

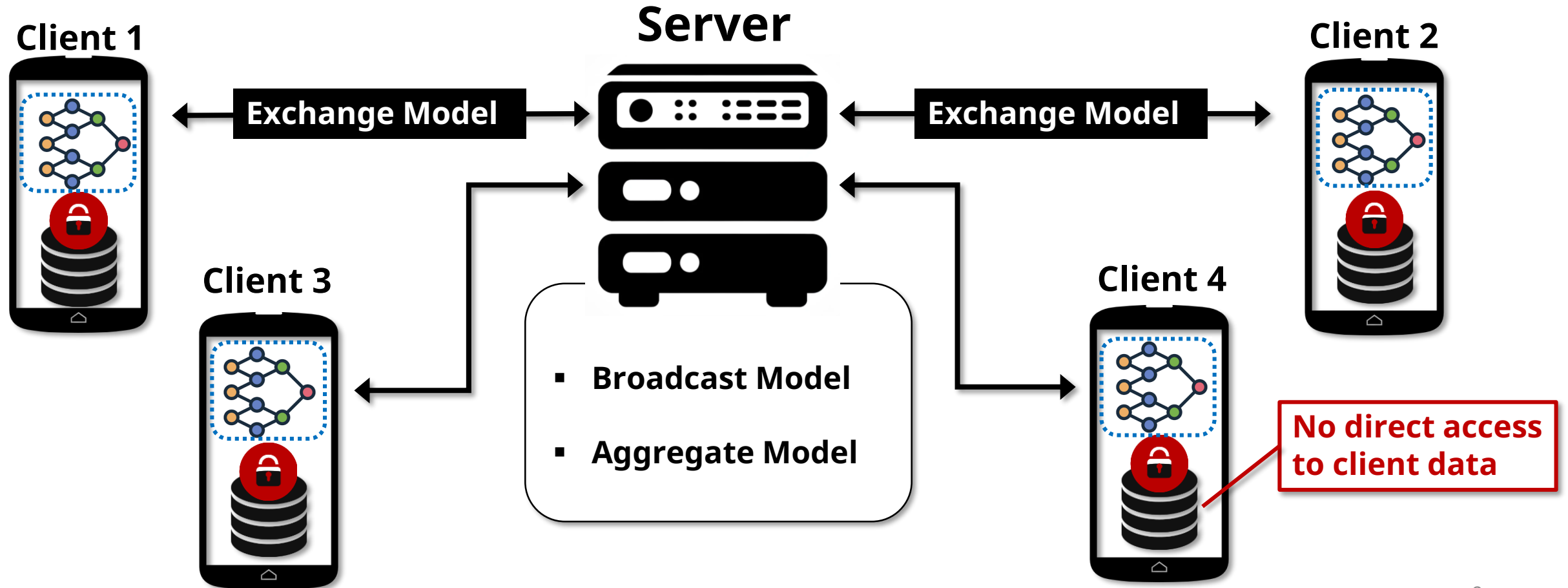
# **FedSOL: Stabilized Orthogonal Learning with Proximal Restrictions in Federated Learning**

**Gihun Lee, Minchan Jeong, Sangmook Kim, Jaehoon Oh,  
and Se-Young Yun**

# 1-1. Federated Learning

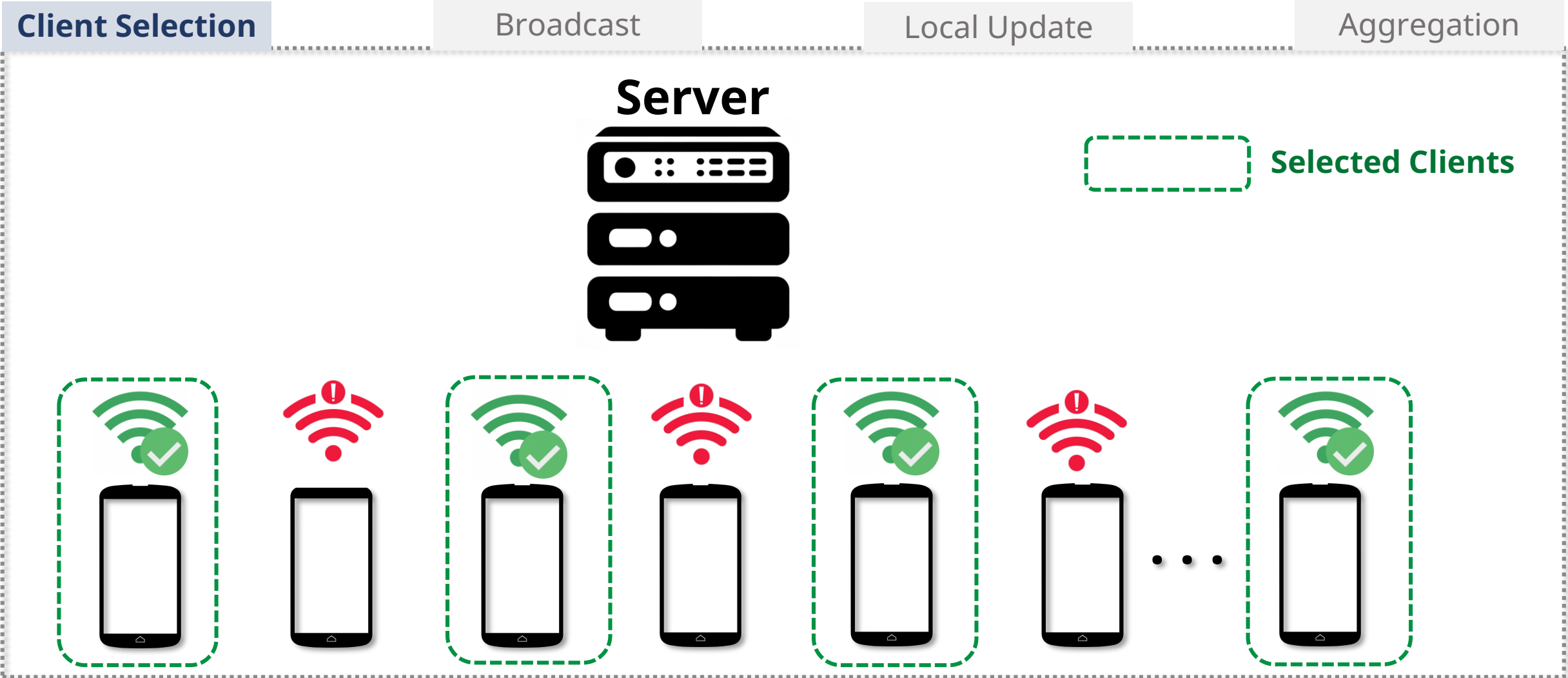
## 1 Background

**Federated Learning (FL)** is a **Distributed Learning** paradigm in which multiple clients collaboratively learn a machine learning model while preserving their **Data Privacy**.



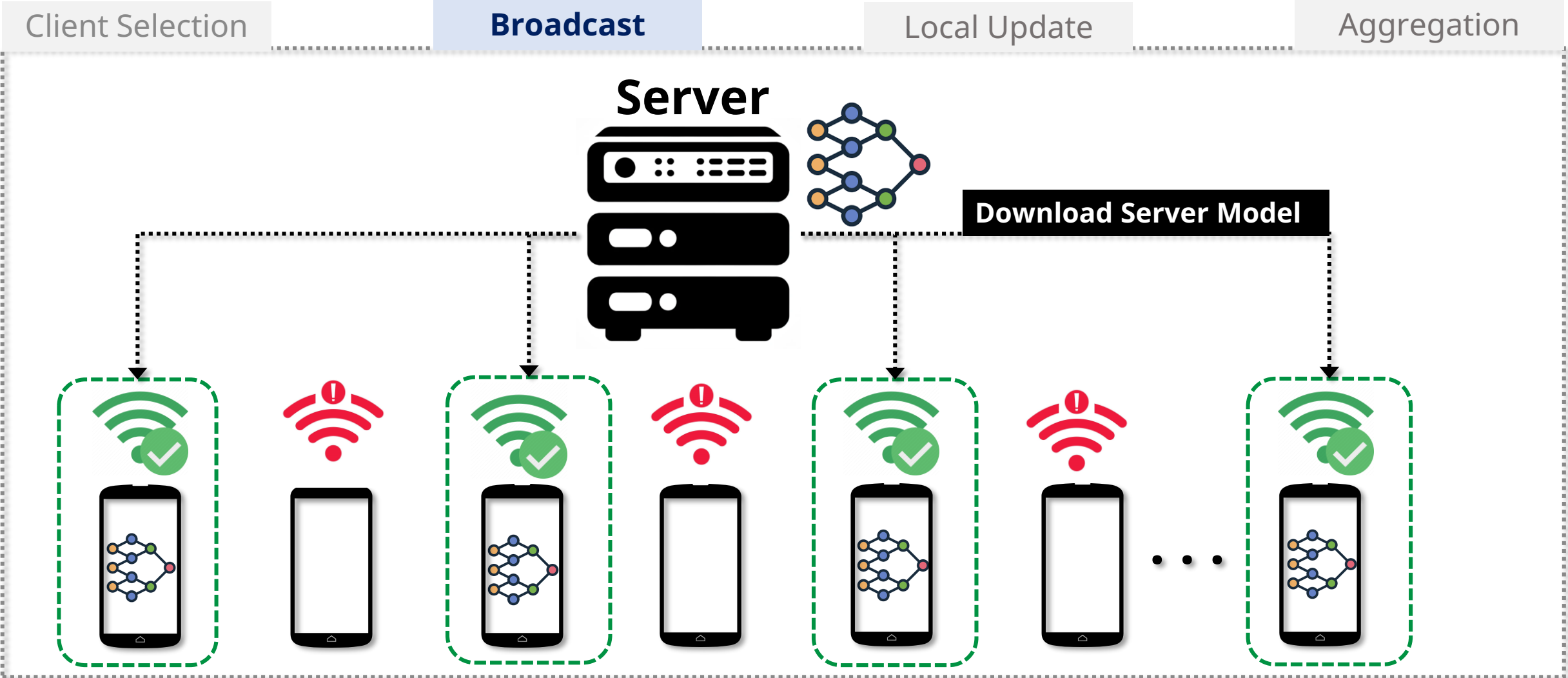
# 1-1. Federated Learning

## 1 Background

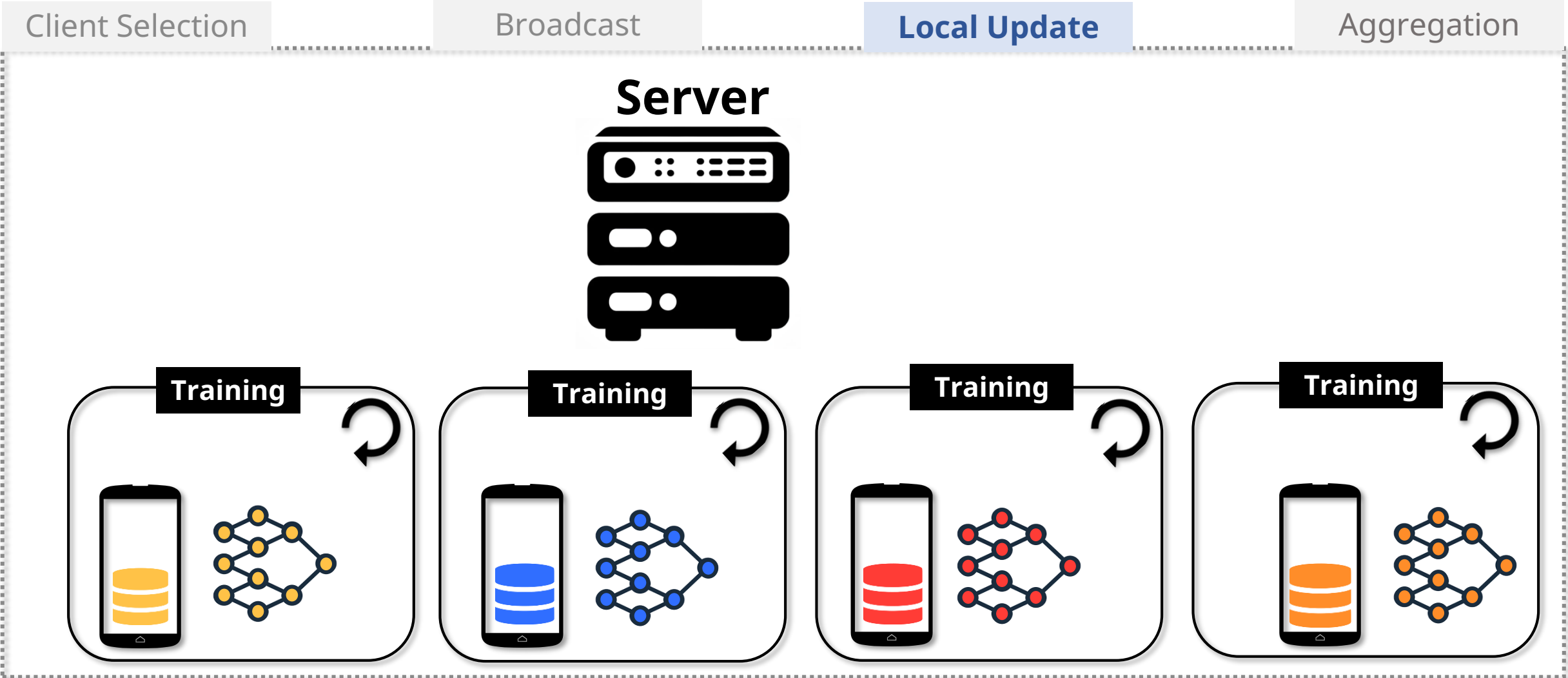


# 1-1. Federated Learning

## 1 Background



# 1-1. Federated Learning



# 1-1. Federated Learning

## 1 Background

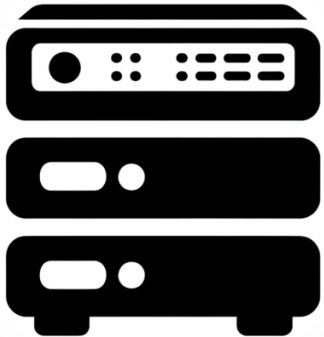
Client Selection

Broadcast

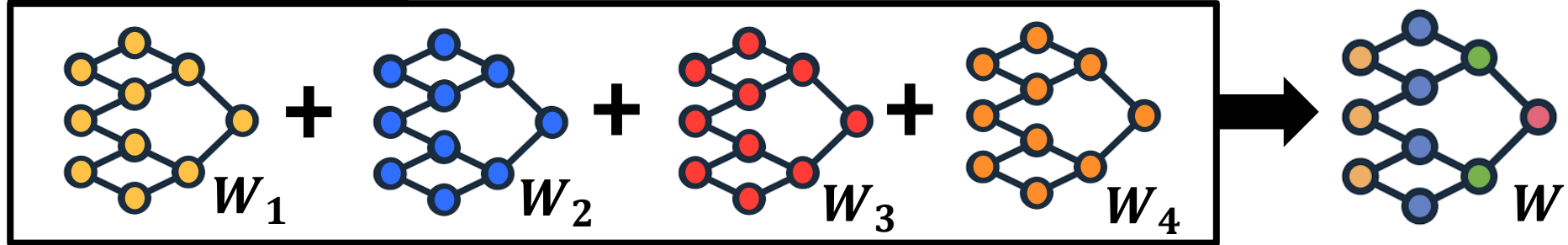
Local Update

Aggregation

### Server



#### Model Aggregation



**FedAvg** (McMahan et al., 2018)

$$W^t = \frac{1}{K} \sum_{i=1}^K p_i W_i^{t-1} \quad \text{Weighted Average of Local Models}$$

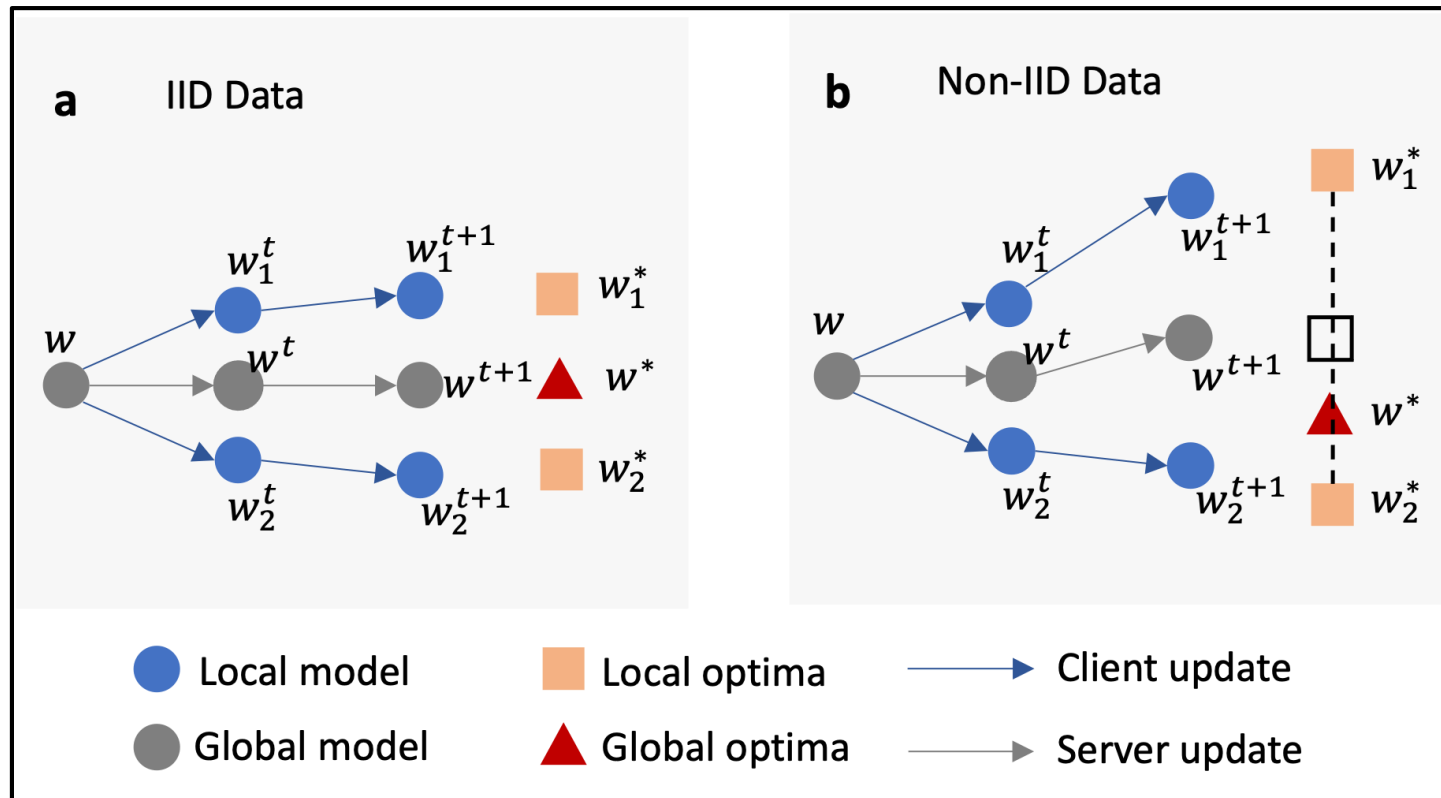
Konečný, Jakub, et al. "Federated learning: Strategies for improving communication efficiency." arXiv preprint arXiv:1610.05492 (2016).  
Konečný, Jakub, et al. "Federated optimization: Distributed machine learning for on-device intelligence." arXiv preprint arXiv:1610.02527 (2016).

# 1-2. Data Heterogeneity

## 1 Background

When client data is **heterogeneous**, FL suffers from "**Client Drift**".

→ How to accumulate knowledge from heterogeneous clients in a single **Global Model**?



Tan, Alysia Ziyang, et al. "Towards personalized federated learning." IEEE Transactions on Neural Networks and Learning Systems (2022).

# 1-2. Data Heterogeneity

- **Sharding:** Uniform Size, Different Distributions  
(Higher Heterogeneity at **lower  $s$** )
- **Latent Dirichlet Allocation (LDA):** Varying Size & Distributions  
(Higher Heterogeneity at **lower  $\alpha$** )

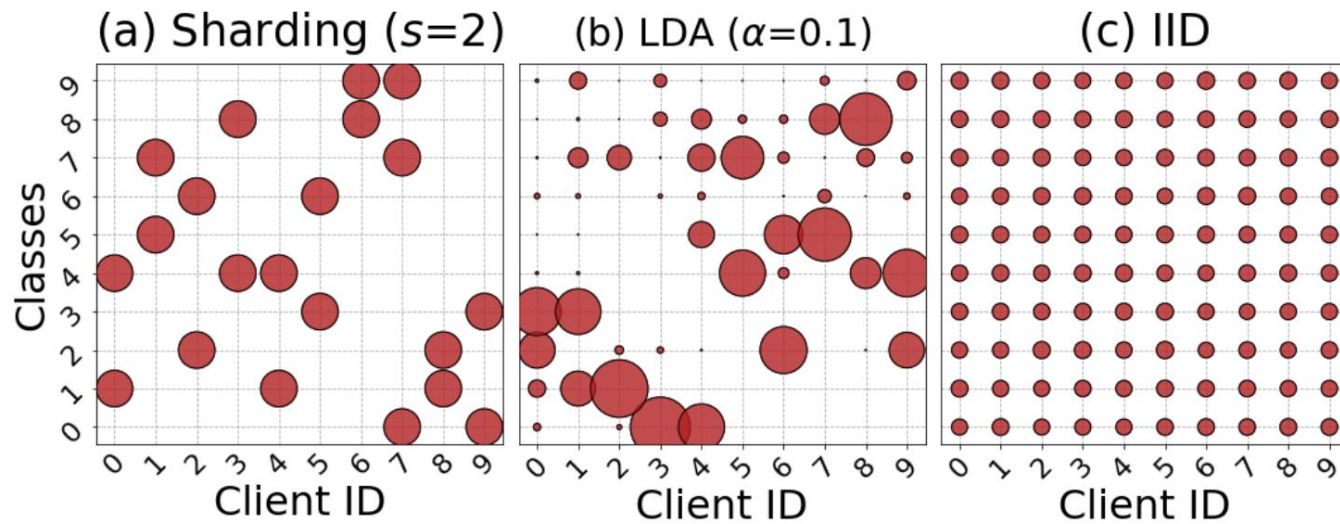


Figure 2. CIFAR-10 partition examples across 10 clients.



### Typical Local Objective:

$$\mathcal{L}^k(\mathbf{w}_k) = \mathcal{L}_{\text{local}}^k(\mathbf{w}_k) + \beta \cdot \mathcal{L}_p^k(\mathbf{w}_k; \mathbf{w}_g)$$

**Original Objective**  
(e.g., CE Loss)

**Proximal Objective**  
(e.g., FedProx, SCAFFOLD, FedKD...)

However, those two losses **conflicts** each other.

**Knowledge  
Acquisition**



**Knowledge  
Preservation**

## 2. Motivation – Proximal Objectives

## 2 Motiv

### Examples of Proximal Objectives:

**FedGKT** (NeurIPS 2020)

$$\mathcal{L}_c = \mathcal{L}_{CE} + D_{KL}(p_s || p_k)$$

**FedProx** (ICML 2019)

$$F_k(w) + \frac{\mu}{2} \|w - w^t\|^2$$

**SCAFFOLD** (ICML 2020)

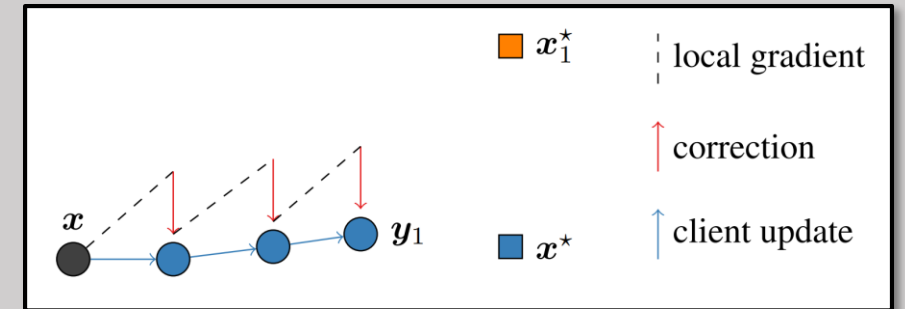
$$y_i \leftarrow y_i - \eta_l (g_i(y_i) + c - c_i)$$

**FedDyn** (ICLR 2021)

$$\mathcal{L}_k(\theta) - \langle \nabla \mathcal{L}_k(\theta_k^{t-1}), \theta \rangle + \frac{\alpha}{2} \|\theta - \theta^{t-1}\|^2$$

**MOON** (CVPR 2021)

$$\ell_{sup}(w_i^t; (x, y)) + \mu \ell_{con}(w_i^t; w_i^{t-1}; w_t; x)$$

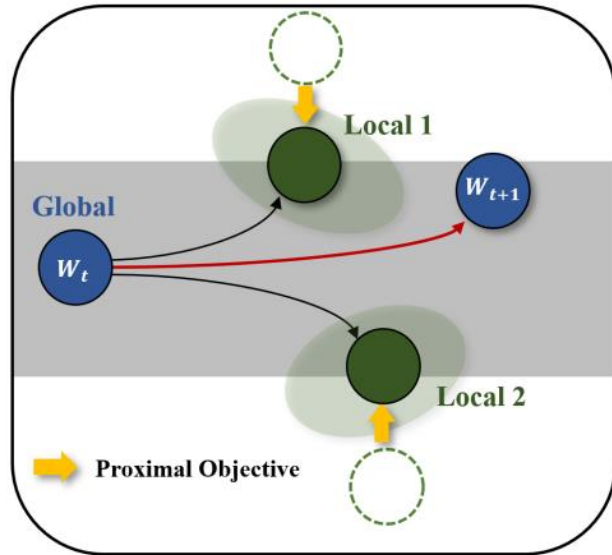


$$\mathcal{L}^k(w_k) = \mathcal{L}_{\text{local}}^k(w_k) + \beta \cdot \mathcal{L}_p^k(w_k; w_g)$$

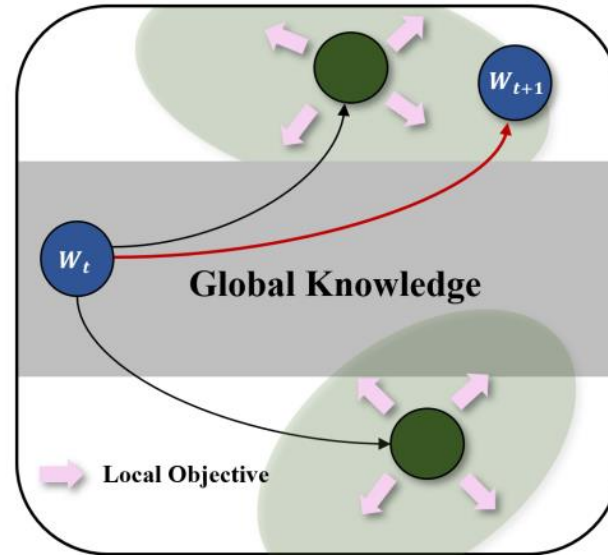
**Original Objective** **Proximal Objective**

## 2. Motivation – Global vs Local Knowledge

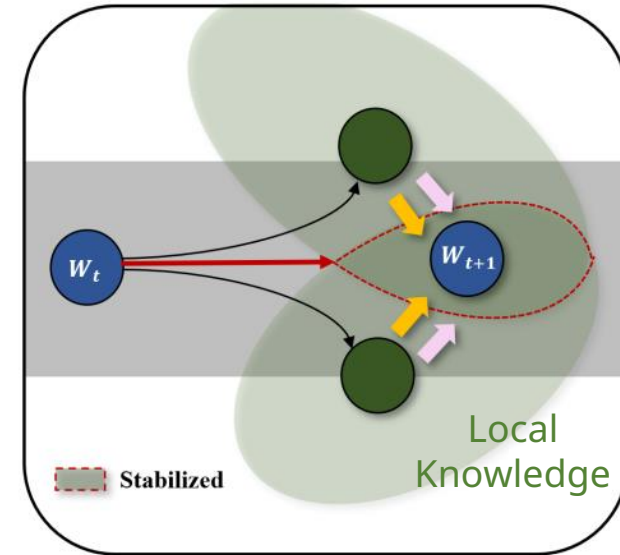
FL must navigate the **balance between those two conflicting objectives.**



(a) Global Knowledge Preservation

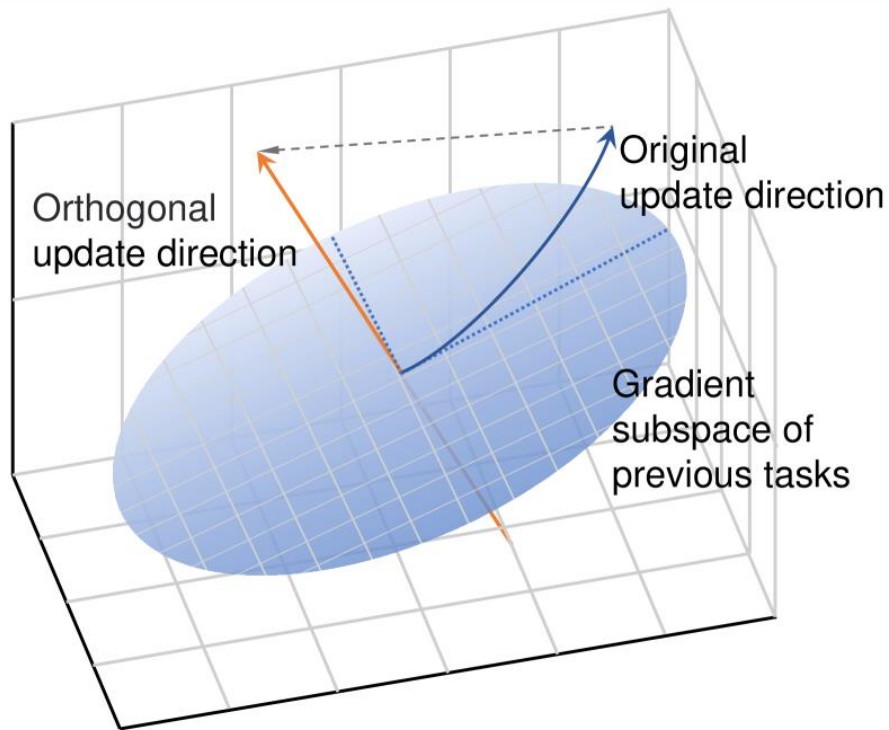


(b) Local Knowledge Acquisition



(c) FedSOL (Ours)

“How should the local learning **acquire local knowledge** while *minimizing its conflicts* with **global knowledge**?”



### Orthogonal Learning in CL

Find update direction that is orthogonal to the previous task.



### Challenges to implement to FL

- Cannot retain past data or gradients for reference for global knowledge.
- The global distribution overlaps with individual local distributions.

## 2. Motivation – Gradient Projection

## 2 Motiv

### Proximal Gradient Projection

$$\mathbf{g}_u^{\text{Proj}} = \mathbf{g}_l - \frac{\mathbf{g}_l^T \mathbf{g}_p}{\mathbf{g}_p^T \mathbf{g}_p} \mathbf{g}_p \quad \text{if } \mathbf{g}_l^T \mathbf{g}_p < 0$$

$\mathbf{g}_l$ : Gradient on **Original Local Loss**

$\mathbf{g}_p$ : Gradient on **Proximal Loss**

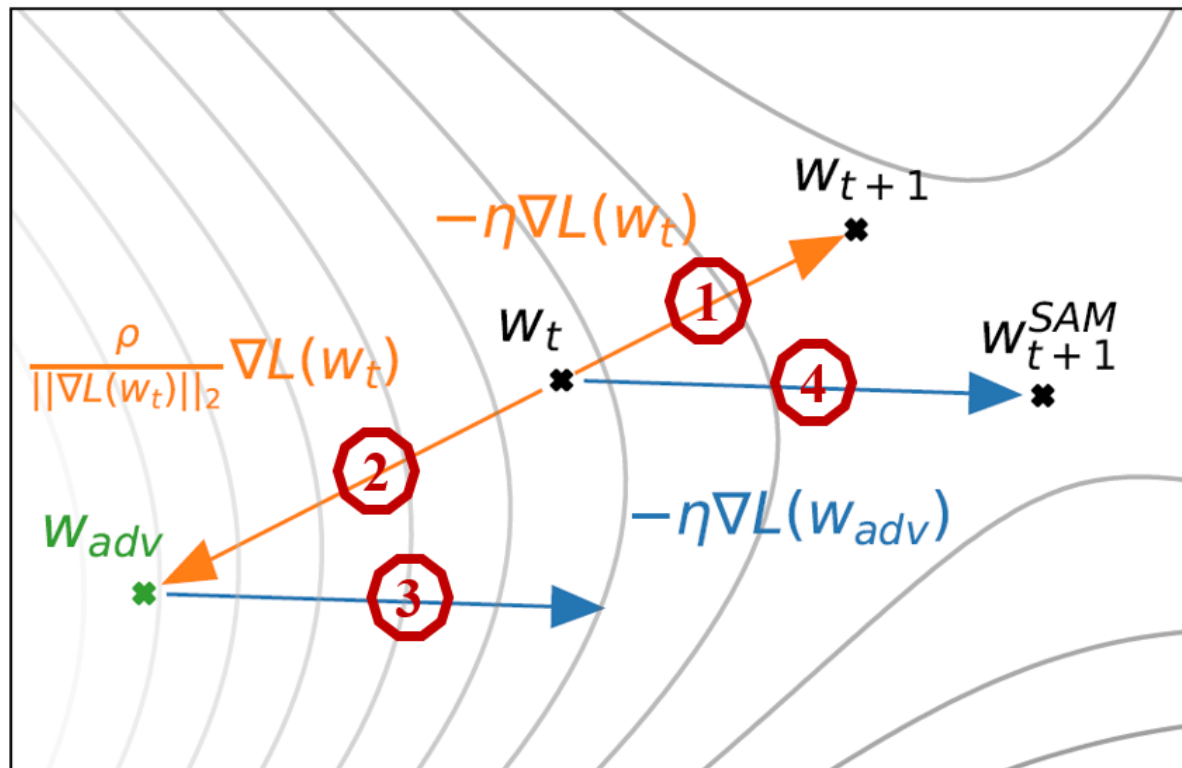
$\mathbf{g}_u$ : Gradient for Parameter Update

### Experiment: CIFAR-10 ( $\alpha = 0.1$ )

Proximal Loss	Usage	
	Base	Projection
None (FedAvg)		56.13 $\pm$ 0.78
L2 Distance	<b>59.80</b> $\pm$ 1.12	56.35 $\pm$ 2.85 (- 3.45)
KL-Divergence	<b>60.31</b> $\pm$ 2.07	50.88 $\pm$ 3.55 (- 9.43)

**Directly negating the proximal gradient ( $\mathbf{g}_p$ )** rather *undermines* the local learning.

## Sharpness-Aware Minimization (SAM)



**Perturb & Update** (*using same  $\mathcal{L}$* )

$$\min_w \max_{\|\epsilon\|_2 < \rho} \mathcal{L}(w + \epsilon)$$

Use single gradient step.

$$\epsilon^* = \rho \frac{\nabla_w L(w)}{\|\nabla_w L(w)\|_2}$$

# 3. Main Approach - Stabilized Orthogonal Learning

[Main Idea]: Perturb using “Global Knowledge” & Update to acquire “Local Knowledge”

## Step 1: Weight Perturbation

$$\epsilon_p^* = \underset{\|\epsilon\|_2 \leq \rho}{\operatorname{argmax}} \mathcal{L}_p^k(\mathbf{w}_k + \epsilon; \mathbf{w}_g) \approx \rho \frac{\mathbf{g}_p}{\|\mathbf{g}_p\|_2}$$

Perturbation Strength

Adversarial perturb using **Proximal Loss  $L_p^k$**

## Step 2: Parameter Update

$$\mathbf{w}_k \leftarrow \mathbf{w}_k - \gamma \cdot \nabla_{\mathbf{w}_k} \mathcal{L}_{\text{local}}^k(\mathbf{w}_k + \epsilon_p^*)$$

Update the model with the **Original Local Loss  $L_{\text{local}}^k$** .

# 3. Main Approach - Stabilized Orthogonal Learning

## 3 FedSOL

[Main Idea]: Regularize **Local Loss Surface** using "Global Perspective"

Local gradient and its direction

$$g_l = \nabla_w \mathcal{L}_{local}^k(w_k) \quad \hat{g}_l = \frac{g_l}{\|g_l\|_2}$$

Proximal gradient and its direction

$$g_p = \nabla_w \mathcal{L}_p^k(w_k) \quad \hat{g}_p = \frac{g_p}{\|g_p\|_2}$$

**FedSAM**

$$g_u^{FedSAM} = \nabla_{w_k} \mathcal{L}_{local}^k(w_k + \epsilon^*) \\ \approx g_l + \rho \nabla_{w_k}^2 \mathcal{L}_{local}^k(w_k) \hat{g}_l$$

Regularize Loss Surface w.r.t. "Local Loss"

**FedSOL**

$$g_u^{FedSOL} = \nabla_{w_k} \mathcal{L}_{local}^k(w_k + \epsilon_p^*) \\ \approx g_l + \rho \nabla_{w_k}^2 \mathcal{L}_{local}^k(w_k) \hat{g}_p$$

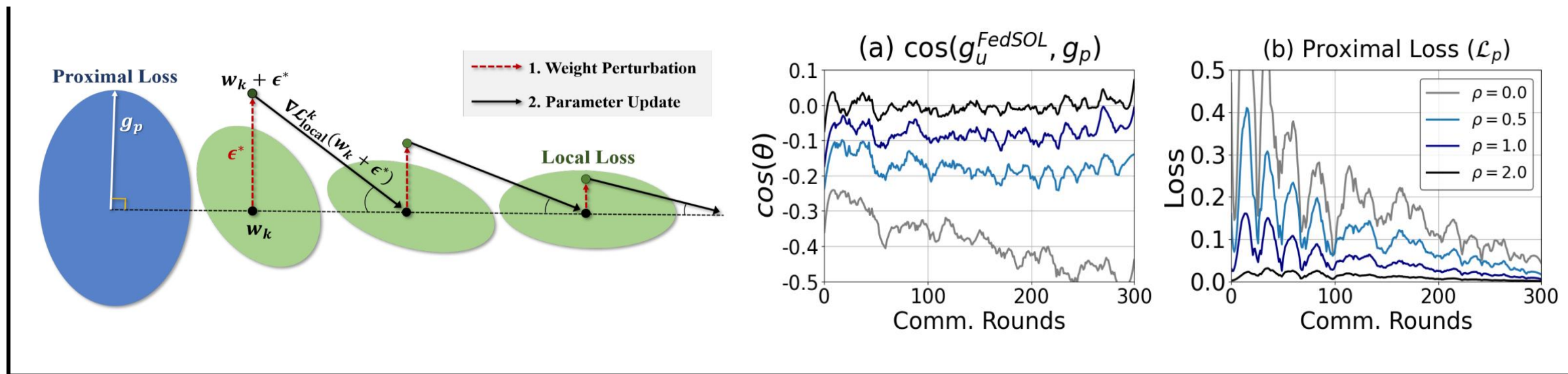
Regularize Loss Surface w.r.t. "Proximal Loss"



# 3. Main Approach - Stabilized Orthogonal Learning

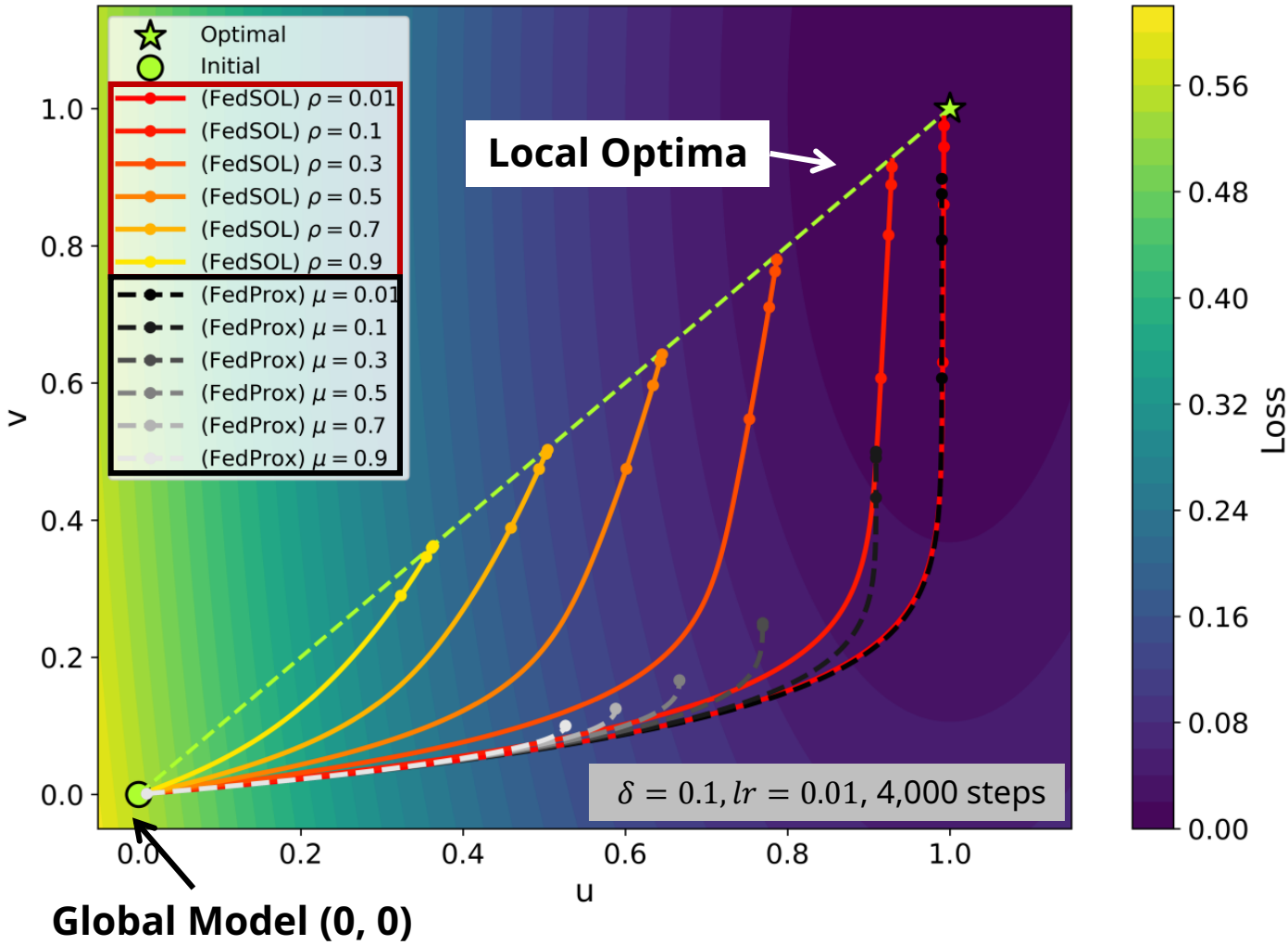
[Main Idea]: Identify the **local gradient** is minimally affected by **proximal gradient** finds the *orthogonal direction* of local knowledge acquisition.

## Experiment: CIFAR-10 ( $\alpha = 0.1$ )



# 3. Main Approach – Toy Example (2-dimensional)

**Optimization Path for FedProx and FedSOL**



**Local Loss** Distribution coefficient

$$\mathcal{L}_{local}(\mathbf{u}, \mathbf{v}) = \frac{1}{2}(\mathbf{u} - \mathbf{u}_l)^2 + \frac{\delta}{2}(\mathbf{v} - \mathbf{v}_l)$$

**Proximal Loss**

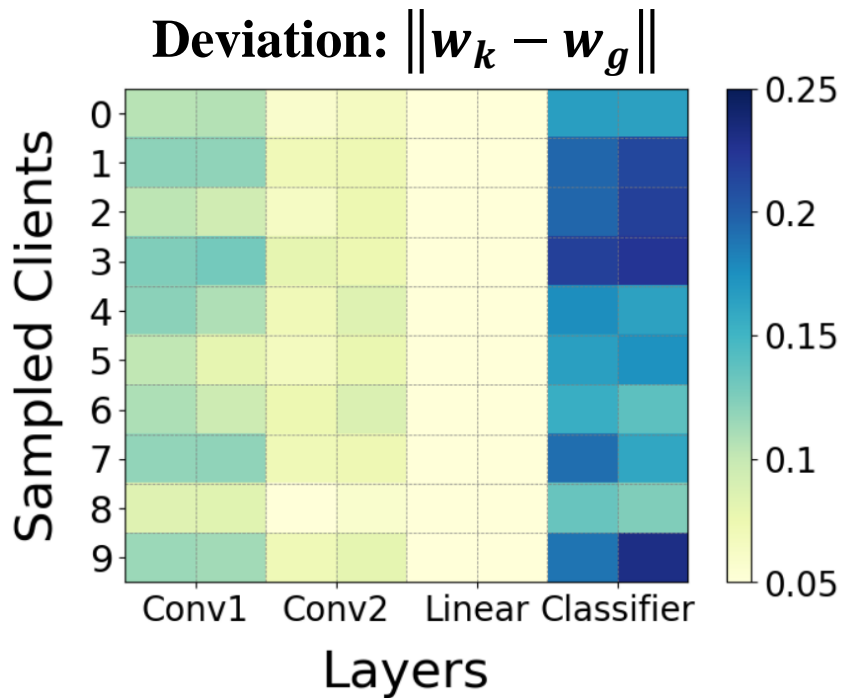
$$\mathcal{L}_p(\mathbf{u}, \mathbf{v}) = \frac{\mu}{2}(\mathbf{u}^2 + \mathbf{v}^2)$$

- **FedProx** results in *skewed* update that influenced by local distribution.
- **FedSOL** *aligns* parallel to the local optima regardless of  $\delta$ .

### 3. Main Approach – C. Partial Perturbation

**[Main Idea]: Perturbing only the last classifier layer is sufficient** for FedSOL. The performance reaches as high as full-model perturbation.

Experiment: CIFAR-10 ( $\alpha = 0.1$ )



Target Position	Perturbation ( $\rho$ )					FLOPs
	0.0	0.5	1.0	1.5	2.0	
All ( <i>full</i> )		61.17	<b>64.16</b>	<b>64.38</b>	<b>63.94</b>	$2\times +\delta$
Body ( <i>partial</i> )	56.13	60.98	62.95	63.94	63.80	$1.96\times +\delta$
Head ( <i>partial</i> )		<b>62.65</b>	63.62	64.13	63.25	$1.33\times +\delta$

- $\delta$  : Computation for the proximal loss.
- **FLOPs** shows relative computation w.r.t. **FedAvg**.

# 4. Main Experiment – Data Heterogeneity

## 4 Results

**Experiment: Test Accuracy @1 (%)**. The values in ( ) are the standard deviations.

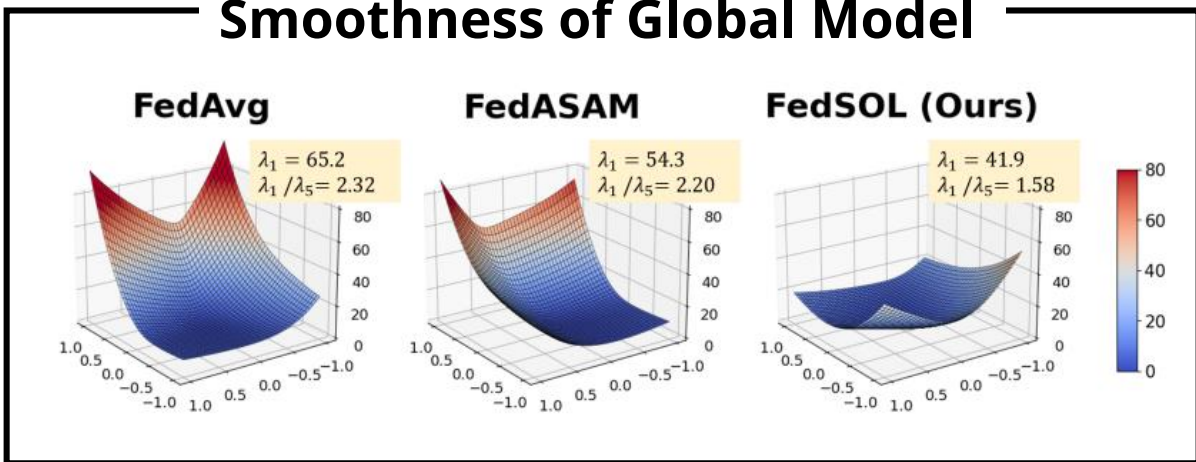
Method	MNIST	CIFAR-10				SVHN	CINIC-10	PathMNIST	TissueMNIST
		$\alpha = 0.05$	$\alpha = 0.1$	$\alpha = 0.3$	$\alpha = 0.5$				
FedAvg [46]	96.11 <sub>(0.19)</sub>	42.27 <sub>(1.34)</sub>	56.13 <sub>(0.78)</sub>	67.32 <sub>(0.94)</sub>	73.90 <sub>(0.66)</sub>	55.36 <sub>(4.85)</sub>	36.49 <sub>(4.37)</sub>	65.98 <sub>(4.76)</sub>	42.78 <sub>(2.03)</sub>
FedProx [38]	96.05 <sub>(0.13)</sub> ↓	50.58 <sub>(0.57)</sub> ↑	59.80 <sub>(1.12)</sub> ↑	68.39 <sub>(0.81)</sub> ↑	72.87 <sub>(0.55)</sub> ↑	72.40 <sub>(3.15)</sub> ↑	40.09 <sub>(3.97)</sub> ↑	70.44 <sub>(1.92)</sub> ↑	52.25 <sub>(1.40)</sub> ↑
FedNova [64]	88.24 <sub>(1.37)</sub> ↓	10.00 <sub>(Failed)</sub> ↓	10.00 <sub>(Failed)</sub> ↓	64.67 <sub>(0.77)</sub> ↓	70.04 <sub>(0.45)</sub> ↓	53.07 <sub>(3.30)</sub> ↓	21.89 <sub>(1.71)</sub> ↓	38.94 <sub>(2.34)</sub> ↓	15.03 <sub>(3.74)</sub> ↓
Scaffold [28]	94.18 <sub>(0.32)</sub> ↓	10.00 <sub>(Failed)</sub> ↓	10.00 <sub>(Failed)</sub> ↓	<b>71.92</b> <sub>(0.17)</sub> ↑	<b>75.49</b> <sub>(0.21)</sub> ↑	21.46 <sub>(1.75)</sub> ↓	16.89 <sub>(2.25)</sub> ↓	18.07 <sub>(0.04)</sub> ↓	32.04 <sub>(0.07)</sub> ↓
FedNTD [33]	96.97 <sub>(0.27)</sub> ↑	58.08 <sub>(0.48)</sub> ↑	<b>63.16</b> <sub>(1.02)</sub> ↑	71.56 <sub>(0.26)</sub> ↑	74.91 <sub>(0.33)</sub> ↑	79.25 <sub>(0.61)</sub> ↑	50.22 <sub>(3.71)</sub> ↑	74.26 <sub>(1.25)</sub> ↑	44.55 <sub>(1.95)</sub> ↑
FedSAM [55]	95.72 <sub>(0.43)</sub> ↓	36.14 <sub>(1.21)</sub> ↓	52.14 <sub>(0.94)</sub> ↓	64.83 <sub>(0.56)</sub> ↓	70.74 <sub>(0.40)</sub> ↓	13.27 <sub>(2.78)</sub> ↓	36.70 <sub>(4.28)</sub> ↑	66.64 <sub>(3.76)</sub> ↑	44.07 <sub>(3.02)</sub> ↑
FedASAM [6]	96.60 <sub>(0.10)</sub> ↑	43.12 <sub>(1.25)</sub> ↑	57.00 <sub>(0.30)</sub> ↑	67.45 <sub>(0.92)</sub> ↑	73.91 <sub>(0.51)</sub> ↑	60.25 <sub>(4.56)</sub> ↑	36.93 <sub>(4.60)</sub> ↑	69.45 <sub>(3.19)</sub> ↑	42.73 <sub>(2.35)</sub> ↑
<b>FedSOL (Ours)</b>	<b>97.44</b> <sub>(0.11)</sub> ↑	<b>60.01</b> <sub>(0.30)</sub> ↑	<b>64.13</b> <sub>(0.46)</sub> ↑	<b>71.94</b> <sub>(0.57)</sub> ↑	<b>75.60</b> <sub>(0.32)</sub> ↑	<b>83.92</b> <sub>(0.29)</sub> ↑	<b>55.07</b> <sub>(1.48)</sub> ↑	<b>78.88</b> <sub>(0.46)</sub> ↑	<b>53.40</b> <sub>(0.85)</sub> ↑

**FedSOL** achieves best results in most cases, showing its robustness against data heterogeneity.

# 4. Main Experiment – Analysis

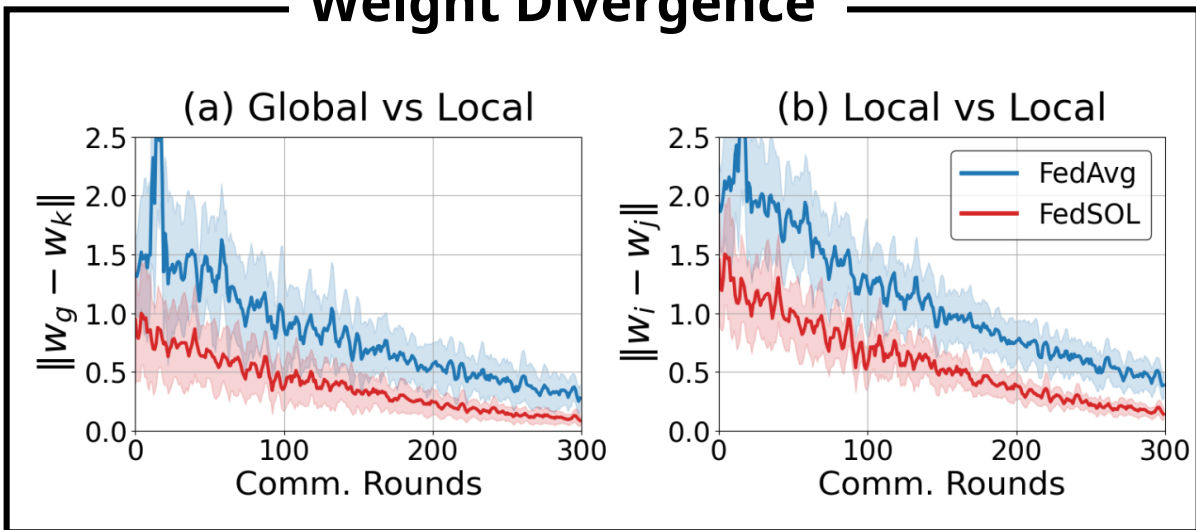
## 4 Results

### Smoothness of Global Model

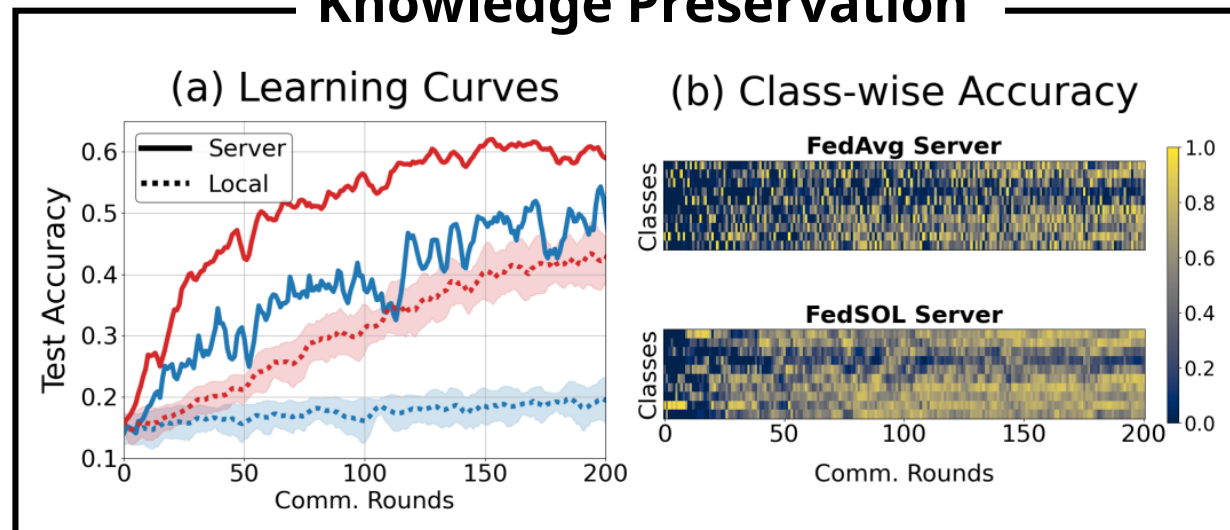


**FedSOL** stabilizes federated learning in various perspectives.

### Weight Divergence



### Knowledge Preservation



- FL must strike the balance between **global knowledge preservation** (Proximal Loss) and **local knowledge acquisition** (Original Local Loss).
- We suggest that **orthogonal learning in CL** could be an effective strategy in FL, by **resolving conflicts of local knowledge on global knowledge**.
- We propose **FedSOL**, which aims to obtain the **local gradient which is orthogonal to the proximal gradient**.
- **FedSOL** acquire local knowledge during local learning that less conflicts with the global knowledge.

**Thank You!**