



# FEDBE: Making Bayesian Model Ensemble Applicable to Federated Learning

Hong-You Chen, Wei-Lun Chao  
The Ohio State University

## Highlights

- A novel federated **aggregation** framework based on **Bayesian ensemble**
- Construct global posterior with client models and perform Bayesian model ensemble by sampling
- Leverage **unlabeled data** on the server for **knowledge distillation** for aggregation, bypassing weight average
- Apply stochastic weight average (SWA) [Izmailov et al., 2018] for robust distillation

## Introduction

- *Background: FEDAVG* [McMahan et al., 2017]:

Weight average:

$$\bar{w} \leftarrow \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} w_i.$$

- Non-i.i.d local data  $\rightarrow$  *clients drift away from the ideal global model*
- Element-wise average  $\rightarrow$  *negative effects on over-parameterized and permutation-invariant models like neural networks*

- *Our observation: model ensemble is significantly better than FEDAVG Bayesian model ensemble:*

$$p(y|x; \mathcal{D}) \approx \frac{1}{M} \sum_{m=1}^M p(y|x; w^{(m)}), \text{ where } w^{(m)} \sim p(w|\mathcal{D}).$$

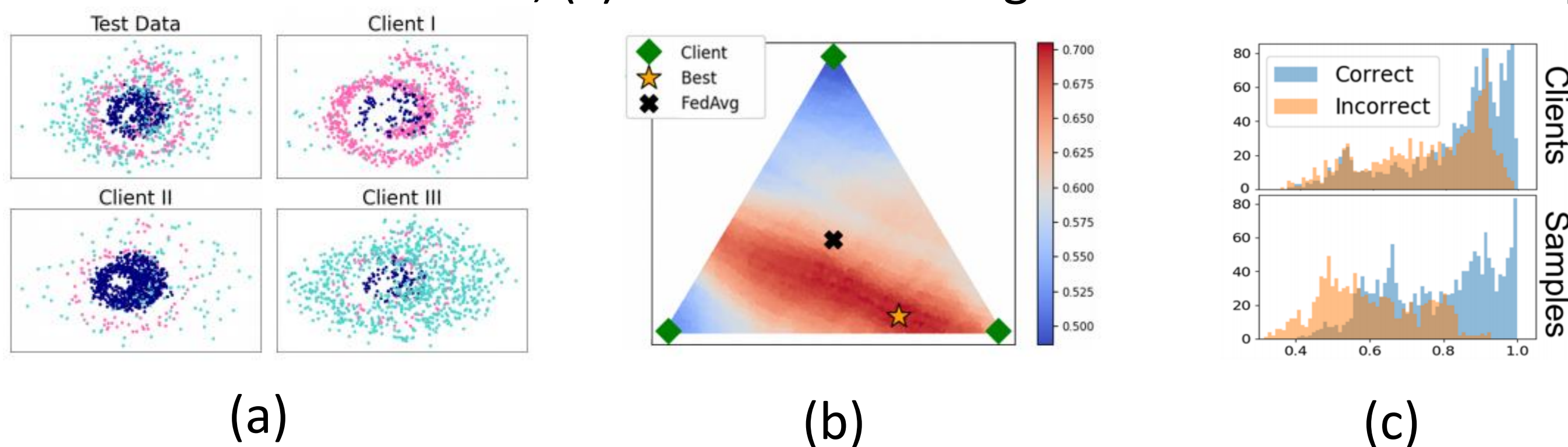
- How to construct global model distribution?  $\rightarrow$  *stochastic weight average-Gaussian (SWAG)* [Maddox et al., 2019]
- How to translate the ensemble to a global model for multi-round federated learning?  $\rightarrow$  *knowledge distillation on server unlabeled data*
- Distillation on noisy pseudo labels?  $\rightarrow$  *stochastic weight average (SWA)*

**Incorporate Bayesian ensemble into federated learning by distilling on more robust aggregation of the outputs of a wide spectrum of global models**

## Federated Bayesian Ensemble

- *A swiss roll three-client federated learning toy example*

(a) non-i.i.d clients, balanced test set; (b) accuracy heatmap of convex combinations of clients; (c) confidence histograms of clients and samples.



- *Construct global model distribution*  $p(w|\mathcal{D})$

**Gaussian**  $w^{(m)} \sim \mathcal{N}(\mu, \Sigma_{\text{diag}})$ :

$$\mu = \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} w_i, \quad \Sigma_{\text{diag}} = \text{diag} \left( \sum_i \frac{|\mathcal{D}_i|}{|\mathcal{D}|} (w_i - \mu)^2 \right)$$

**Dirichlet:**

$$w^{(m)} = \sum_i \frac{\gamma_i^{(m)} |\mathcal{D}_i|}{\sum_{i'} \gamma_{i'}^{(m)} |\mathcal{D}_{i'}|} w_i, \quad \gamma^{(m)} \sim p(\gamma) = p(\gamma_1, \dots, \gamma_{|S|}) = \frac{1}{B(\alpha)} \prod_i \gamma_i^{\alpha_i - 1}$$

- *Construct pseudo labels on server unlabeled dataset*

$$\mathcal{T} = \{(\mathbf{x}_j, \hat{\mathbf{p}}_j)\}_{j=1}^J \quad \hat{\mathbf{p}}_j = \frac{1}{M} \sum_{m=1}^M p(y|\mathbf{x}_j; w^{(m)})$$

- *Robust knowledge distillation with SWA*

Momentum SGD with cyclical learning rate schedule, averaging checkpoints of ends of cycles.

## Algorithm

### Algorithm 1: FEDBE (Federated Bayesian Ensemble)

**Server input** : initial global model  $w$ , SWA scheduler  $\eta_{\text{SWA}}$ , unlabeled data  $\mathcal{U} = \{\mathbf{x}_j\}_{j=1}^J$ ;  
**Client  $i$ 's input** : local step size  $\eta_l$ , local labeled data  $\mathcal{D}_i$ ;  
**for**  $r \leftarrow 1$  **to**  $R$  **do**  
    **Sample** clients  $\mathcal{S} \subseteq \{1, \dots, N\}$ ;  
    **Communicate**  $w$  to all clients  $i \in \mathcal{S}$ ;  
    **for each client**  $i \in \mathcal{S}$  **in parallel do**  
        **Initialize** local model  $w_i \leftarrow w$ ;  
         $w_i \leftarrow$  **Client training**( $w_i, \mathcal{D}_i, \eta_l$ );  
        **Communicate**  $w_i$  to the server;  
    **end**  
    **Construct**  $\bar{w} = \sum_{i \in \mathcal{S}} \frac{|\mathcal{D}_i|}{\sum_{i' \in \mathcal{S}} |\mathcal{D}_{i'}|} w_i$ ;  
    **Construct** global model distribution  $p(w|\mathcal{D})$  from  $\{w_i; i \in \mathcal{S}\}$ ;  
    **Sample**  $M$  global models  $\{w^{(m)}\}_{m=1}^M \sim p(w|\mathcal{D})$ ;  
    **Construct**  $\{w^{(m')}\}_{m'=1}^{M'} = \{\bar{w}\} \cup \{w_i; i \in \mathcal{S}\} \cup \{w^{(m)}\}_{m=1}^M$ ;  
    **Construct**  $\mathcal{T} = \{\mathbf{x}_j, \hat{\mathbf{p}}_j\}_{j=1}^J$ , where  $\hat{\mathbf{p}}_j = \frac{1}{M'} \sum_{m'} p(y|\mathbf{x}_j; w^{(m')})$ ;  
    **Knowledge distillation:**  $w \leftarrow \text{SWA}(\bar{w}, \mathcal{T}, \eta_{\text{SWA}})$ ;  
**end**  
**Server output** :  $w$ .

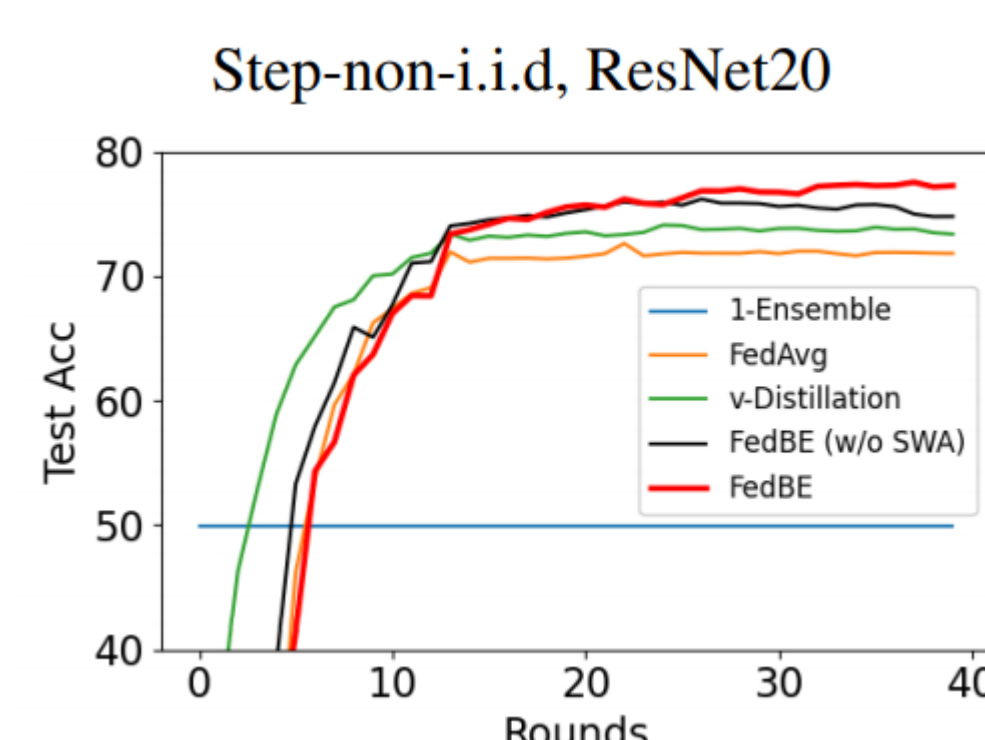
## Experiments and Analysis

- *Non-i.i.d CIFAR-10* (10 clients, 40 rounds, 20 local epochs, 10K unlabeled CIFAR-10 as the server data, Gaussian global distribution)

Non-i.i.d. Type	Method	ConvNet	ResNet20	ResNet32	ResNet44	ResNet56
Step	1-Ensemble	60.5±0.28	49.9±0.46	35.5±0.50	32.8±0.38	23.3±0.52
	FEDAVG	72.0±0.25	70.2±0.17	66.5±0.36	60.5±0.26	51.4±0.15
	v-Distillation	69.2±0.18	72.6±0.62	68.4±0.33	60.4±0.53	56.4±1.10
	FEDBE (w/o SWA)	72.1±1.21	74.9±1.41	71.1±0.75	61.0±0.75	56.6±0.85
	<b>FEDBE</b>	<b>74.5±0.51</b>	<b>77.5±0.42</b>	<b>72.7±0.27</b>	<b>65.5±0.32</b>	<b>60.7±0.45</b>
Dirichlet	1-Ensemble	63.3±0.56	45.2±1.06	39.5±0.78	31.5±0.77	27.2±0.65
	FEDAVG	72.3±0.12	74.4±0.36	73.4±0.23	67.1±0.54	62.2±0.45
	v-Distillation	67.7±0.98	73.1±0.78	70.8±0.64	66.9±0.85	62.8±0.66
	FEDBE (w/o SWA)	70.1±0.42	75.9±0.56	73.9±0.55	68.2±0.72	63.2±0.71
	<b>FEDBE</b>	<b>73.9±0.45</b>	<b>78.2±0.36</b>	<b>77.7±0.45</b>	<b>71.5±0.38</b>	<b>67.0±0.30</b>
Centralized*	SGD	84.5	91.7	92.6	93.1	93.4

- *Compatibility of FEDBE with FEDAVGM and FEDPROX*

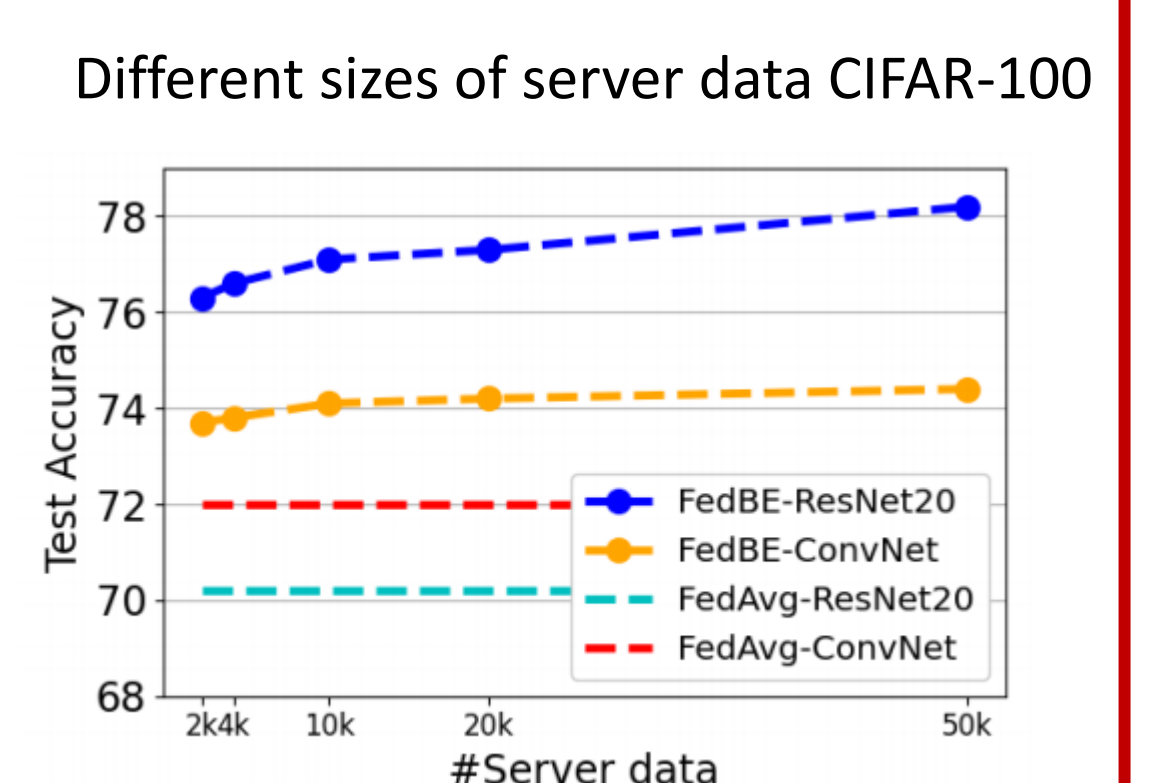
Non-i.i.d. Type	Method	ConvNet	ResNet20	ResNet32	ResNet44	ResNet56
Step	FEDPROX	72.5±0.71	71.1±0.52	67.7±0.26	60.4±0.71	54.9±0.66
	<b>FEDBE +FEDPROX</b>	<b>74.9±0.38</b>	<b>77.7±0.45</b>	<b>72.9±0.44</b>	<b>64.5±0.37</b>	<b>60.1±0.62</b>
	FEDAVGM	72.3±0.55	73.2±0.57	70.0±0.62	59.9±0.65	52.7±0.49
	<b>FEDBE +FEDAVGM</b>	<b>74.5±0.47</b>	<b>78.0±0.46</b>	<b>73.6±0.50</b>	<b>65.5±0.40</b>	<b>59.7±0.51</b>
Dirichlet	FEDPROX	72.6±0.38	76.1±0.49	73.4±0.51	68.1±0.79	60.9±0.46
	<b>FEDBE +FEDPROX</b>	<b>74.6±0.35</b>	<b>78.7±0.49</b>	<b>77.3±0.60</b>	<b>71.7±0.43</b>	<b>66.5±0.41</b>
	FEDAVGM	73.0±0.43	76.5±0.44	75.5±0.79	67.7±0.46	58.9±0.72
	<b>FEDBE +FEDAVGM</b>	<b>74.4±0.49</b>	<b>78.5±0.66</b>	<b>78.5±0.26</b>	<b>72.0±0.51</b>	<b>67.0±0.55</b>



- FEDBE with SWA consistently outperforms baselines and improves on other methods
- We observe deeper models suffer from performance drop in federated learning
- FEDBE performs well with small and out-of-domain unlabeled server data

- *Effects of unlabeled data for distillation*

Non-i.i.d. Type	$\mathcal{U}$	$ \mathcal{U} $	ConvNet	ResNet20
Step	CIFAR-10	10K	74.5±0.51	77.5±0.42
	CIFAR-100	50K	74.4±0.45	78.2±0.58
	Tiny-ImageNet	100K	74.5±0.64	77.1±0.51
Dirichlet	CIFAR-10	10K	73.9±0.45	78.2±0.36
	CIFAR-100	50K	73.5±0.41	78.6±0.63
	Tiny-ImageNet	100K	74.0±0.35	78.2±0.72



- *FL with systems heterogeneity*

Local training epochs  $E \sim (0, 20]$

Method	ConvNet	ResNet20	ResNet32
FEDAVG	70.6±0.46	69.9±0.59	64.0±0.50
FEDPROX	71.2±0.55	69.4±0.48	65.9±0.63
<b>FEDBE</b>	<b>73.3±0.56</b>	<b>77.1±0.61</b>	<b>70.2±0.39</b>
<b>+FEDPROX</b>	<b>73.7±0.24</b>	<b>77.5±0.51</b>	<b>71.6±0.37</b>

- **Conclusion:** FEDBE is robust to non-i.i.d. clients and network architectures, compatible to regularized client training and server momentum methods