. Liu, F. Wu, C. Wu, Y. Wang, L. Lyu, H. Chen, and X. Xie, “No one left behind: Inclusive  
 235 federated learning over heterogeneous devices,” *arXiv preprint arXiv:2202.08036*, 2022.
- 236 [10] A. Khaled, K. Mishchenko, and P. Richtárik, “First analysis of local GD on heterogeneous  
 237 data,” *arXiv preprint arXiv:1909.04715*, 2019.
- 238 [11] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, “On the convergence of fedavg on non-iid  
 239 data,” in *International Conference on Learning Representations*, 2019.
- 240 [12] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, “Federated optimization  
 241 in heterogeneous networks,” *arXiv preprint arXiv:1812.06127*, 2018.
- 242 [13] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, “Scaffold:  
 243 Stochastic controlled averaging for federated learning,” in *International Conference on Ma-*  
 244 *chine Learning*. PMLR, 2020, pp. 5132–5143.
- 245 [14] X. Zhang, M. Hong, S. Dhople, W. Yin, and Y. Liu, “Fedpd: A federated learning framework  
 246 with optimal rates and adaptivity to non-iid data,” *arXiv preprint arXiv:2005.11418*, 2020.
- 247 [15] T. Lin, L. Kong, S. U. Stich, and M. Jaggi, “Ensemble distillation for robust model fusion in  
 248 federated learning,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 2351–  
 249 2363, 2020.
- 250 [16] N. Shi, F. Lai, R. A. Kontar, and M. Chowdhury, “Fed-ensemble: Improving generalization  
 251 through model ensembling in federated learning,” *arXiv preprint arXiv:2107.10663*, 2021.
- 252 [17] D. H. Wolpert, “Stacked generalization,” *Neural networks*, vol. 5, no. 2, pp. 241–259, 1992.
- 253 [18] R. Polikar, “Ensemble learning,” in *Ensemble machine learning*. Springer, 2012, pp. 1–34.
- 254 [19] B. Hanin, “Universal function approximation by deep neural nets with bounded width and relu  
 255 activations,” *Mathematics*, vol. 7, no. 10, p. 992, 2019.
- 256 [20] T. Nguyen, M. Raghu, and S. Kornblith, “Do wide and deep networks learn the same things?  
 257 uncovering how neural network representations vary with width and depth,” in *International*  
 258 *Conference on Learning Representations*, 2020.
- 259 [21] X. Zhang, W. Yin, M. Hong, and T. Chen, “Hybrid federated learning: Algorithms and imple-  
 260 mentation,” *arXiv preprint arXiv:2012.12420*, 2020.
- 261 [22] Y. Idelbayev, “Proper ResNet implementation for CIFAR10/CIFAR100 in PyTorch,” [https://github.com/akamaster/pytorch\\_resnet\\_cifar10](https://github.com/akamaster/pytorch_resnet_cifar10),  
 262 accessed: 20xx-xx-xx.
- 263 [23] Y. Wu and K. He, “Group normalization,” in *Proceedings of the European conference on com-*  
 264 *puter vision (ECCV)*, 2018, pp. 3–19.

265 **A Algorithm Details**

266 In this section, we provide the detailed explanation of step 2 and 3 in Algorithm 1.

For clarity of the notations, we put a list of notations with their meanings in Table 1.

Notation	Notation meaning
$(\cdot)_0$	Server index
$(\cdot)_k$	Client index
$K$	Total client number
$(\cdot)^t$	Iteration index
$T$	Total training iteration
$(\cdot)^P$	Matching index
$P$	Total layer matching iteration
$R$	Layer matching frequency
$(\cdot)_m$	Block index
$M$	Total block number
$(\cdot)[l]$	Layer index
$L, \bar{L}$	Total layer number, maximum layer number
$\Theta$	Model parameter
$\Pi$	Matching pattern

Table 1: Notation Table

267

268 **Step 2:** The server-side layer matching is the key step for the proposed FLM algorithm. At the iter-  
 269 ation to perform layer matching  $t \bmod R = 0$ , for each block  $m \in [M]$  in parallel, we sequentially  
 270 match each client layers to the server layers for  $P$  iterations. The matching procedure is illustrated  
 271 in Algorithm 2, which consist four major stages discussed below.

272 **Stage 2-1:** After randomly choosing a client  $k$  to perform layer matching, the first stage is to con-  
 273 struct the server model  $\Theta_{0,m}^{t,p}$  with the layers of the other clients and their matching patterns, i.e.,

$$\Theta_{0,m}^{t,p} = \frac{\sum_{k' \neq k} \Pi_{k',m}^{t,p} \Theta_{k',m}^t}{\sum_{k' \neq k} \Pi_{k',m}^{t,p} \mathbb{1}_{L_{k',m}}},$$

274 which removes the impact of the parameters of client  $k$ . Therefore, the layer similarities between the  
 275 client and server models computed in the following matching stages is not affected by the previous  
 276 matching pattern  $\Pi_{k,m}^{t-1}$  of client  $m$ .

277 **Stage 2-2:** In this stage, we compute the cost for each client layer  $l_k = 1, \dots, L_{k,m}$  in the block to  
 278 a layer in the server. Assume that we set the maximum number of layers in block  $m$  in the server  
 279 model as  $\bar{L}_{0,m}$ . Then the costs for matching all layers in block  $m$  for client  $M$  form a cost matrix  
 280  $\mathbf{C}_{k,m}^{t,p} \in \mathbb{R}^{\bar{L}_{0,m} \times L_{k,m}}$ , where its entries are defined as:

$$\mathbf{C}_{k,m}^{t,p}[l_0, l_k] = c_1(\Theta_{k,m}^t[l_k], \Theta_{0,m}^{t,p}[l_0]) + c_2(\Theta_{k,m}^t[l_k], l_0, \bar{L}_{0,m}), \quad (6)$$

281 where  $c_1(\cdot)$  is the fidelity score to match a client layer to an existing server layer and  $c_2(\cdot)$  is the  
 282 penalty to increase the number of the server layer. A reasonable choice for  $c_1(\cdot)$  is a distance metric  
 283 (e.g., Euclidean distance or cosine angle) between the server and client layers, and the choice for  
 284  $c_2(\cdot)$  can be model selection metric (e.g., BIC) on the model size.

285 **Stage 2-3** With the cost matrix  $\mathbf{C}_{k,m}^{t,p}$  that describe the cost to match each client layer to the server  
 286 layers, we can optimize the corresponding linear assignment problem:

$$\min_{\Pi_{k,m}^{t,p+1}} \text{Trace} \left( (\Pi_{k,m}^{t,p+1})^\top \mathbf{C}_{k,m}^{t,p} \right), \quad \text{s.t. } \Pi_{k,m}^{t,p+1} \in \{0, 1\}^{\bar{L}_{0,m} \times L_{k,m}}, (\Pi_{k,m}^{t,p+1})^\top \mathbb{1}_{\bar{L}_{0,m}} = \mathbb{1}_{L_{k,m}}, \quad (7)$$

287 and obtain the updated matching pattern  $\Pi_{m,k}^{t,p+1}$ . If any of the rows  $L_{0,m} + 1, \dots, \bar{L}_{0,m}$  in the  
 288 matching pattern are not all-zero, then some of the client layers are added to the server model as  
 289 new layers, and the server model becomes deeper. Otherwise, the server model keeps the same  
 290 depth  $L_{0,m}$  or even gets shallower.



291 **Stage 2-4:** After updating the matching pattern from  $\Pi_{m,k}^{t,p}$  to  $\Pi_{k,m}^{t,p+1}$ , the server layer order and the  
 292 total layer number may change. So in this stage, we need to update the matching patterns of the  
 293 other clients  $\{\Pi_{k',m}^{t,p+1}\}_{k' \neq k}$  accordingly. First, we compute the ‘‘average layer index’’ for each server  
 294 layer  $l_0^{t,p}$ :

$$\bar{l}_0 = \frac{1}{|\{l_k | \Pi_{k,m}^t[l_0, l_k] = 1\}|} \sum_{l_k \in \{l_k | \Pi_{k,m}^t[l_0, l_k] = 1\}} l_k.$$

295 Then, we reorder the server layers based on the order of the average layer index  $\bar{l}_0$  and obtain the  
 296 new server layer indices  $l_0^{t,p+1}$ . Finally, the layer matching patterns for the clients are updated with  
 297  $\Pi_{k,m}^{t,p+1}[l_0^{t,p+1}, l_k] = \Pi_{k,m}^{t,p}[l_0^{t,p}, l_k]$ .

298 By repeating Stage 2-1 to 2-4 for  $P$  iterations ( $P \geq K$ ), all the client matching patterns  $\{\Pi_k^t\}_{k=1}^K$   
 299 will be updated and that finishes the whole procedure of Step 2 in Algorithm 1.

300 **Step 3:** In this step, the server construct the server model with all the updated client matching  
 301 patterns  $\{\Pi_k^t\}_{k=1}^K$  and the client models  $\{\Theta_k^t\}_{k=1}^K$  as

$$\Theta_0^t = \frac{\sum_{k=1}^K \Pi_k^t \Theta_k^t}{\sum_{k=1}^K \Pi_k^t \mathbb{1}_{L_k}},$$

302 which averages the client layers that match to the same server layer, i.e.,

$$\Theta_{0,m}^t[l_0] = \frac{1}{|\{l_k | \Pi_{k,m}^t[l_0, l_k] = 1\}|} \sum_{l_k \in \{l_k | \Pi_{k,m}^t[l_0, l_k] = 1\}} \Theta_{k,m}^t[l_k].$$

303 With the server model, the server then reconstruct the client models with  $\Theta_k^t = \Pi_k^{t \top} \Theta_0^t$  that serves  
 304 as the initial client model and the regularizer for the local training at the next iteration  $t + 1$ .

## 305 B Experiment Details

306 **Model:** The structure of the model is summarized in Table 2. Each residual block is constructed with  
 307 two Conv2d layers, each followed by one GroupNorm layer [23] and a ReLU activation layer, and  
 308 finally add an identical mapping to the end. In the experiment, we use ResNet[ $a, b, c$ ] to denote the  
 309 ResNet model using certain number of residual blocks of block 2, 3 and 4 in the table. During model  
 310 aggregation phase, we perform layer matching to the parameters in block 2, 3B and 4B separately,  
 and average the parameters in the other layers.

Block #	Shape	Parameter #
1	Conv2d( $3 \times 3 \times 16$ )	448
	BatchNorm2d(16)	–
2	ResidualLayer(16, 16)	4,608
3A	ResidualLayer(16, 32)	13,824
3B	ResidualLayer(32, 32)	18,432
4A	ResidualLayer(32, 64)	55,296
4B	ResidualLayer(64, 64)	73,728
5	AvgPool()	–
	Linear( $64 \times 10$ )	650

Table 2: ResNet model structure for Cifar-10 dataset.

311

312 **Dataset:** In the experiment, we use Cifar-10 dataset. There are two data distribution settings: 1)  
 313 IID setting: in this setting the samples are uniformly distributed to clients; 2) Non-IID setting: in  
 314 this setting, 50% of the samples on each client belongs to two classes while the other 50% of the  
 315 samples are uniformly picked from the rest of the classes.