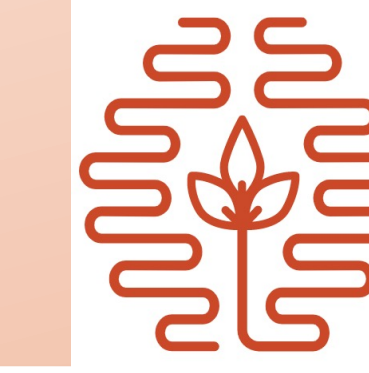


# Adversarial Robustness through Bias Variance Decomposition: A New Perspective for Federated Learning



Yao Zhou<sup>\*1,2</sup>, Jun Wu<sup>\*1</sup>, Haixun Wang<sup>2</sup>, Jingrui He<sup>1</sup>

<sup>1</sup>University of Illinois at Urbana-Champaign, <sup>2</sup>Instacart  
yaozhou3@illinois.edu, junwu3@illinois.edu, haixun@gmail.com, jingrui@illinois.edu



**AIFARMS**  
Artificial Intelligence for Future Agricultural Resilience, Management, and Sustainability



## Problem Definition

### Motivation

- FedAvg
- Client update with local SGD

$$w_k \leftarrow w_k - \alpha \frac{1}{n_k} \sum_{i=1}^{n_k} L(f_{D_k}(x_i^k; w_k), t_i^k)$$

- Server update

$$w_G = \sum_{k=1}^K \frac{n_k}{n} w_k$$

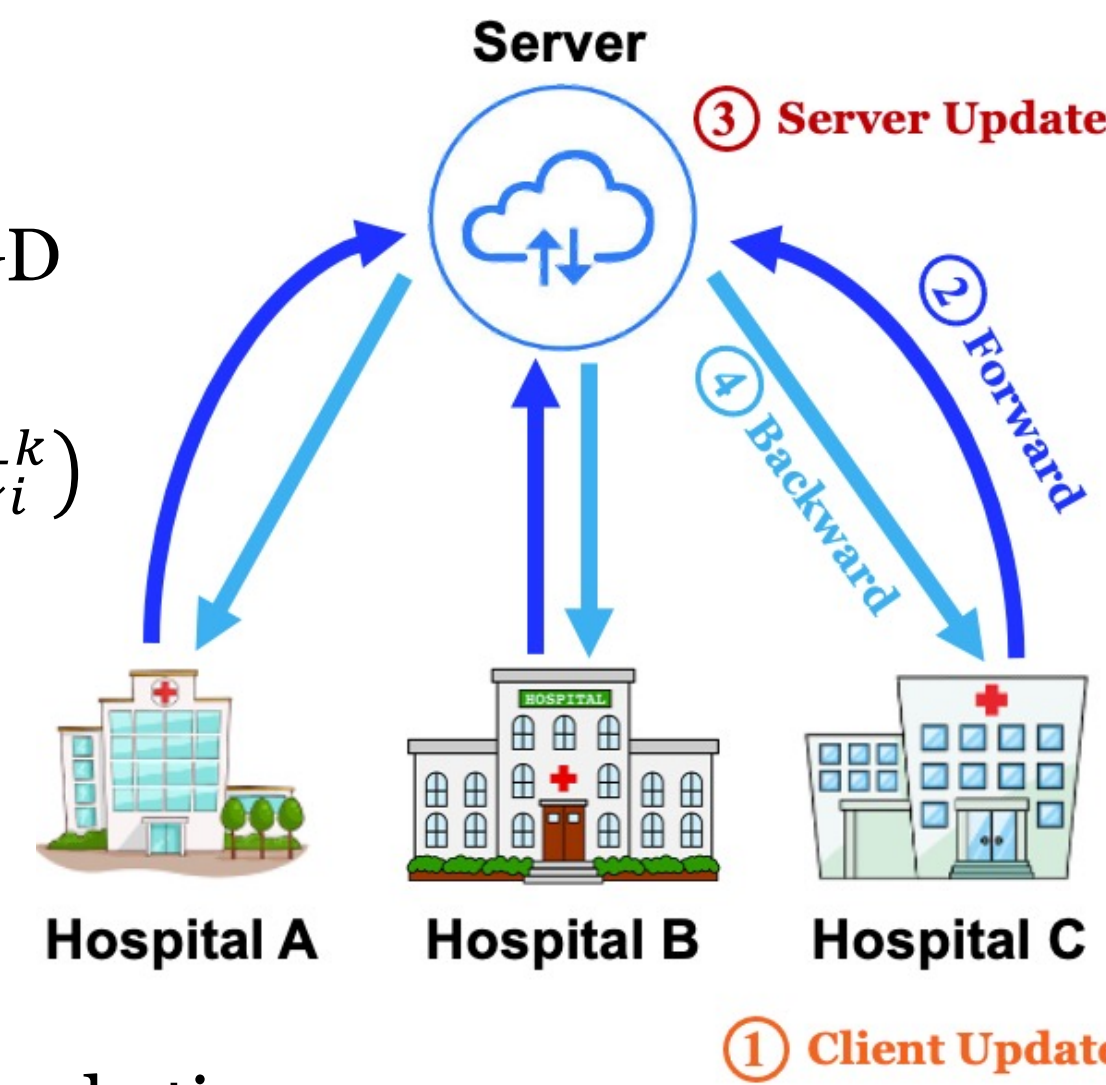
- Vulnerable to adversarial perturbation for model inference

### Adversarially-robust federated learning

- Given
  - $K$  clients with local data  $\{D_k\}_{k=1}^K$
  - A learning algorithm  $f(\cdot)$
  - Loss function  $L(\cdot, \cdot)$
  - A public auxiliary training set  $D_s$

- Output

- A trained model on the central server that is **robust against adversarial perturbations** on the test set  $D_{test}$



## Algorithm: Fed BVA

### Server update

- Model aggregation:  $w_G = \text{Aggregate}(w_1, w_2, \dots, w_K)$
- Adversarial example generation: For any  $x \in D_s$ 

$$\max_{\hat{x} \in \Omega(x)} B(\hat{x}; w_1, w_2, \dots, w_K) + V(\hat{x}; w_1, w_2, \dots, w_K)$$
  - BV-FGSM:
 
$$\hat{x} \leftarrow x + \epsilon \cdot \text{sign}(\nabla_x (B(x; w_1, w_2, \dots, w_K) + V(x; w_1, w_2, \dots, w_K)))$$
  - For cross-entropy loss function,

$$\nabla_x B_{CE}(x; w_1, w_2, \dots, w_K) = \frac{1}{K} \sum_{k=1}^K \nabla_x L(f_{D_k}(x; w_k), t)$$

$$\nabla_x V_{CE}(x; w_1, w_2, \dots, w_K) = \frac{1}{K} \sum_{k=1}^K \sum_{c=1}^C (\log y_m^{(j)} + 1) \cdot \nabla_x f_{D_k}(x; w_k)$$

### Backward communication

- Send global model parameters  $w_G$  and poisoned examples  $\hat{x}$  to candidate client

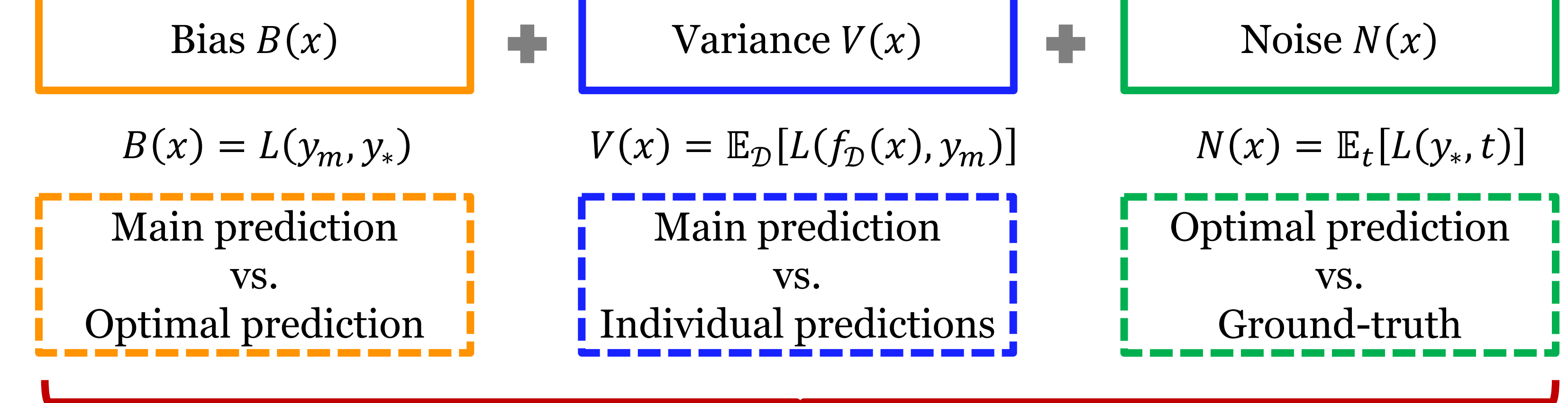
### Client update

$$\text{Robust training: } \min_{w_k} \frac{1}{n_k} \sum_{i=1}^{n_k} L(f_{D_k}(x_i^k; w_k), t_i^k) + \frac{1}{n_s} \sum_{j=1}^{n_s} L(f_{D_k}(\hat{x}_j^s; w_k), t_j^s)$$

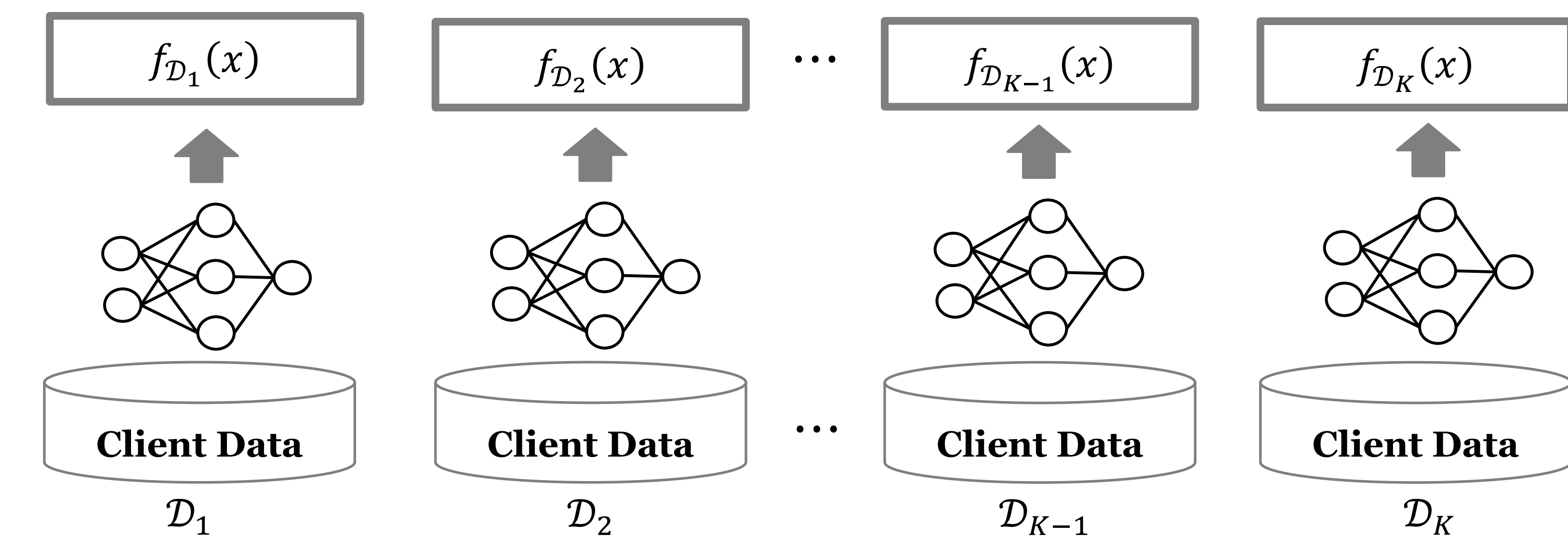
### Forward communication

- Upload local parameter updates to the server

Generalization performance  $\mathbb{E}_{D,t}[L(f_D(x), t)]$



Local client models



## Conclusion

- Problem:** The adversarial robustness of federated learning is studied under the observation that federated learning model is vulnerable to evasion attacks when it is deployed.
- Algorithm:** By investigating the generalization error of clients' local models, we propose a bias-variance oriented adversarial training algorithm Fed\_BVA for robust federated learning.
- Evaluation:** Extensive experiments confirm the effectiveness and efficiency of the Fed\_BVA algorithm.

## Acknowledgments

This work is supported by National Science Foundation under Award No. IIS-1947203, IIS-2117902, IIS-2137468, and Agriculture and Food Research Initiative (AFRI) grant no. 2020-67021-32799/project accession no.1024178 from the USDA National Institute of Food and Agriculture. The views and conclusions are those of the authors and should not be interpreted as representing the official policies of the funding agencies or the government.



## Experimental Results

### Performance comparison

Method	IID			non-IID		
	Clean	FGSM	PGD-20	Clean	FGSM	PGD-20
Centralized	<b>0.991</b> $\pm$ 0.000	0.689 $\pm$ 0.000	0.182 $\pm$ 0.000	n/a	n/a	n/a
FedAvg	0.989 $\pm$ 0.001	0.669 $\pm$ 0.009	0.267 $\pm$ 0.014	<b>0.980</b> $\pm$ 0.002	0.491 $\pm$ 0.067	0.158 $\pm$ 0.074
FedAvg_AT	0.988 $\pm$ 0.000	0.802 $\pm$ 0.001	0.512 $\pm$ 0.042	0.974 $\pm$ 0.005	0.649 $\pm$ 0.066	0.363 $\pm$ 0.066
Fed_Bias	0.986 $\pm$ 0.000	0.812 $\pm$ 0.009	0.583 $\pm$ 0.036	0.971 $\pm$ 0.004	0.679 $\pm$ 0.040	0.394 $\pm$ 0.103
Fed_Variance	0.985 $\pm$ 0.001	0.803 $\pm$ 0.007	0.572 $\pm$ 0.019	0.973 $\pm$ 0.005	0.684 $\pm$ 0.004	0.395 $\pm$ 0.049
Fed_BVA	0.986 $\pm$ 0.001	0.818 $\pm$ 0.003	0.613 $\pm$ 0.020	0.969 $\pm$ 0.002	0.705 $\pm$ 0.009	0.469 $\pm$ 0.031
EAT	0.981 $\pm$ 0.000	<b>0.902</b> $\pm$ 0.001	0.811 $\pm$ 0.004	0.972 $\pm$ 0.002	0.789 $\pm$ 0.016	0.415 $\pm$ 0.035
EAT+Fed_BVA	0.980 $\pm$ 0.001	0.901 $\pm$ 0.006	<b>0.821</b> $\pm$ 0.013	0.965 $\pm$ 0.005	<b>0.811</b> $\pm$ 0.020	<b>0.670</b> $\pm$ 0.014

Robustness on MNIST under IID and non-IID settings

### Ablation study

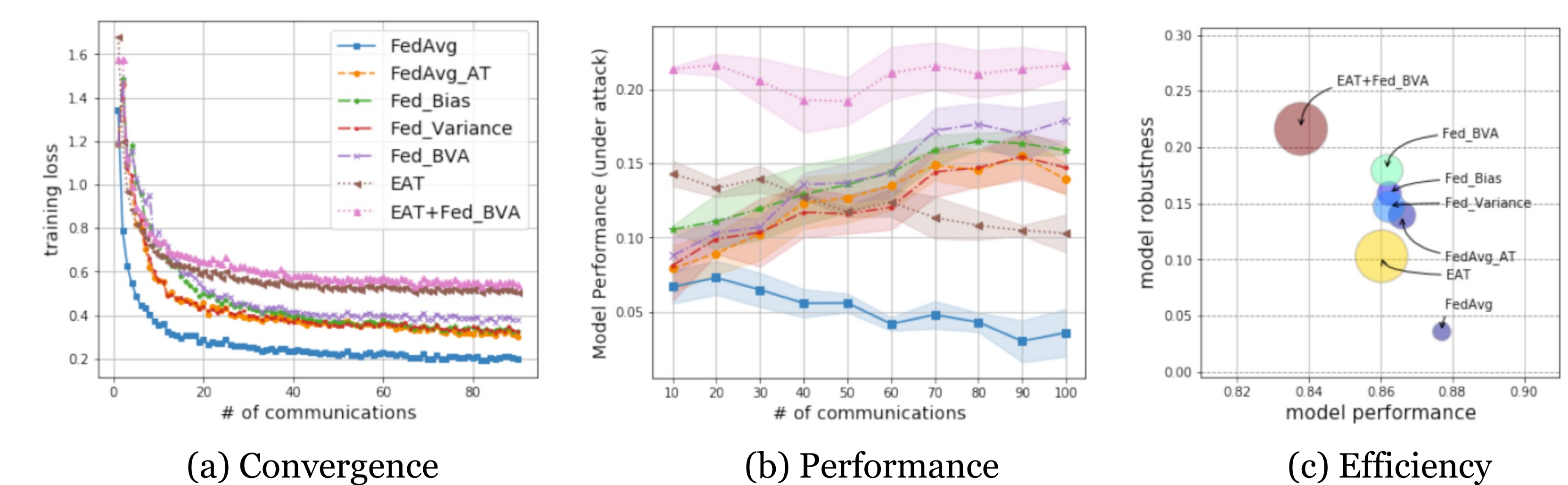
Loss	Clean	Fed_BVA		
		BiasOnly	VarianceOnly	BVA
CE	0.588(38.13s)	<b>0.763</b> (47.58s)	<b>0.759</b> (63.46s)	<b>0.776</b> (63.67s)
MSE	<b>0.601</b> (39.67s)	0.711(65.03s)	0.711(162.40s)	0.712(179.60s)

(a) Cross-Entropy (CE) vs. Mean Squared Error (MSE)

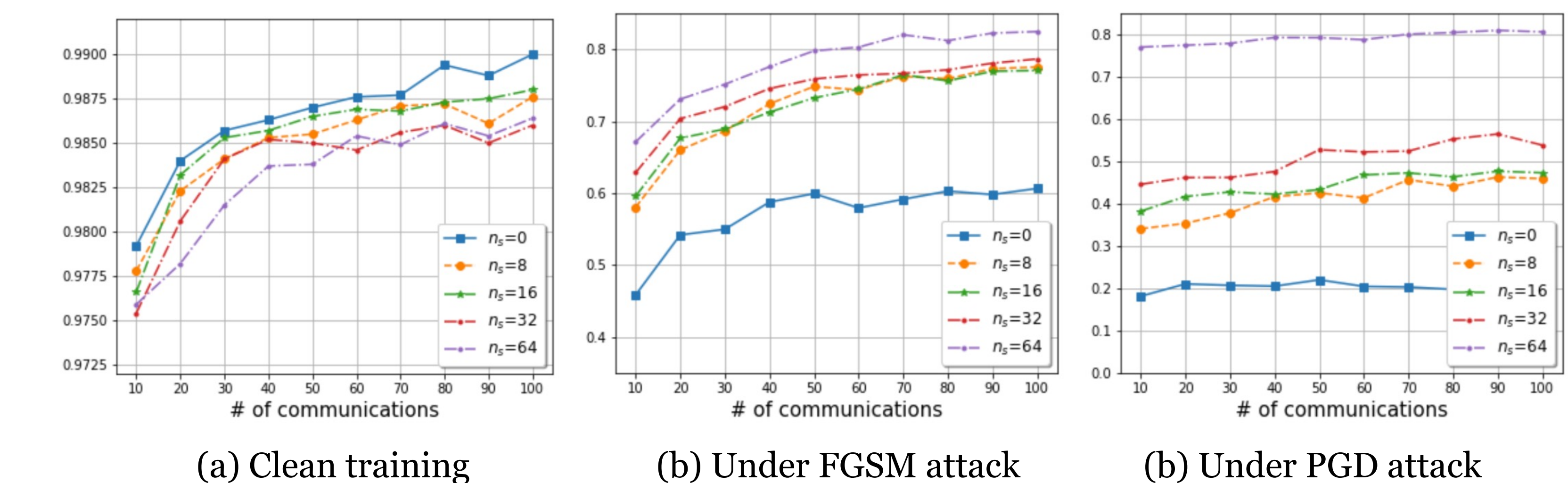
Method	IID			non-IID		
	FGSM	PGD-10	PGD-20	FGSM	PGD-10	PGD-20
FedAvg	0.588	0.620	0.205	0.147	0.525	0.089
Fed_BVA(BV-FGSM)	<b>0.776</b>	0.793	0.570	<b>0.670</b>	0.695	0.472
Fed_BVA(BV-PGD)	0.757	<b>0.840</b>	<b>0.632</b>	0.659	<b>0.784</b>	<b>0.575</b>

(b) BV-FGSM vs. BV-PGD

### Model analysis



### Hyperparameter sensitivity – size of public data set $n_s$



(a) Clean training

(b) Under FGSM attack

(c) Under PGD attack