

Supplementary Material: Data-Free Diversity-Based Ensemble Selection For One-Shot Federated Learning

1 Execution process of DeDES

Fig. 1 gives an case study of DeDES when we use the last layer of model parameters as the model representation and *Kernel-PCA* as the dimension reduction method. Under this case, we assume that all models are well-trained and no low-quality model(s) so that we don't need to do the model filter; the choice of clustering method will depend on the data partition which are shown in section 3.

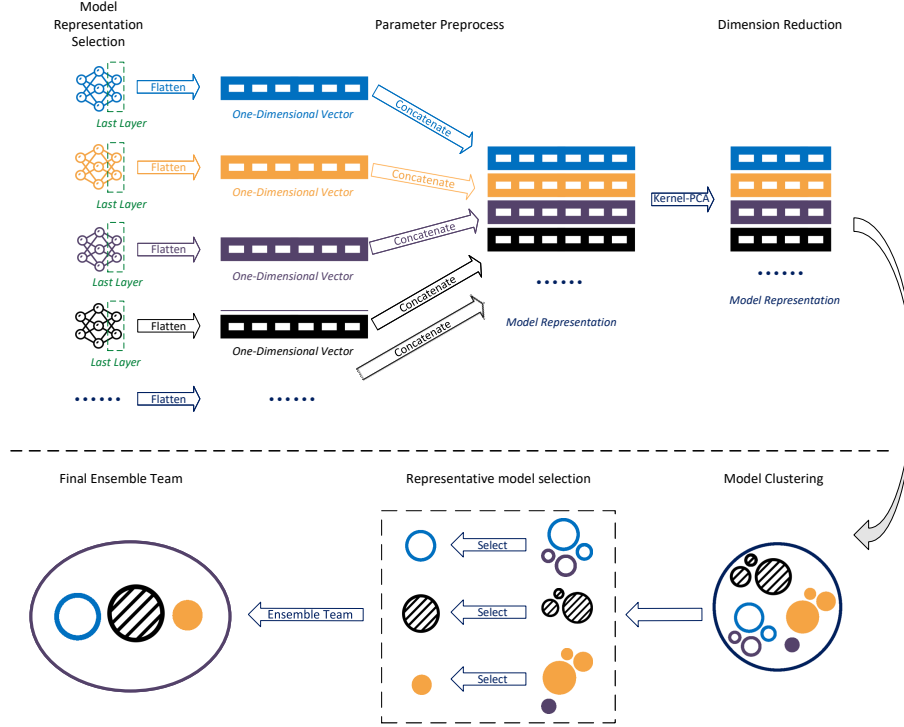


Fig. 1: A case study of the execution process of DeDES framework, where Kernel-PCA is used as the dimension reduction method and no model filtering is needed.

Algorithm 1: Model filter `OutlierFilter` method

Input: model set \mathcal{M} , truncated threshold pair (p_{low}, p_{high}) , interval scale s ,
and model scores $\mathcal{S} = \{s_i\}_{i=1}^m$.
Output: Outlier model set \mathcal{O} .
 \triangleright *Sort the score set \mathcal{S} , such as local validation accuracy, by ascending order.*
1 $\mathcal{S} \leftarrow \text{AscendingSort}(\mathcal{S})$
 \triangleright *Get the value of the p_{low} -th, p_{high} -th element of \mathcal{S} .*
2 $q_{low}, q_{high} \leftarrow \mathcal{S}_{p_{low}}, \mathcal{S}_{p_{high}}$
3 $interval \leftarrow q_{high} - q_{low}$
4 $outlier_threshold \leftarrow q_{low} - s * interval$
5 $\mathcal{O} \leftarrow \emptyset$
6 **for** $i = 1$ **to** m **do**
7 **if** $s_i < outlier_threshold$ **then**
8 $\mathcal{O} \leftarrow \mathcal{O} \cup \{M_i\}$
9 **return** \mathcal{O}

2 MODEL FILTER ALGORITHM

As in Alg. 1 from the *main paper*, we use the `OutlierFilter` to obtain the outlier models \mathcal{O} based on the model scores \mathcal{S} provided from each party, which can be the local validation accuracy or prediction confidence. `OutlierFilter` can be any score-based unsupervised outlier detection methods, we apply a variation of the commonly-used box-plot in our experiment, which is show in Alg. 1 of this *supplementary material*.

3 Experiment Setup

3.1 Data partition method explanation

Examples of the four types of data partitions are shown below.

1. Homogeneous (*homo*): all parties have the same amount of data and the same distribution of data, i.e., this data distribution is independent and identical distribution (IID).

E.g, if we divide the MNIST dataset which has 280,000 images into 100 parties with *homo* setting, then all parties will have the same amount of data (all round $280,000 / 100 = 2,800$) and the same label distribution of data (every party will have around $2,800 / 10 = 280$ images for every label from 0 to 9).

2. IID but with different quantity (*iid-dq*): all parties have the same distribution of data but with different data amount. E.g., every party will have similar amount of images of 0 to 9, but the amount of total images is different for each party.

3. Non-IID label distribution skew (*noniid-lds*): the label distributions $P(y_i)$ vary across parties. Each party will be allocated a proportion of the samples of each label according to *Dirichlet Distribution*. That is, the label distribution

of each party is skewed and data amount among parties is different with each other.

4. Extreme Non-IID setting with only k labels of each party (*noniid-lk*): Under this partition setting, each party will only have k labels of data. For example, if we set $k = 3$, then each party will only have 3 labels of data, e.g., 0, 3 and 7.

Figure 2 shows the distribution of data amount and label distribution of each party under these four different data partition settings.

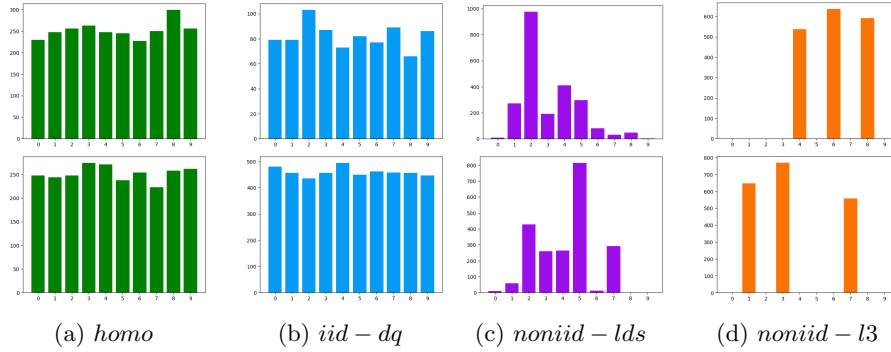


Fig. 2: Example distribution of four dataset partition strategies for the *EMNIST digits* dataset with party number of 100. For each of subfigures (a), (b), (c), and (d), we randomly picked two parties from a total of 100 parties and depicted their distributions of data quantity and label distribution.

3.2 Component Configurations

In this subsection, we will describe the default configurations of our DeDES framework for the experiments in the main paper.

For local model training, we utilize the *SGD* optimizer to train model for every party for 200 epochs with learning rate started at 0.1 and decrease at later epochs. After the training finished, we select the model with the highest local validation accuracy of each party to upload to the server of model market.

In our experiments, we select the final model layer as the model representation in our experiments; we utilize the *MINMAX* scaler to preprocess the *Model Representation Matrix*, and the *Gaussian Normalization* scaler to preprocess the *Label Distribution* ground-truth input data; we do not utilize any dimension reduction strategies for the model representation matrix because the final layer of model parameters are already few in number.

Spectral clustering is used for the *homo* and *iid-dq* distribution partition methods; *K-Means* clustering is used for the *noniid-ls* and *noniid-lk* partition methods.

Our proposed model representative selection approach is utilized to determine the representative model within each cluster; as we have frequently stated, we use plurality voting to perform ensemble learning; we use the test accuracy as the evaluation metric for all the methods.

To compare the efficiency of different methods, we also recorded the running time of all methods, details will be shown in the experiment section. More details about the experimental settings, such as the devices we use or examples about four data partitions, are mentioned in the *supplementary material*.

3.3 Environment

All our experiments are running on a single machine with 1TB RAM and 256 cores AMD EPYC 7742 64-Core Processor @ 3.4GHz CPU. The GPU we use is NVIDIA A100 SXM4 with 40GB memory. The OS we use is Ubuntu 20.04.4 LTS.

Our software environment settings are: Python 3.9.12, PyTorch 1.12.1 with CUDA 11.6.

All experiments were tested 3 times, and the average results are presented in the experiment section.

4 ADDITIONAL EXPERIMENTS

4.1 The importance of ensemble learning

Fig. 3 shows the data distributions of the EMNIST balanced dataset with local validation accuracy and test accuracy for the whole test set when $m=10$. We can see that the capacity of a single model is weak, even their local validation accuracy (validation accuracy for their own dataset) is high, therefore validates the importance of ensemble learning which can utilize the collaborative power of multiple models.

4.2 Impact on Efficiency

Table. 1 shows that in some cases, DeDES is the second best method after *All Selection (AS)*. Note that the efficiency of *AS* is quite poor, and the performance gap between these two approaches is small, validating that our method can reduce ensemble time to a large extent with minimal performance loss.

It is easy to know that the inference time for ensemble learning (weighted voting) increases linearly with K , i.e., the total inference time T for one test sample is $K \times c$, where c is constant inference time for one sample by one model. The experimental results depicted in Fig.4 indicate that when K reaches a certain value, the test accuracy will not increase significantly, sometimes even decrease. Therefore, with a suitable K (usually 50% of m), we can substantially reduce our inference time for ensemble learning while achieving good ensemble performance. And we do not need to concern too much about the running duration of DeDES compared to others because the ensemble selection process will only run once and will finish in a few minutes, therefore it is of little consequence.

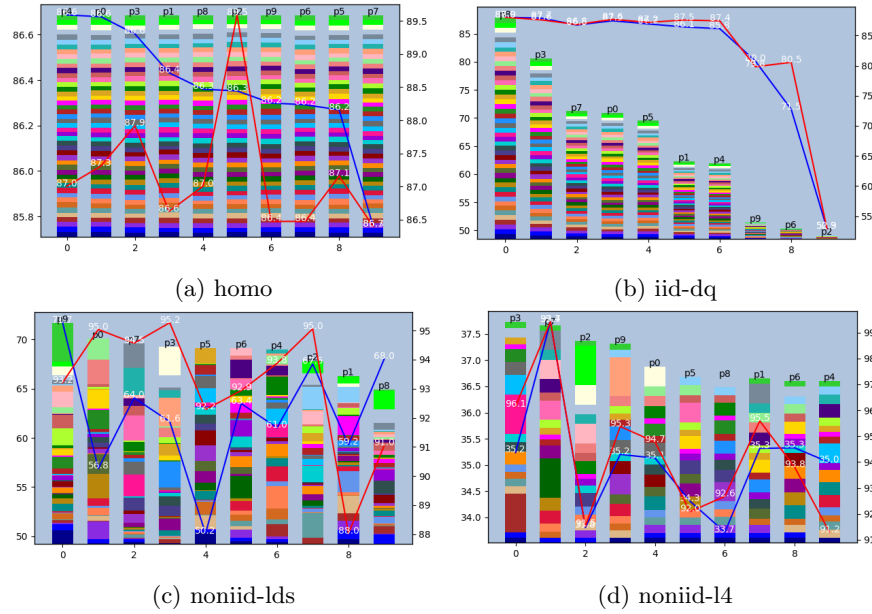


Fig. 3: Different data distributions of the EMNIST balanced dataset with local validation accuracy and test accuracy for the whole test set when $m=10$. Every dot in the red line shows the best local validation accuracy after training for 200 epochs while the blue line shows the whole test accuracy of the k -th party p_k .

4.3 Ablation Studies

- **Performance Comparison on different model structures and datasets**
Our method is solid for various model structures and datasets.
- **Performance Comparison on different model representation** It is better to use the models' later layer's parameters for representation than utilizing their front layer's parameters.
- **Importance of Dimension Reduction Methods.** *Kernel-PCA* is better than other dimension reduction methods such as *PCA* and *non-compression*.
- **Clustering/Diversity validation** Our method can really cluster similar models together and the whole team's diversity is higher than other methods.

Table 1: Test accuracy (%) comparison for different dataset on different data partitions and model structures. The best and next best methods are **bolded** and underlined, respectively. If our DeDES method is better than the *LD* ground-truth method, the value of *LD* method will be marked in skyblue.

Dataset	Partition	m	K	DeDES	AS	CV	DS	RS	FedAvg	MeanAvg	LD	Oracle
EMNIST Digits (VGG-5 Spinal FC)	homo	400	150	98.03	<u>98.10</u>	98.10	98.08	98.07	10.28	10.26	98.10	99.74
	iid-dq	400	150	99.27	98.75	<u>98.93</u>	98.88	98.72	10.51	10.48	99.27	99.71
	noniid-l1	400	150	97.67	<u>96.99</u>	95.47	91.70	96.67	10.01	9.89	<u>92.86</u>	99.72
	noniid-l3	400	150	98.21	<u>97.96</u>	97.87	63.59	94.35	10.11	10.09	<u>98.13</u>	99.61
EMNIST Letters (VGG-5 Spinal FC)	homo	200	120	88.64	88.77	88.88	88.82	88.68	3.72	3.71	88.77	95.12
	iid-dq	200	120	<u>92.32</u>	92.19	91.97	92.33	92.13	3.84	3.82	92.33	95.12
	noniid-l1	200	120	87.93	<u>87.74</u>	86.52	83.45	87.45	4.03	4.02	<u>85.01</u>	94.90
	noniid-l8	200	120	89.10	<u>87.93</u>	84.40	86.98	85.95	3.85	3.84	<u>87.54</u>	95.06
EMNIST Balanced (VGG-5 Spinal FC)	homo	100	50	85.19	84.94	<u>85.10</u>	84.96	84.96	2.10	2.11	<u>84.83</u>	89.70
	iid-dq	100	50	<u>87.34</u>	87.28	87.31	87.35	86.90	2.04	2.04	87.35	89.25
	noniid-l1	100	50	83.43	82.72	78.65	79.44	81.89	2.19	2.16	<u>77.28</u>	89.48
	noniid-l18	100	50	85.43	<u>82.99</u>	81.22	81.02	81.93	2.09	2.08	<u>82.87</u>	89.52
CIFAR10 (Resnet-50)	homo	200	100	<u>32.08</u>	32.09	32.07	30.78	30.30	10.18	9.69	32.08	88.68
	iid-dq	200	100	36.97	38.49	<u>38.84</u>	39.03	36.66	10.04	10.03	38.81	88.10
	noniid-l1	200	100	29.71	<u>29.23</u>	26.02	29.10	26.67	9.89	9.88	<u>28.94</u>	87.31
	noniid-l4	200	100	34.40	<u>33.50</u>	32.24	30.00	33.05	10.02	9.87	<u>34.15</u>	89.67
CIFAR100 (Resnet-50)	homo	20	12	<u>20.84</u>	22.84	20.58	20.65	20.48	0.99	0.99	<u>20.85</u>	59.81
	iid-dq	20	12	<u>47.38</u>	47.37	<u>47.38</u>	<u>47.38</u>	25.10	1.00	0.94	47.38	60.35
	noniid-l1	20	12	<u>16.31</u>	18.71	15.97	16.15	15.78	0.96	0.97	<u>15.32</u>	60.38
	noniid-l45	20	12	<u>21.29</u>	23.68	20.56	20.26	19.97	0.92	0.91	<u>19.61</u>	61.74

Method	Rank	Accuracy (%)
DeDES	34/1024	84.34
AS	114/1024	83.39
CV	241/1024	82.29
DS	348/1024	80.86
LD	608/1024	77.77
RS	669/1024	76.54

Table 2: Complete inspection on ensemble teams for *EMNIST Balanced* dataset with $m = 10$, *noniid-l1s* partition.

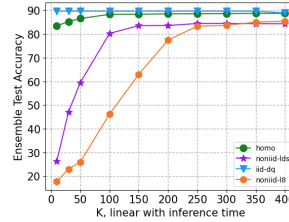


Fig. 4: The relationship of K and Ensemble Test Accuracy of DeDES for the *EMNIST Letters* Dataset when $m=400$.