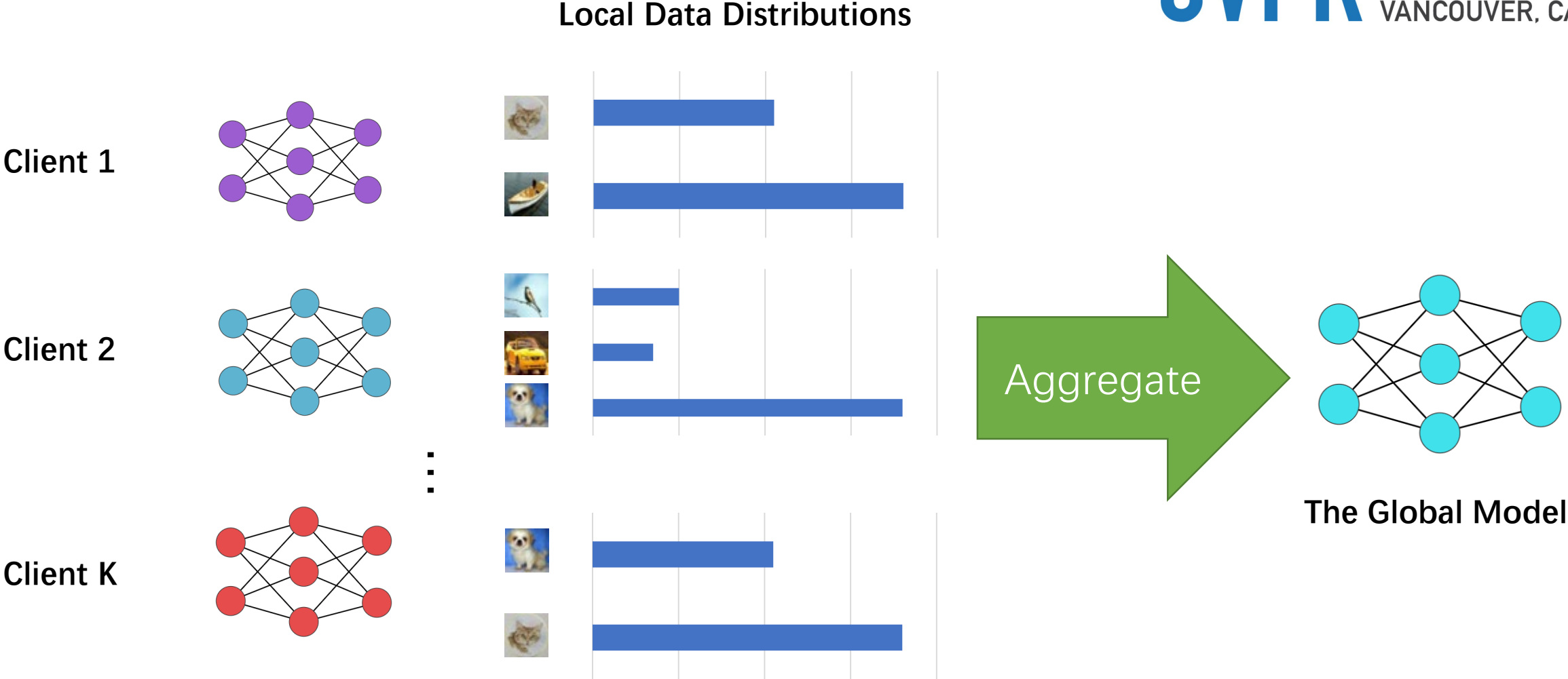


# Federated Learning with Data-Agnostic Distribution Fusion

Jian-hui Duan, Wenzhong Li\*, Derun Zou, Ruichen Li, Sanglu Lu  
State Key Laboratory for Novel Software Technology, Nanjing University  
Nanjing, China

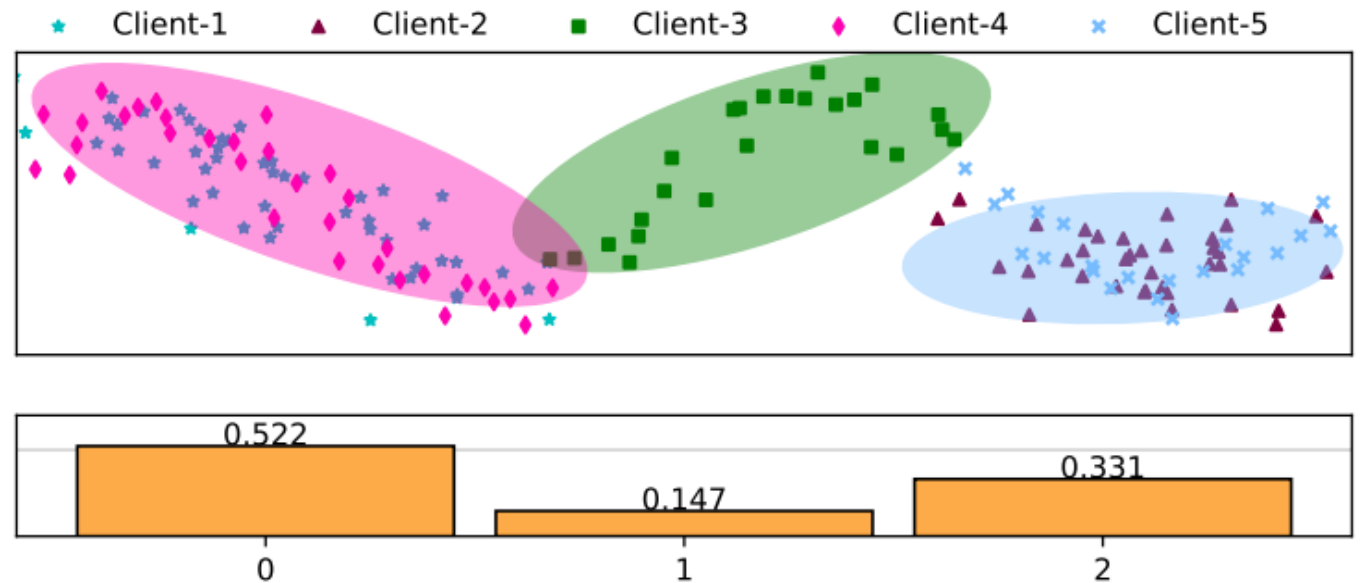
# Introduction



Data distribution heterogeneity in federated learning usually cause accuracy drop in global model

# Introduction

- Data distribution of five clients can be represented by a distribution fusion model with three virtual components.
- Client models can be aggregated based on components to better approach centralized training.



# Problem Definition

JUNE 18-22, 2023

CVPR



VANCOUVER, CANADA

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \sum_{k=1}^K \frac{n_k}{n} \mathcal{L}_k(\mathbf{w})$$

First step is to modify the optimizing target:

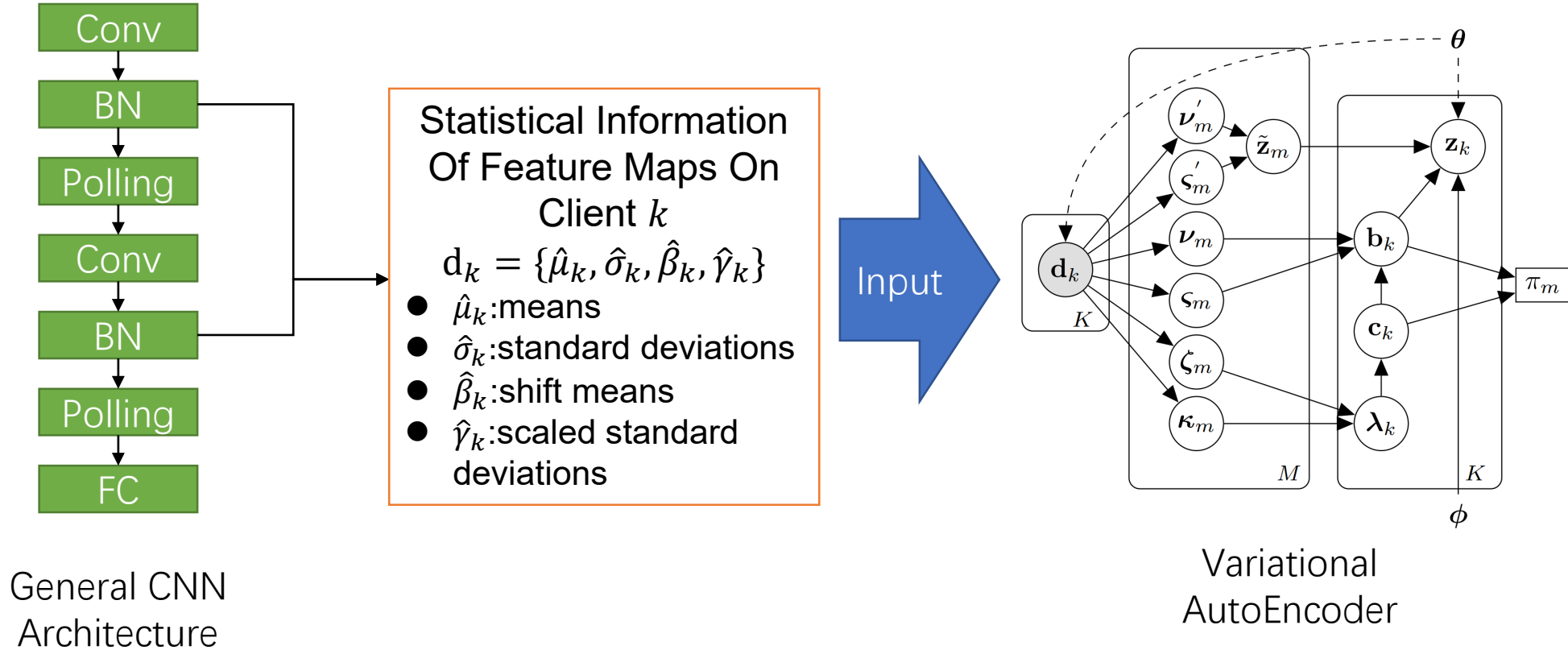
reallocate local models with weight  $b_{km}$  to  $m$  component, then aggregate component with weight  $\pi_m$ .



$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \underbrace{\sum_{m=1}^M \pi_m}_{\text{Virtual Component Fusion Weight}} \underbrace{\sum_{k=1}^K b_{km} \mathcal{L}_k(\mathbf{w})}_{\text{Virtual Component Allocation Weight}}$$

Virtual Component Fusion Weight  
Virtual Component Allocation Weight

# Variational AutoEncoder



Optimal aggregation weights require accurate local data distributions, we gather parameters of normalization layers, and use Variational AutoEncoder to infer local distribution parameters, in order to construct local data distributions.

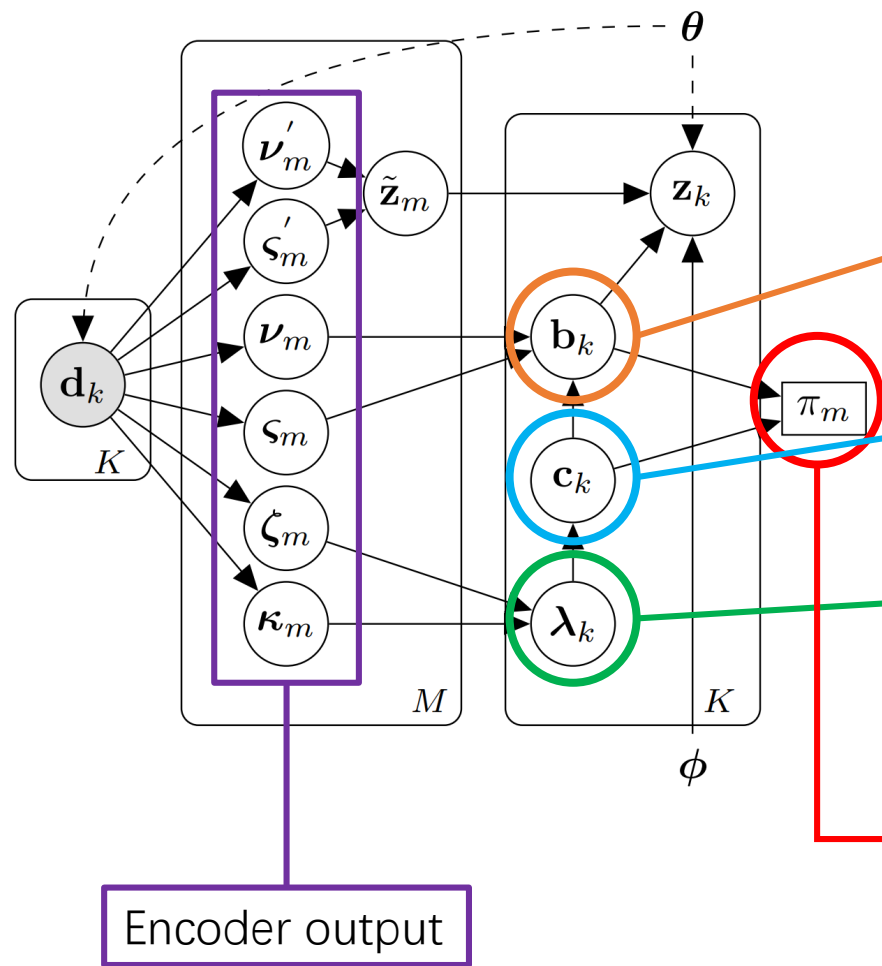
# Variational AutoEncoder

JUNE 18-22, 2023

CVPR



VANCOUVER, CANADA



$$\mathbf{b}_k \sim \mathcal{N}(\boldsymbol{\nu}_m, \boldsymbol{\Sigma}_m)$$

$$\mathbf{c}_k \sim \text{Bernoulli}\left(\prod_{m=1}^M \lambda_{km}\right)$$

$$\lambda_k \stackrel{i.i.d.}{\sim} \text{Beta}(\boldsymbol{\zeta}_m, \boldsymbol{\kappa}_m)$$

$$\pi_m = \frac{\exp\left(\frac{1}{K} \sum_{k=1}^K q_\phi(c_{km}) \cdot b_{km}\right)}{\sum_{m=1}^M \exp\left(\frac{1}{K} \sum_{k=1}^K q_\phi(c_{km}) \cdot b_{km}\right)}$$

# Optimizing Procedure

To better optimize distribution parameters, we design following sampling methods:

$$\lambda_k \stackrel{i.i.d.}{\sim} \text{Beta}(\zeta_m, \kappa_m) \quad \xrightarrow{\text{Sampling}} \quad \lambda_k \sim (1 - \xi^{\frac{1}{\kappa_k}})^{\frac{1}{\zeta_k}}, \text{ where } \xi \sim \text{Uniform}(0, 1)$$

$$\mathbf{c}_k \sim \text{Bernoulli}\left(\prod_{m=1}^M \lambda_{km}\right) \quad \xrightarrow{\text{Sampling}} \quad \mathbf{c}_{km} = \arg \max_i (g_i + \log \prod_{i=1}^2 \lambda_{ki}) \text{ Where } g_i \sim \text{Gumbel}(0, 1)$$

# Experiment

## Datasets:

Datasets	Data Type	Train	Test	Total
MNIST[25]	1 channel image	60,000	10,000	70,000
Fashion-MNIST[44]	1 channel image	60,000	10,000	70,000
CIFAR-10[22]	3 channel image	50,000	10,000	60,000
Sentiment140[8]	Text data	-	-	1,600,000

## BackBone Models:

ResNet18[9], DenseNet121[11], MobileNetV2[36], LeNet[24], BiLSTM[]

## Benchmarks:

Single-model: FedAvg[30], FedProx[26], Fed-GN[10], FedMA[43]

Multi-model: FeSEM[46], IFCA[7], FedCluster[2], FedGroup[6]



# Experiment

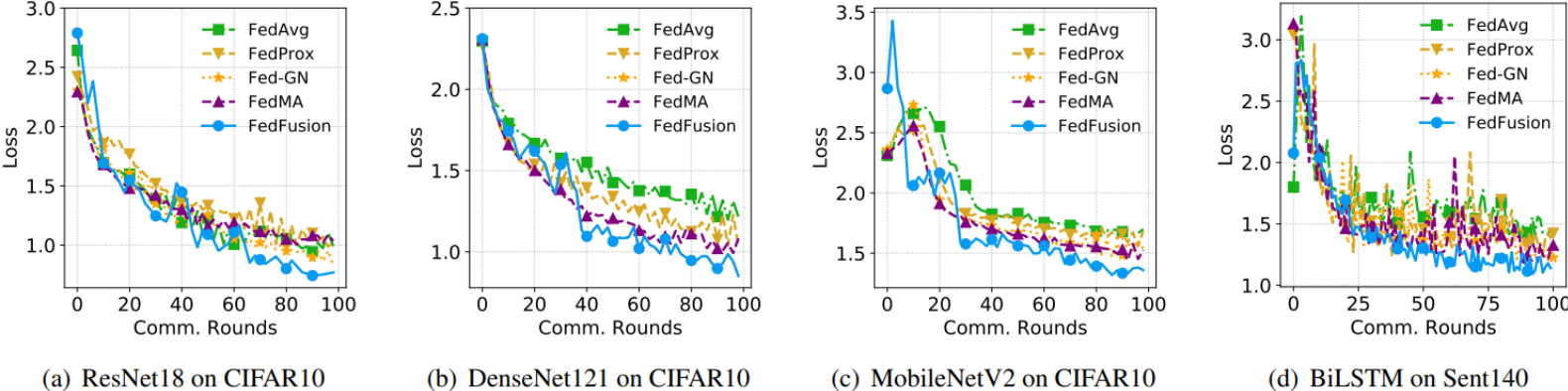


Figure 3. Training loss of different algorithms.

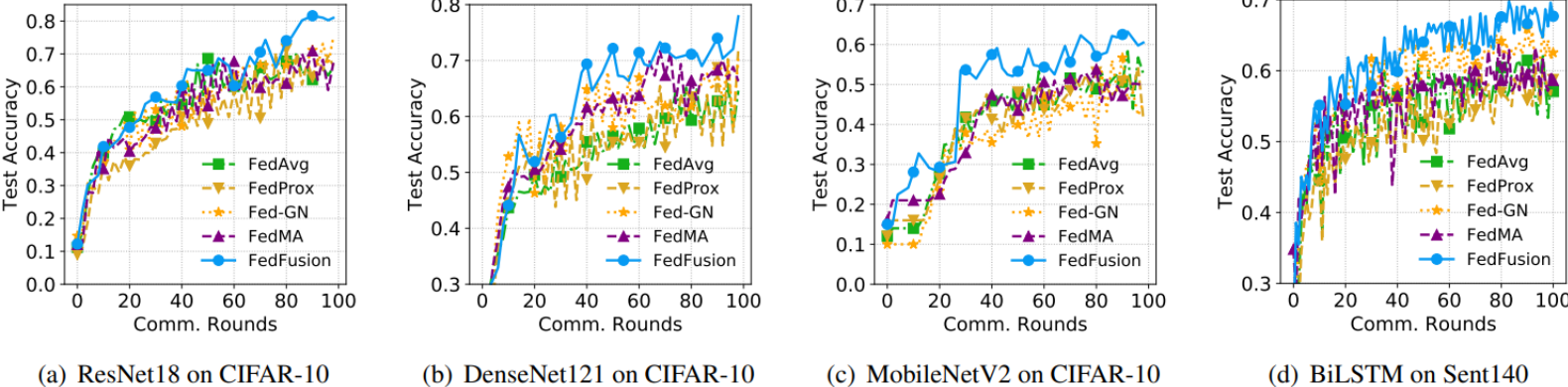


Figure 4. Training efficiency of different algorithms.

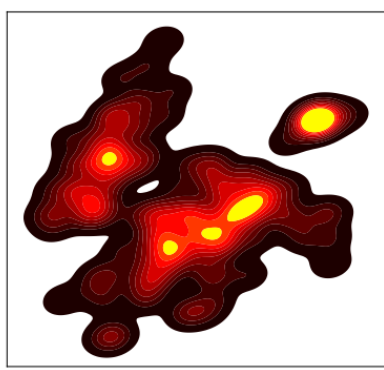
FedFusion(blue line) shows the lowest loss, and converges the fastest among all evaluated algorithms

# Experiment

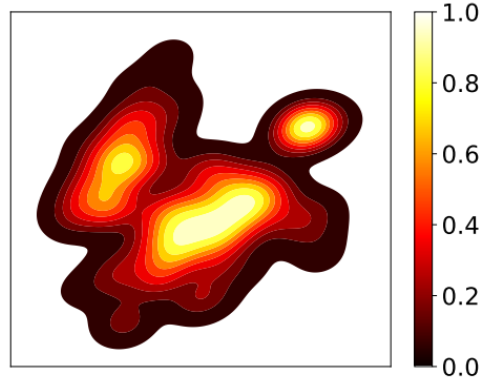
	Dataset	CIFAR-10			FMNIST	MNIST	Sent140
	Model	ResNet18	DenseNet121	MobileNetV2	LeNet	LeNet	BiLSTM
Single-model	FedAvg	68.78 ( $\pm 0.89$ )	63.33 ( $\pm 0.67$ )	54.69 ( $\pm 3.92$ )	79.20 ( $\pm 1.15$ )	97.32 ( $\pm 0.04$ )	58.33 ( $\pm 2.03$ )
	FedProx	70.18 ( $\pm 0.45$ )	66.85 ( $\pm 0.93$ )	55.03 ( $\pm 2.77$ )	80.03 ( $\pm 0.98$ )	97.55 ( $\pm 0.02$ )	59.73 ( $\pm 1.38$ )
	Fed-GN	72.57 ( $\pm 0.78$ )	70.02 ( $\pm 1.36$ )	56.43 ( $\pm 1.92$ )	81.11 ( $\pm 0.74$ )	97.88 ( $\pm 0.02$ )	63.41 ( $\pm 1.94$ )
	FedMA	73.43 ( $\pm 1.03$ )	70.13 ( $\pm 1.71$ )	59.61 ( $\pm 2.01$ )	81.02 ( $\pm 1.35$ )	98.06 ( $\pm 0.03$ )	60.86 ( $\pm 2.42$ )
Multi-model	FeSEM	67.78 ( $\pm 2.58$ )	62.65 ( $\pm 0.82$ )	53.82 ( $\pm 3.69$ )	78.18 ( $\pm 1.45$ )	96.24 ( $\pm 0.17$ )	59.57 ( $\pm 3.41$ )
	IFCA	73.04 ( $\pm 1.45$ )	70.85 ( $\pm 2.03$ )	58.93 ( $\pm 2.45$ )	80.82 ( $\pm 1.29$ )	97.09 ( $\pm 0.11$ )	60.82 ( $\pm 2.74$ )
	FedCluster	72.57 ( $\pm 0.78$ )	68.77 ( $\pm 1.38$ )	58.18 ( $\pm 1.22$ )	79.11 ( $\pm 0.74$ )	97.88 ( $\pm 0.02$ )	63.41 ( $\pm 1.94$ )
	FedGroup	74.38 ( $\pm 1.92$ )	71.63 ( $\pm 0.74$ )	59.86 ( $\pm 2.09$ )	81.32 ( $\pm 2.07$ )	97.37 ( $\pm 0.61$ )	63.61 ( $\pm 3.26$ )
	FedFusion	<b>81.26</b> ( $\pm 0.82$ )	<b>75.92</b> ( $\pm 1.25$ )	<b>62.88</b> ( $\pm 1.21$ )	<b>83.16</b> ( $\pm 0.74$ )	<b>98.49</b> ( $\pm 0.04$ )	<b>67.51</b> ( $\pm 1.71$ )

Comparison of average test accuracy on non-IID datasets

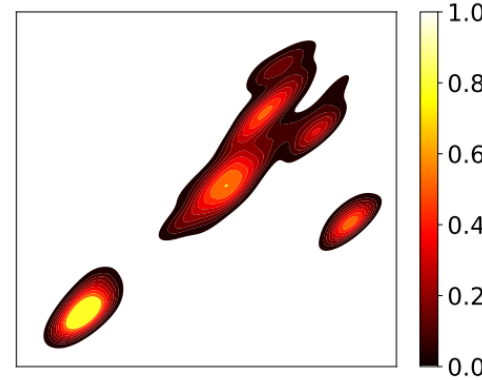
# Experiment



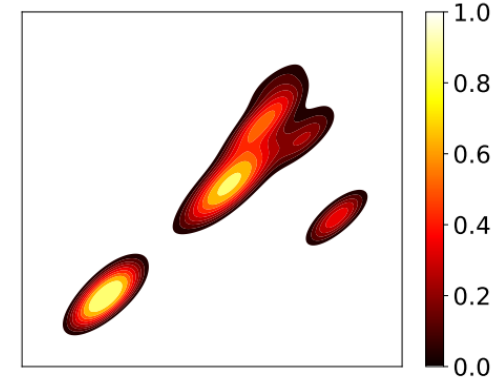
Original Data  
Distribution Of  
MNIST Dataset



Data Distribution Of  
MNIST Dataset  
Inferred By FedFusion



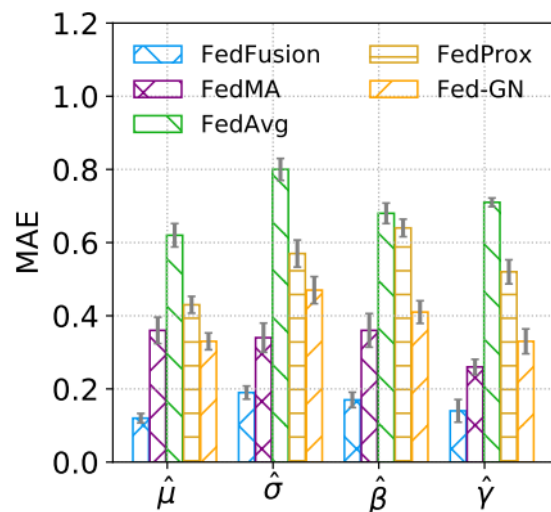
Original Data  
Distribution Of  
CIFAR-10 Dataset



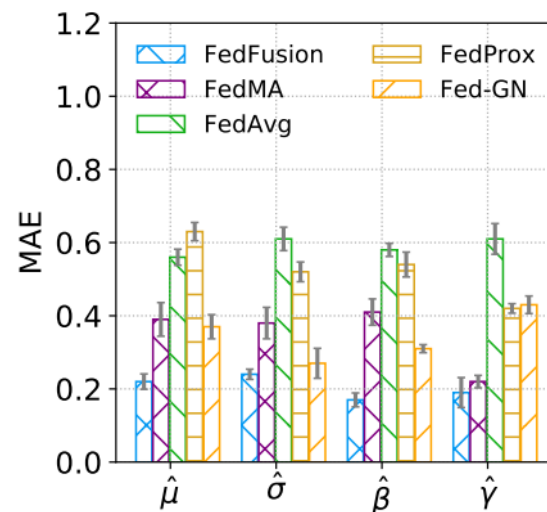
Data Distribution Of  
CIFAR-10 Dataset  
Inferred by FedFusion

FedFusion accurately infer and reconstruct global data distribution, gives Fedfusion ability to approach centralized training.

# Experiment



(a) ResNet18 on CIFAR-10

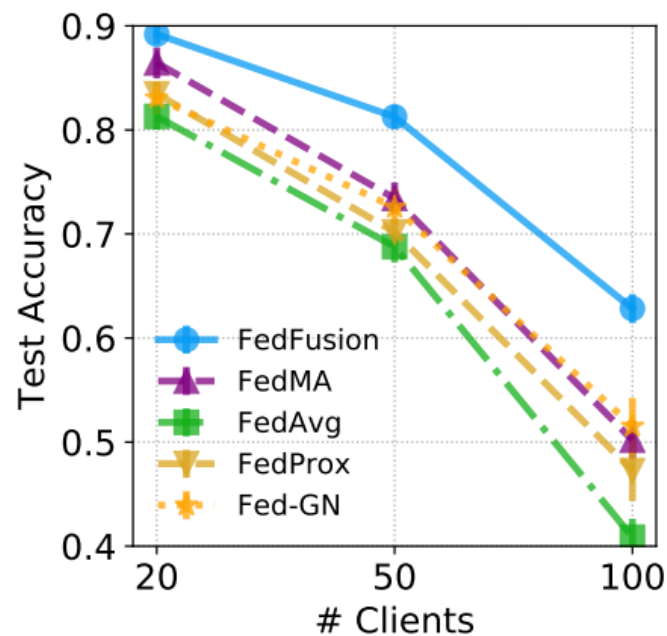
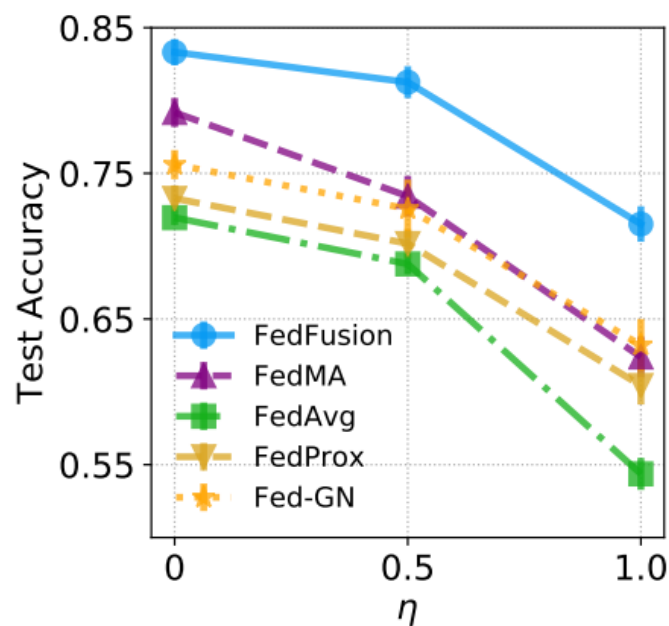


(b) BiLSTM on Sent140

Normalization layer in models trained by FedFusion has fewer bias compared with centralized trained model, also shows FedFusion approximate centralized training well.

Comparison of feature distribution bias

# Experiment



With different hyper-parameter settings, FedFusion shows better robustness.

**Thank You!**