



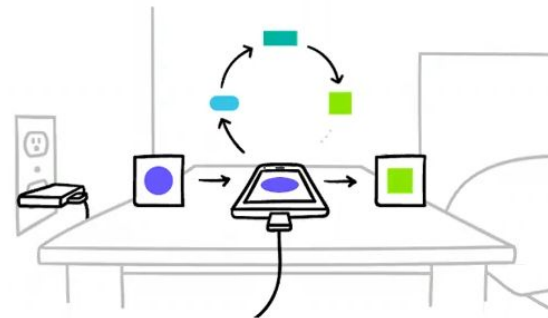
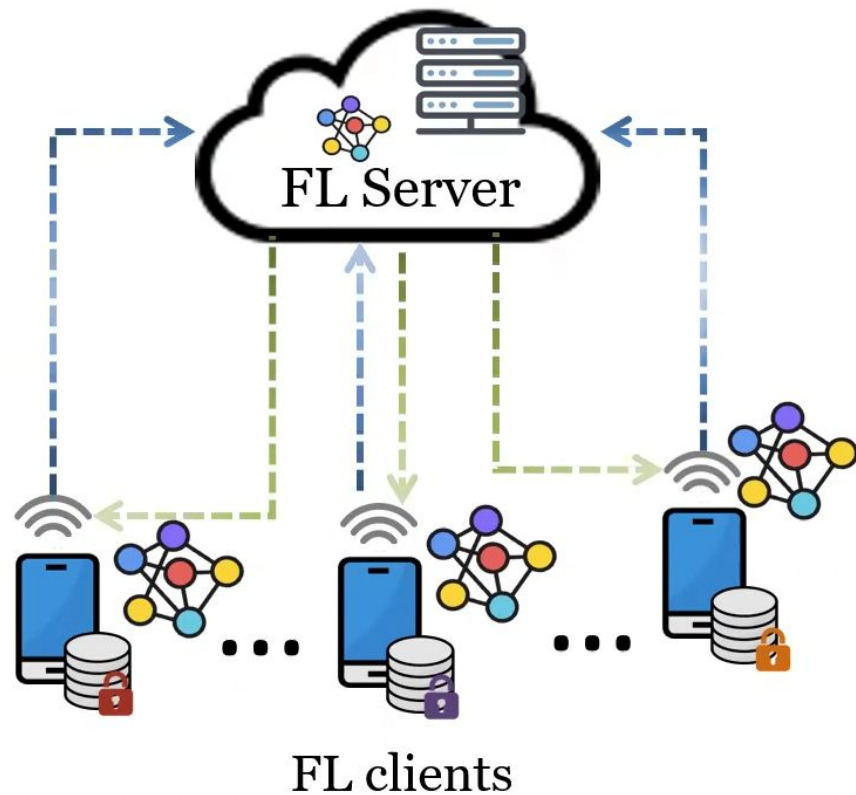
**Politecnico  
di Torino**

# Investigating Statistical Heterogeneity and Privacy Issues In Federated Learning

S277711 GIUSEPPE DESIDERIO  
S301905 VIDA AHMADI  
S303073 SHYAM NANDAN RAI

MLDL  
A.A 2021/2022

# Introduction: Federated Learning



- 🔒 Data never leaves local devices
- 🔒 Learn on fresh real-world data

# Content

Our project investigates two major problems in federated learning:

Our project investigates two major problems in federated learning:

## 1. **Statistical Heterogeneity**

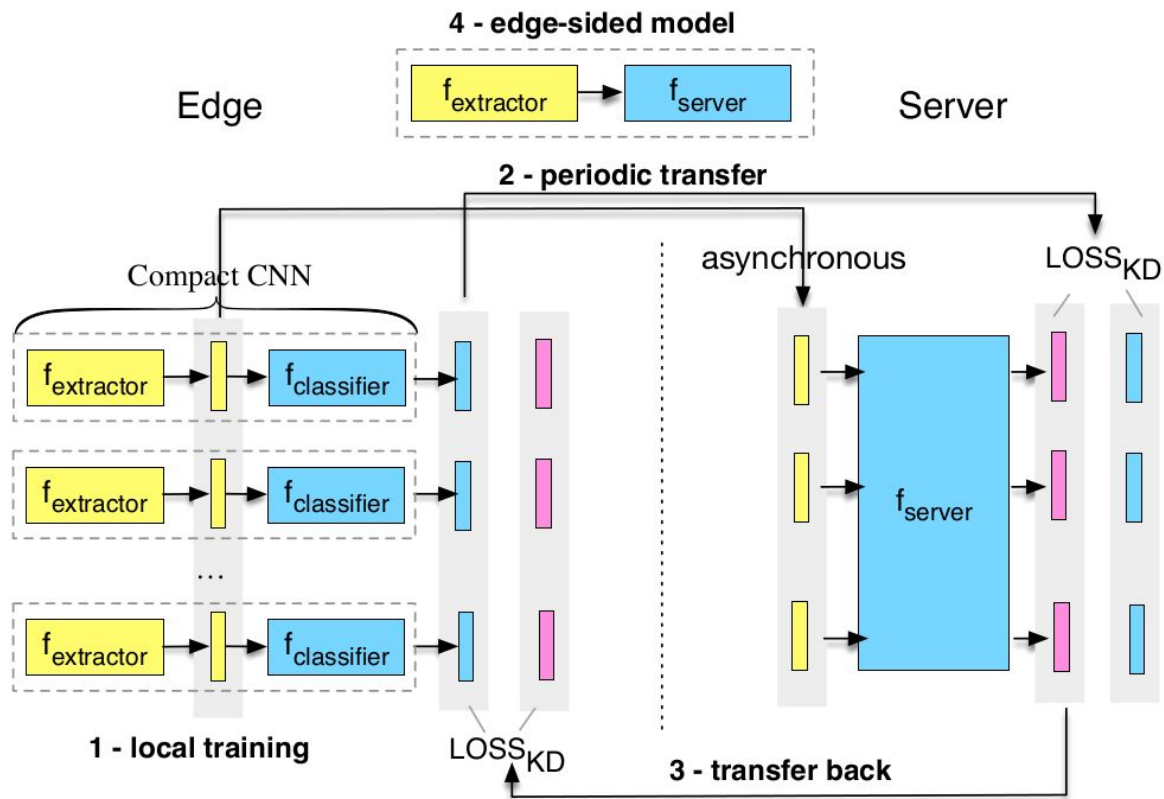
- a. We benchmarked the performance of the federated learning methods: FedAvg and FedGKT in IID and non-IID setting.
- b. Comparing Group Normalization vs Batch Normalization the former shows better performance in the presence of heterogeneity.
- c. We propose MOON-Prox that addresses the problem of data-heterogeneity by introducing a novel loss function.

## 2. **Privacy Concerns**

- a. We show that it is possible to reconstruct client data using gradient inversion attack.
- b. Extensive experiments on FedSGD and FedAvg are conducted showing client data reconstruction.

- **Client**  
receives the global model from server.  
train the model for  $E$  local epochs.  
sends the updated model to the server.
- **Server**  
it takes a weighted average of the resulting models.
- **the amount of computation is controlled by three key parameters :**  
 $C$ , the fraction of clients that perform computation on each round.  
 $E$ , the number of training passes each client makes over its local dataset on each round.  
 $B$ , the local minibatch size used for the client updates.

- FedGKT tackles the resource-constrained problem on edge devices.

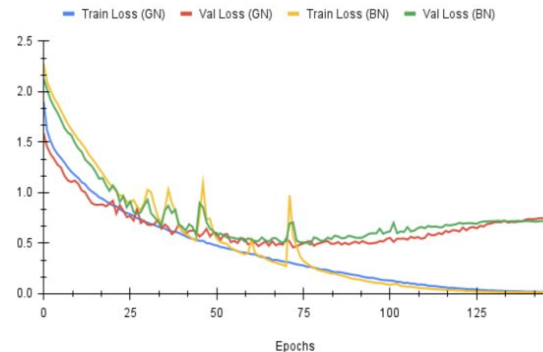


# Experimental Setting

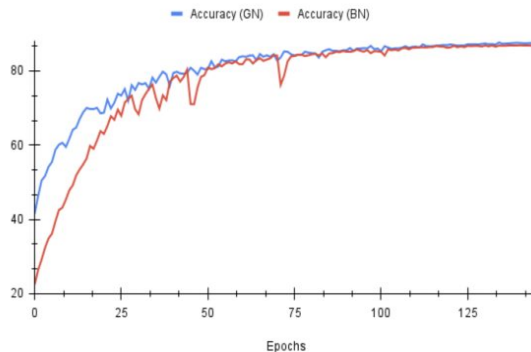
- **Architectures:**
  - FedAvg, Centralized Baseline: ResNet50
  - FedGKT: Resnet-8 (Client), Resnet-49(Server)
- **Normalization Layers:** Group Norm, Batch Norm
- **Dataset:** CIFAR-10
- **Baselines :** FedAvg, FedGKT, Centralized Baseline
- **Non-IID distribution:** Dirichlet [ $\beta=0.5$ ].

# Results: Centralized Baseline

Loss Plots



Accuracy Plots

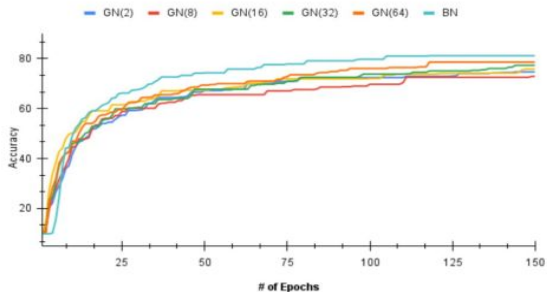


Method	Model	Norm.	Accuracy
Centralized	ResNet-50	BN	86.91
Centralized	ResNet-50	GN	87.48
FedAVG[IID]	ResNet-50	BN	81.14
FedAVG[IID]	ResNet-50	GN	78.60
FedAVG[nonIID]	ResNet-50	BN	60.00
FedAVG[nonIID]	ResNet-50	GN	60.51
FedGKT[IID]	ResNet-49	BN	55.28
FedGKT[IID]	ResNet-49	GN	53.27
FedGKT[nonIID]	ResNet-49	BN	27.01
FedGKT[nonIID]	ResNet-49	GN	19.79

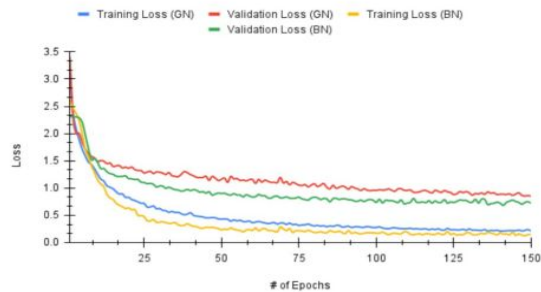


# Results: FedAvg

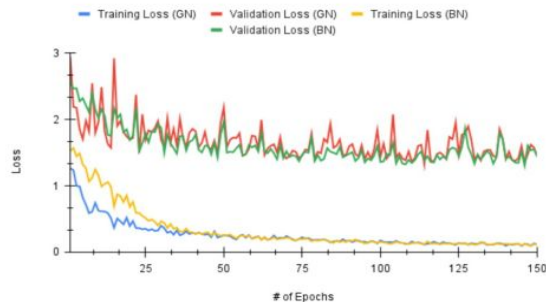
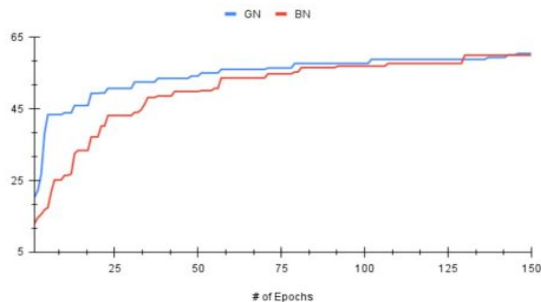
## Accuracy Plots



## Loss Plots



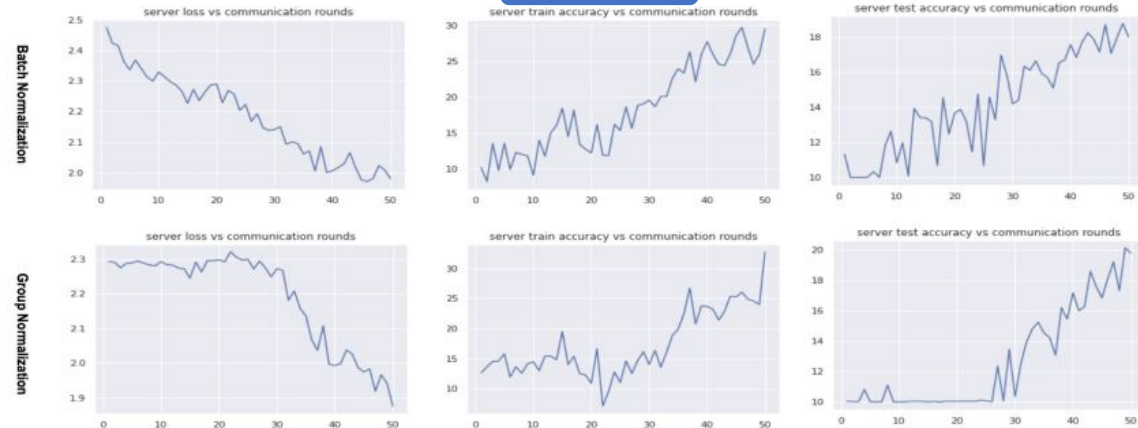
non-IID



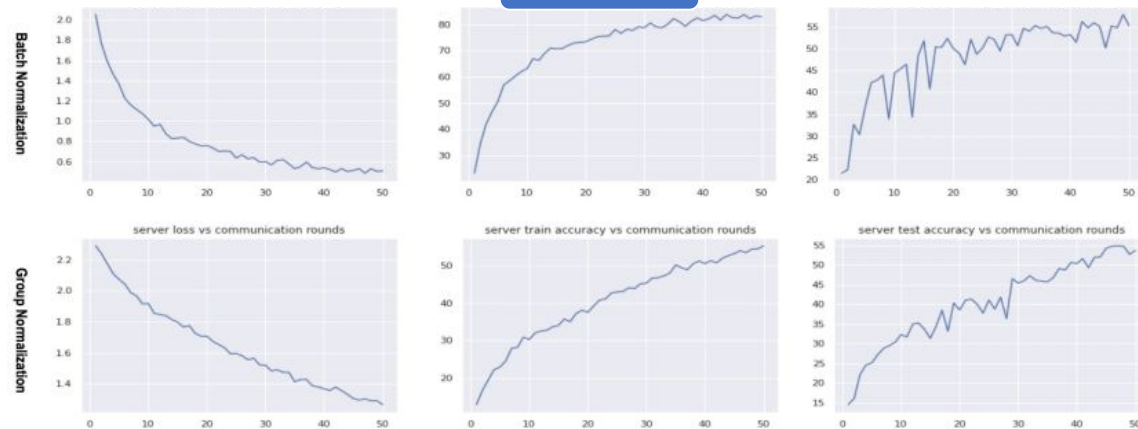
Method	Model	Norm.	Accuracy
Centralized	ResNet-50	BN	86.91
Centralized	ResNet-50	GN	87.48
FedAVG[IID]	ResNet-50	BN	81.14
FedAVG[IID]	ResNet-50	GN	78.60
FedAVG[nonIID]	ResNet-50	BN	60.00
FedAVG[nonIID]	ResNet-50	GN	60.51
FedGKT[IID]	ResNet-49	BN	55.28
FedGKT[IID]	ResNet-49	GN	53.27
FedGKT[nonIID]	ResNet-49	BN	27.01
FedGKT[nonIID]	ResNet-49	GN	19.79

# Results: FedGKT

non-IID



IID

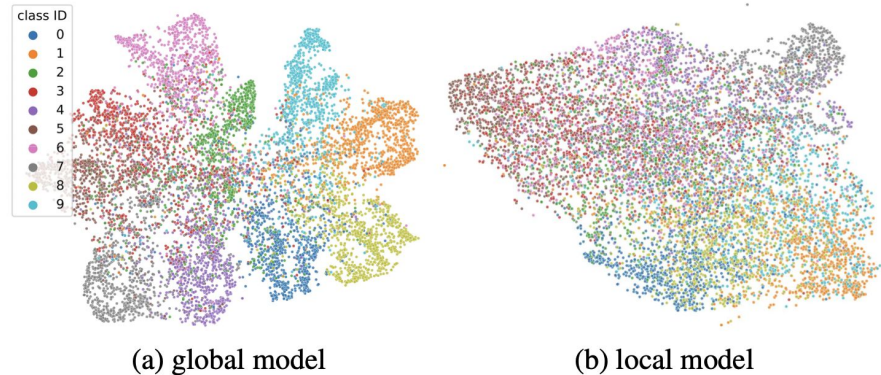


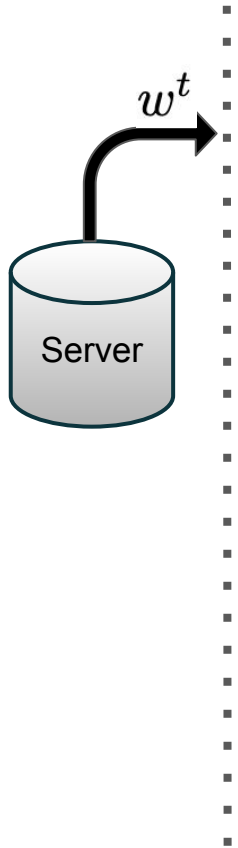
Method	Model	Norm.	Accuracy
Centralized	ResNet-50	BN	86.91
Centralized	ResNet-50	GN	87.48
FedAVG[IID]	ResNet-50	BN	81.14
FedAVG[IID]	ResNet-50	GN	78.60
FedAVG[nonIID]	ResNet-50	BN	60.00
FedAVG[nonIID]	ResNet-50	GN	60.51
FedGKT[IID]	ResNet-49	BN	55.28
FedGKT[IID]	ResNet-49	GN	53.27
FedGKT[nonIID]	ResNet-49	BN	27.01
FedGKT[nonIID]	ResNet-49	GN	19.79

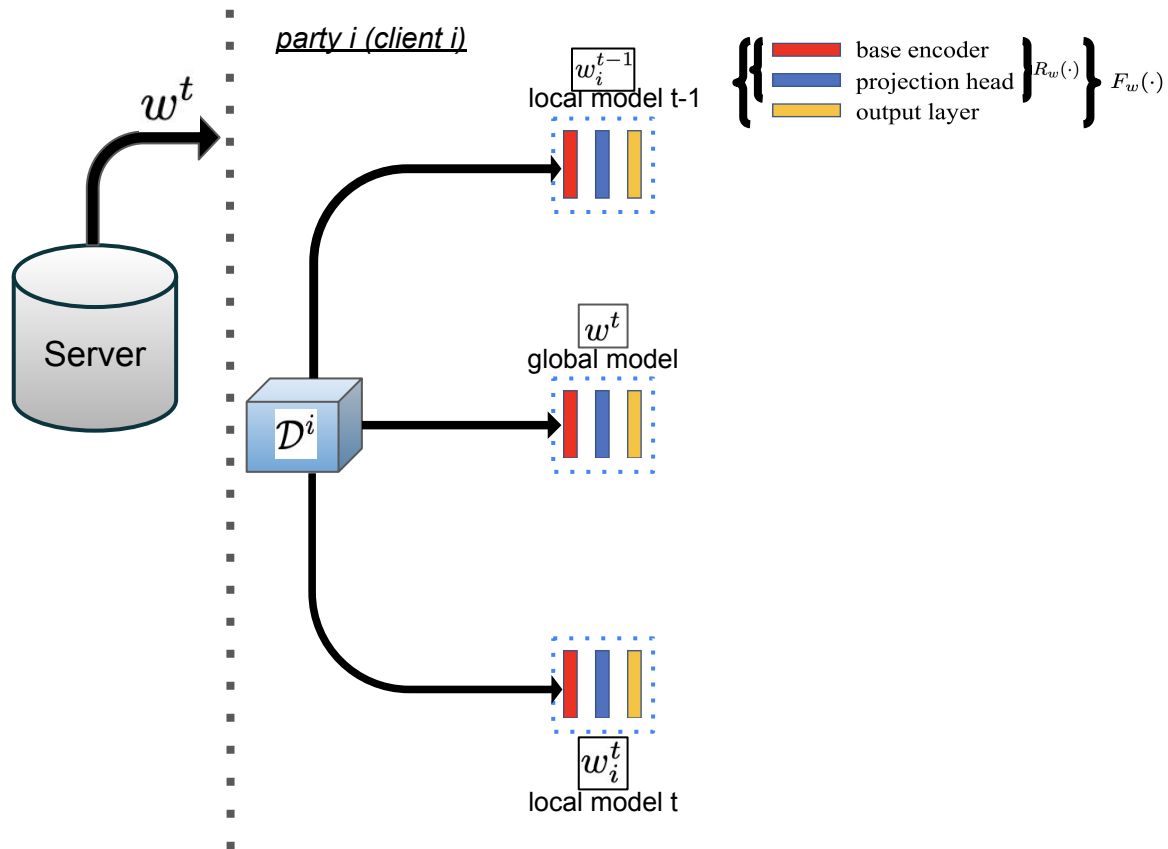
# MOON: Motivation

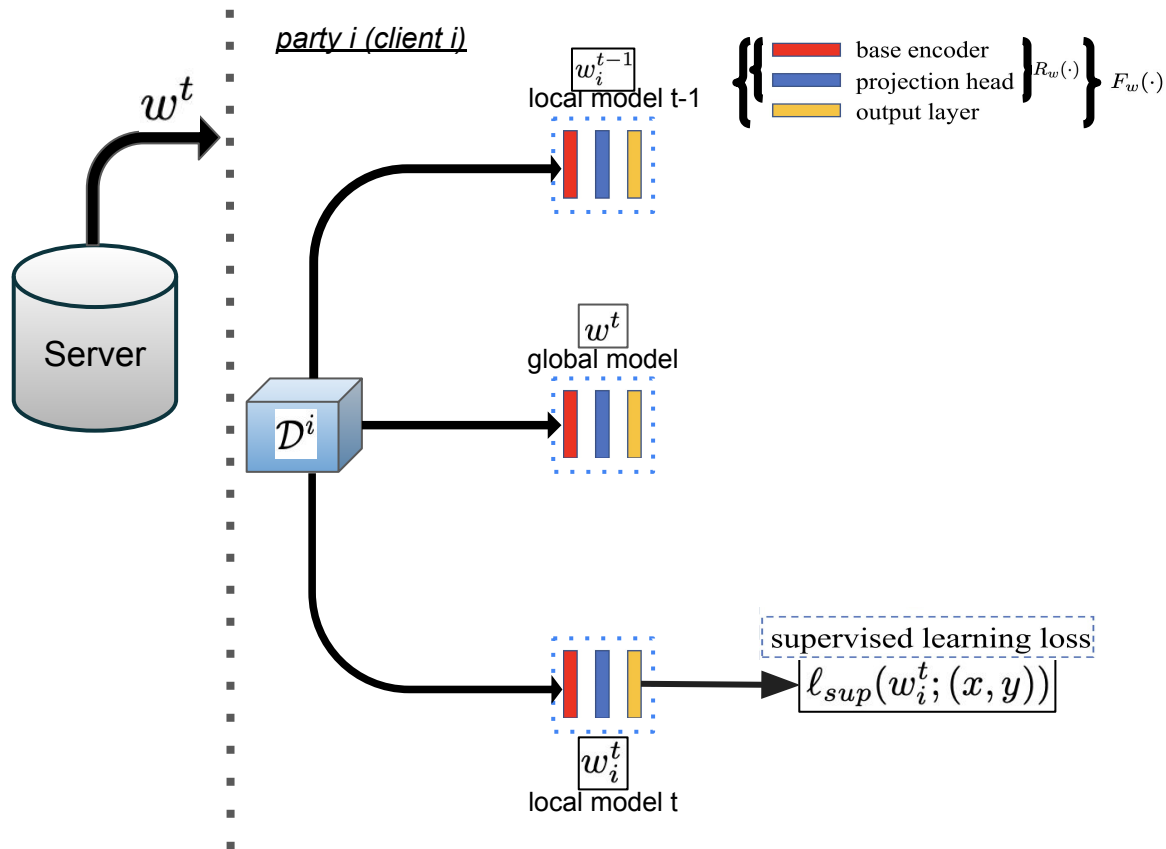
**Intuition:** Model trained on the whole dataset **[Global Model]** extracts better feature representation than model trained on skewed subset **[Local Model]**.

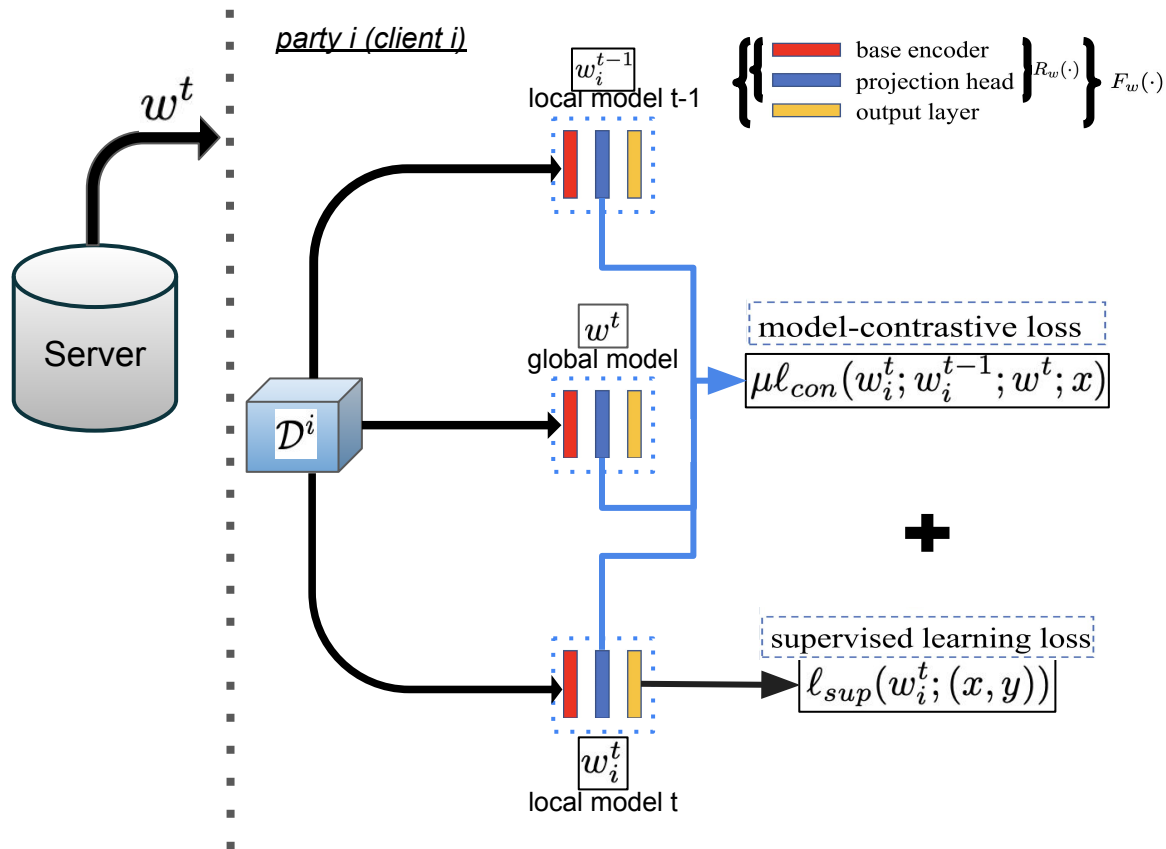
- Bridge the gap between the representation of local model and global model. Minimize the drift caused in the local model updates.
- Contrastive Learning:  $\downarrow$  global and local model distance &  $\uparrow$  local model (t) and local model (t-1).

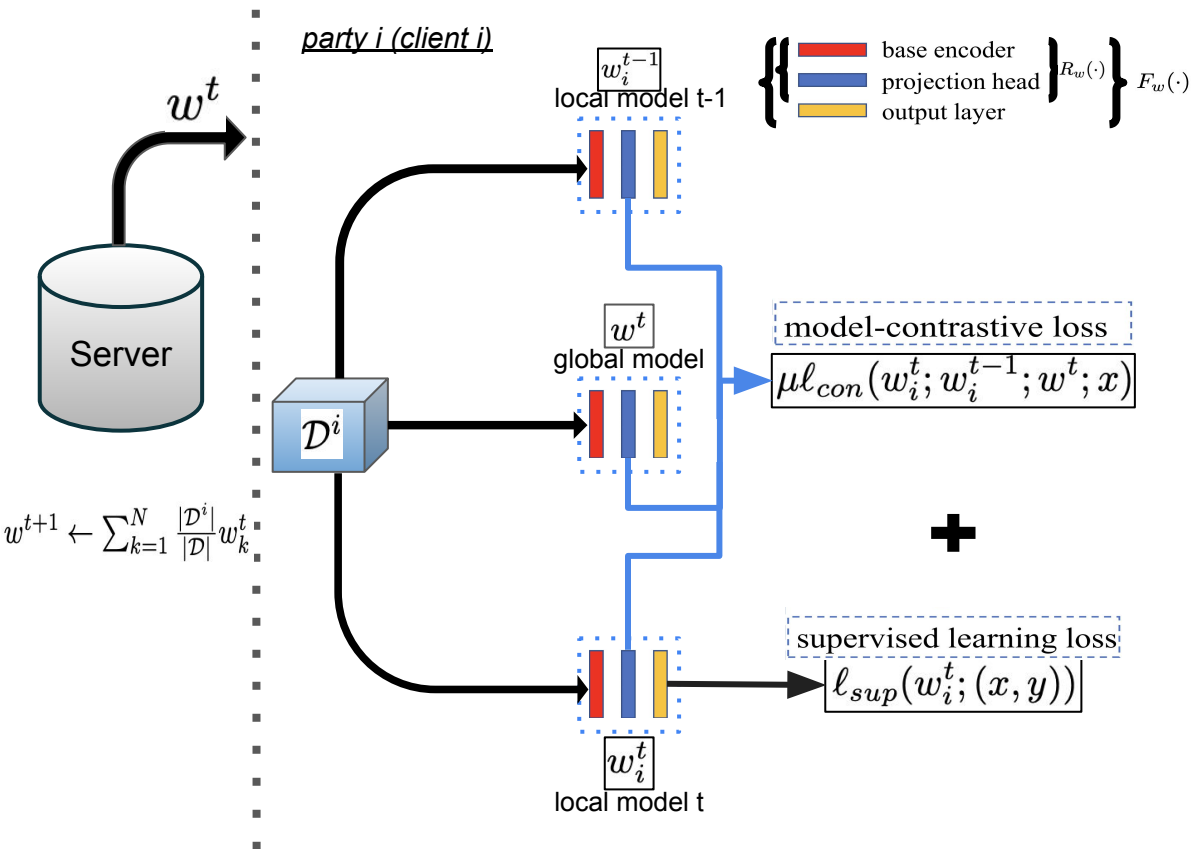














# MOON-Results

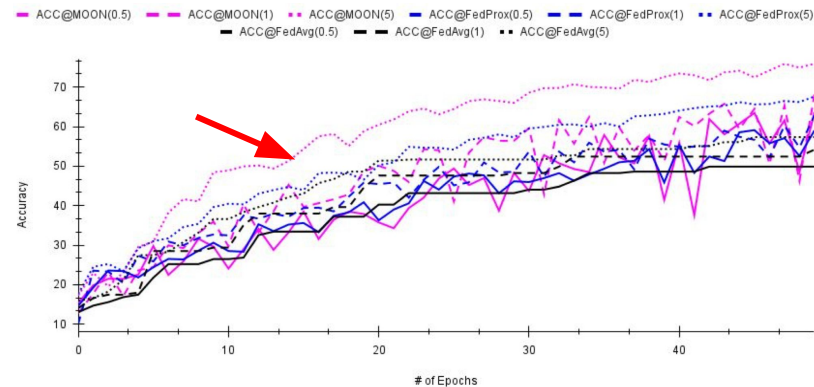
**Dataset:** CIFAR-10

**Network:** ResNet-18

Accuracy Comparison

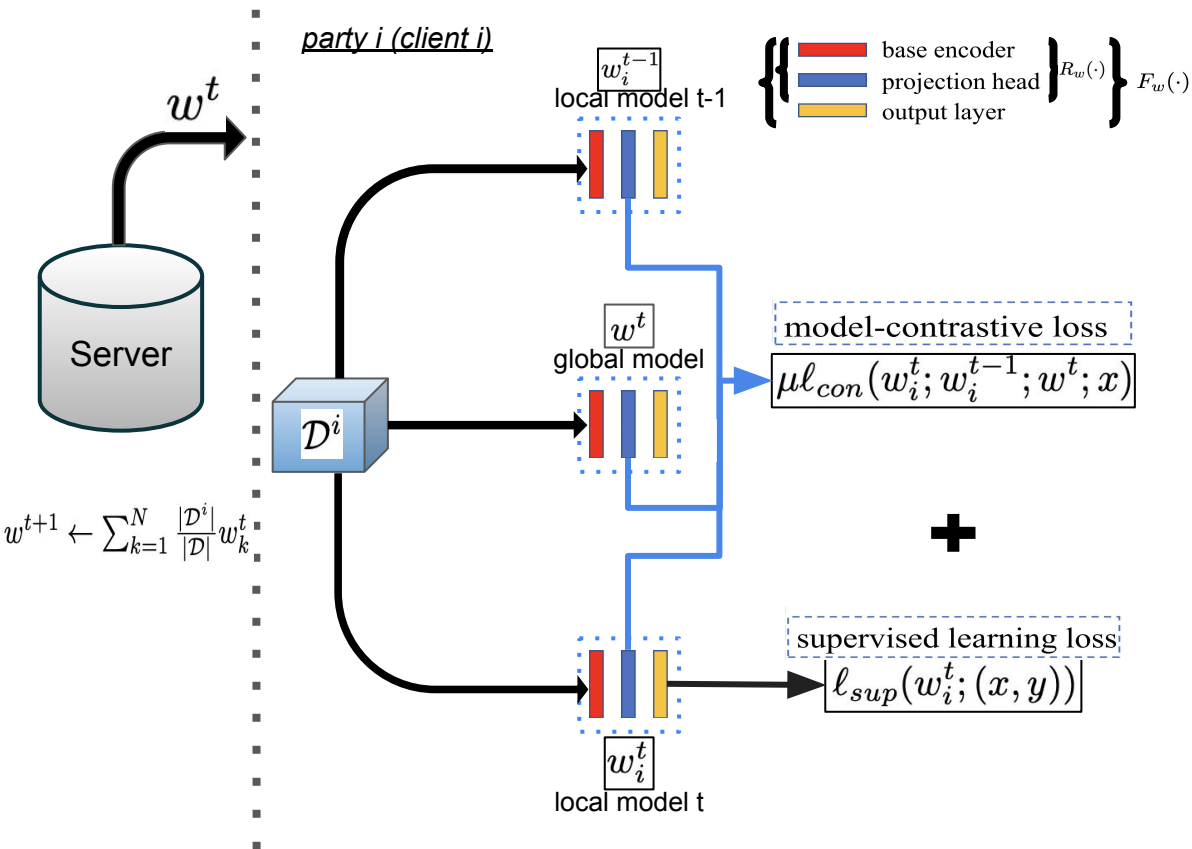
Method	Acc.@ $\beta=0.5$	Acc.@ $\beta=1$	Acc.@ $\beta=5$
FedAvg	49.87	54.17	57.34
FedProx	59.10	62.88	67.61
MOON	63.62	68.23	75.93

Accuracy Plots



Speed-up Comparison

Method	Rounds@ $\beta = 0.5$	Speedup@ $\beta = 0.5$	Rounds@ $\beta = 1$	Speedup@ $\beta = 1$	Rounds@ $\beta = 5$	Speedup@ $\beta = 5$
FedAvg	50	1x	50	1x	50	1x
FedProx	35	1.42x	34	1.47x	27	1.85x
MOON	25	2x	23	2.17x	16	3.12x



Maximize

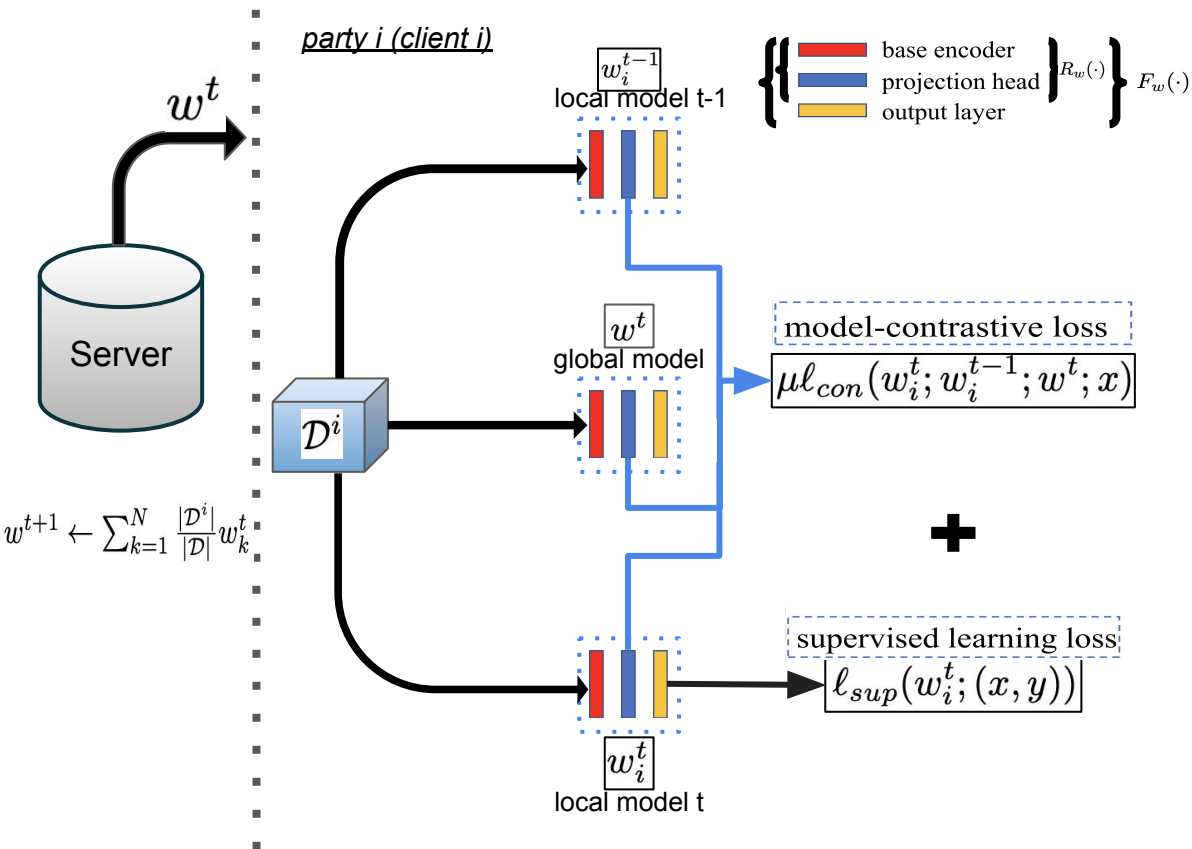
Feature Alignment



Weight Alignment



# Problem



Maximize

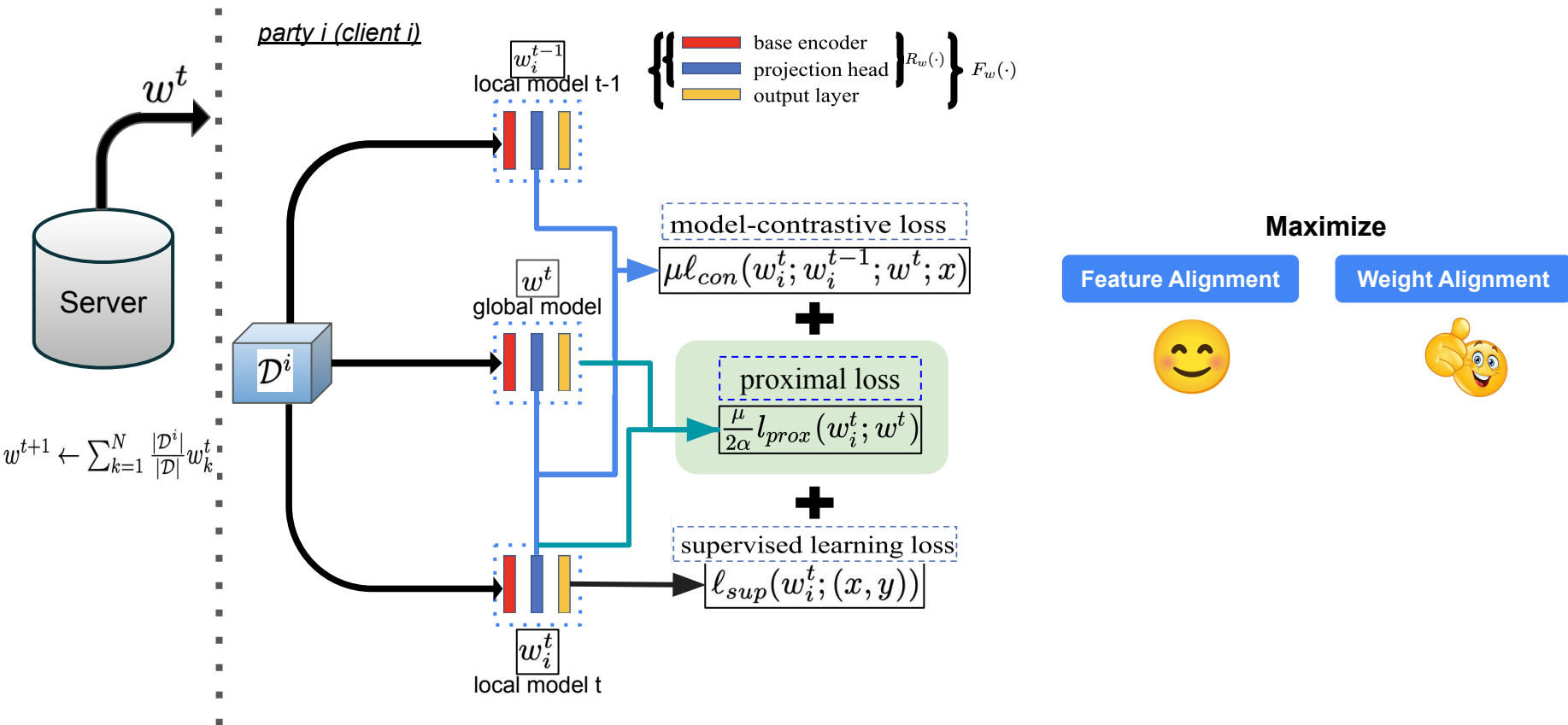
Feature Alignment

Weight Alignment

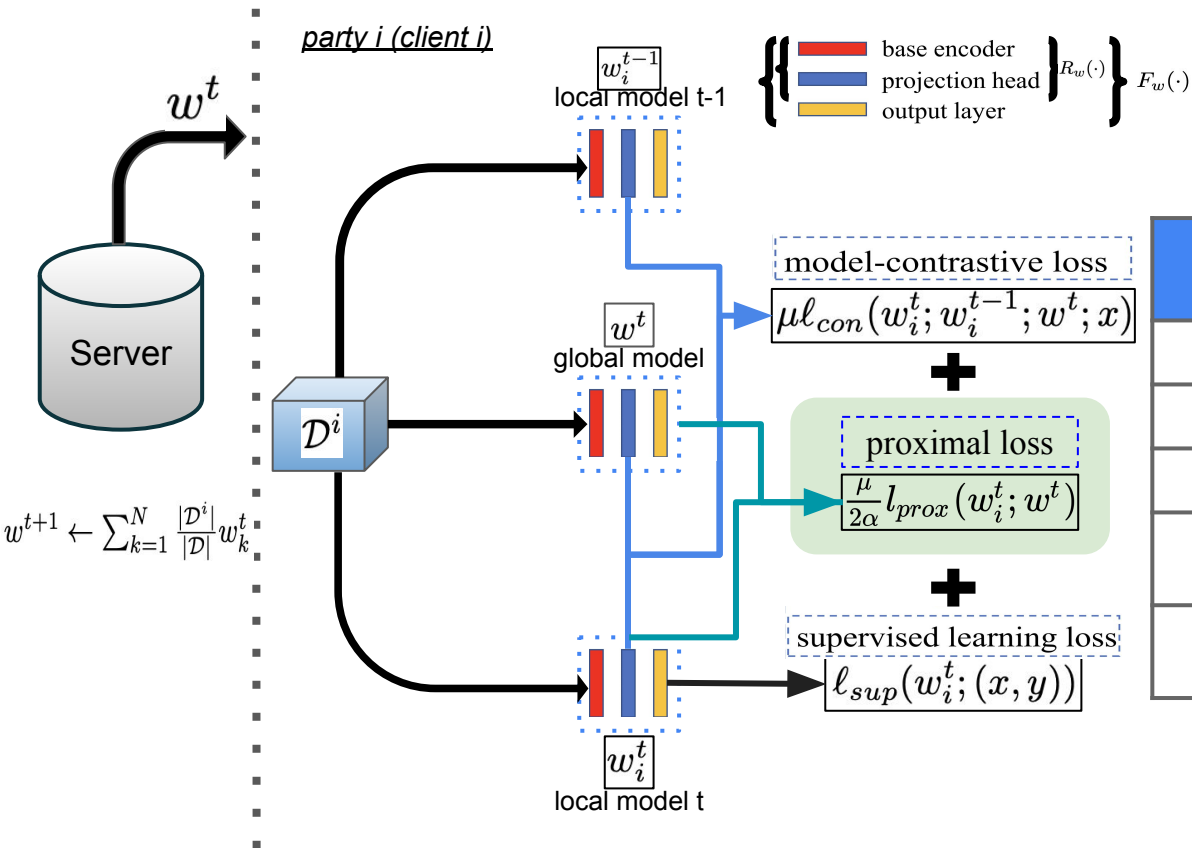


- Large weight updates
- Unstable Training

# MOON-Prox



# MOON-Prox



Method	Acc@CR=50	Acc@CR=100
<i>FedAvg</i>	49.87	56.96
<i>FedProx</i>	59.10	65.21
<i>MOON</i>	63.62	77.03
<b>Ours</b> (alpha=100)	66.50 (+2.88)	77.25
<b>Ours</b> (alpha=500)	63.16	77.59 (+0.56)

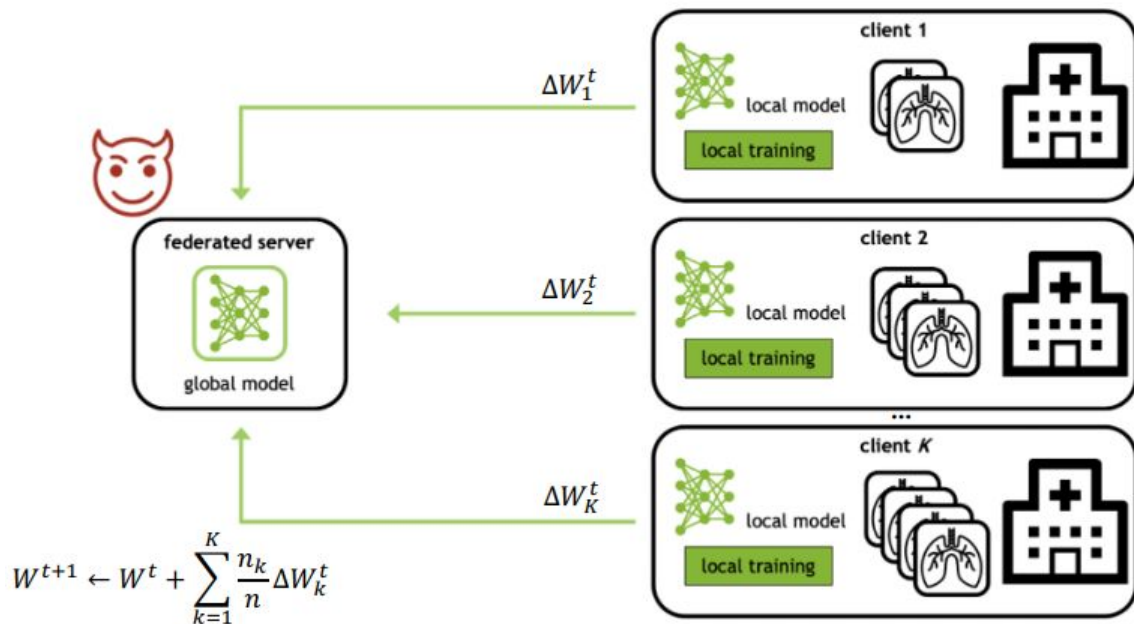
# Privacy Concerns: Gradient Inversion Attack

- Might federated learning give users a false sense of privacy?

# Privacy Concerns: Gradient Inversion Attack

- Might federated learning give users a false sense of privacy?

## Threat model: Honest-but-Curious Server



✓ Allow to separately store and process updates transmitted by individual users.

✗ Not interfere with the collaborative learning algorithm.

✗ Not modify the model architecture to better suit their attack.

✗ Cannot send malicious global parameters that do not represent the actually learned global model.

# Privacy Concerns: Gradient Inversion Attack

- Optimization problem:

$$\arg \min_x \mathcal{L}_{grad}(x; \theta, \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*)) + \alpha \mathcal{R}_{aux}(x)$$

- $\theta$  = neural network parameters
- $\nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*)$  = the gradient computed with a private data batch  $(x^*, y^*) \in \mathcal{R}^{b \times d} \times \mathcal{R}^b$  (b, d being the batch size, image size)



# Privacy Concerns: Gradient Inversion Attack

- Optimization problem:

$$\arg \min_x \mathcal{L}_{grad}(x; \theta, \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*)) + \alpha \mathcal{R}_{aux}(x)$$

- $\theta$  = neural network parameters
- $\nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*)$  = the gradient computed with a private data batch  $(x^*, y^*) \in \mathcal{R}^{b \times d} \times \mathcal{R}^b$  (b, d being the batch size, image size)

- $\mathcal{L}_{grad}(x; \theta, \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*))$  enforces to produce a gradient that matches the observed model updates.

# Privacy Concerns: Gradient Inversion Attack

- Optimization problem:

$$\arg \min_x \mathcal{L}_{grad}(x; \theta, \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*)) + \alpha \mathcal{R}_{aux}(x)$$

- $\theta$  = neural network parameters
- $\nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*)$  = the gradient computed with a private data batch  $(x^*, y^*) \in \mathcal{R}^{b \times d} \times \mathcal{R}^b$  (b, d being the batch size, image size)

- $\mathcal{L}_{grad}(x; \theta, \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*))$  enforces to produce a gradient that matches the observed model updates.
- $\mathcal{R}_{aux}(x)$  regularizes the recovered image based on image prior(s).

# Privacy Concerns: Gradient Inversion Attack

- Optimization problem:

$$\arg \min_x \mathcal{L}_{grad}(x; \theta, \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*)) + \alpha \mathcal{R}_{aux}(x)$$

- $\theta$  = neural network parameters
- $\nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*)$  = the gradient computed with a private data batch  $(x^*, y^*) \in \mathcal{R}^{b \times d} \times \mathcal{R}^b$  (b, d being the batch size, image size)

- $\mathcal{L}_{grad}(x; \theta, \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*))$  enforces to produce a gradient that matches the observed model updates.
- $\mathcal{R}_{aux}(x)$  regularizes the recovered image based on image prior(s).

- Jonas Geiping implementation from [1]:

$$\arg \min_x 1 - \frac{\langle \nabla_{\theta} \mathcal{L}_{\theta}(x, y), \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*) \rangle}{\|\nabla_{\theta} \mathcal{L}_{\theta}(x, y)\| \|\nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*)\|} + \alpha_{TV} \mathcal{R}_{TV}(x)$$

- $\mathcal{L}_{grad}(x; \theta, \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*))$  = cosine similarity loss.
- $\mathcal{R}_{TV}(x)$  = total variation of images.

# Privacy Concerns: Gradient Inversion Attack

- Optimization problem:

$$\arg \min_x \mathcal{L}_{grad}(x; \theta, \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*)) + \alpha \mathcal{R}_{aux}(x)$$

- $\theta$  = neural network parameters
- $\nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*)$  = the gradient computed with a private data batch  $(x^*, y^*) \in \mathcal{R}^{b \times d} \times \mathcal{R}^b$  (b, d being the batch size, image size)

- $\mathcal{L}_{grad}(x; \theta, \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*))$  enforces to produce a gradient that matches the observed model updates.
- $\mathcal{R}_{aux}(x)$  regularizes the recovered image based on image prior(s).

- Jonas Geiping implementation from [1]:

$$\arg \min_x 1 - \frac{\langle \nabla_{\theta} \mathcal{L}_{\theta}(x, y), \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*) \rangle}{\|\nabla_{\theta} \mathcal{L}_{\theta}(x, y)\| \|\nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*)\|} + \alpha_{TV} \mathcal{R}_{TV}(x)$$

- $\mathcal{L}_{grad}(x; \theta, \nabla_{\theta} \mathcal{L}_{\theta}(x^*, y^*))$  = cosine similarity loss.
- $\mathcal{R}_{TV}(x)$  = total variation of images.

**Assumptions:**

- Batch normalization statistics ✓
- Private labels ✓

# Experimental Setup

- **Architectures:** ResNet-18, ResNet-50.
- **Normalization Layers:** Group Norm, Batch Norm.
- **Algorithms:** FedAvg, FedSGD.
- **Dataset:** CIFAR-10.
- **Metrics:** LPIPS, PSNR.
- **Repository:** Breaching at [2]

# Experimental Results: BatchNorm Case

ResNet-18 Batch Norm

Original



Reconstructed



LPIPS: 0.03  
PSNR: 20.43

**FedSGD**

LPIPS: 0.02  
PSNR: 21.08

**FedSGD**

LPIPS: 0.04  
PSNR: 20.77

**FedAVG**

ResNet-50 Batch Norm



LPIPS: 0.17  
PSNR: 13.49

**FedAVG**

LPIPS: 0.15  
PSNR: 16.62

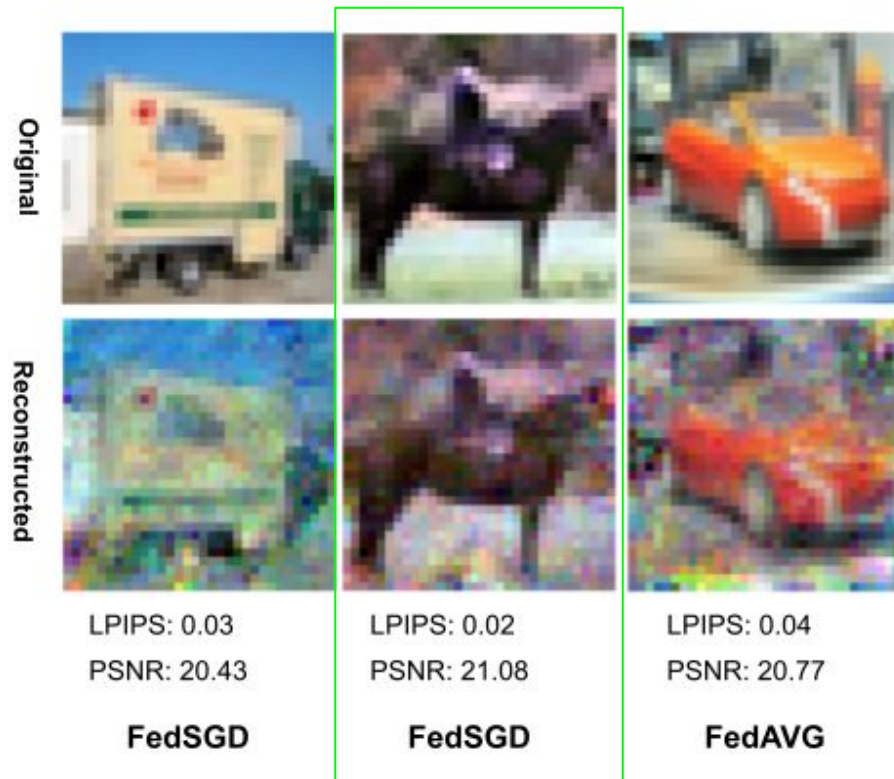
**FedSGD**

LPIPS: 0.09  
PSNR: 19.38

**FedSGD**

# Experimental Results: BatchNorm Case

ResNet-18 Batch Norm



ResNet-50 Batch Norm





# Experimental Results: BatchNorm Case

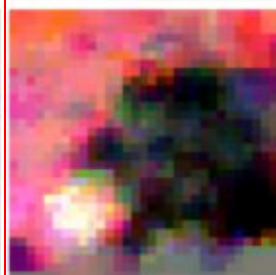
ResNet-18 Batch Norm

ResNet-50 Batch Norm

Original



Reconstructed



LPIPS: 0.03  
PSNR: 20.43

LPIPS: 0.02  
PSNR: 21.08

LPIPS: 0.04  
PSNR: 20.77

LPIPS: 0.17  
PSNR: 13.49

LPIPS: 0.15  
PSNR: 16.62

LPIPS: 0.09  
PSNR: 19.38

**FedSGD**

**FedSGD**

**FedAVG**

**FedAVG**

**FedSGD**

**FedSGD**



# Experimental Results: GroupNorm Case

ResNet-18 Group Norm

Original



Reconstructed



LPIPS : 0.02  
PSNR : 20.17

**FedSGD**



LPIPS : 0.06  
PSNR : 20.77

**FedAVG**



LPIPS : 0.02  
PSNR : 21.35

**FedSGD**

ResNet-50 Group Norm

Original

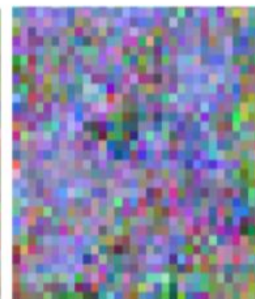


Reconstructed



LPIPS : 0.10  
PSNR : 13.80

**FedSGD**



LPIPS : 0.21  
PSNR : 15.60

**FedAVG**



LPIPS : 0.06  
PSNR : 19.50

**FedSGD**

# Conclusion

- We benchmarked recent federated learning algorithms: FedAvg, FedGKT and MOON and assessed their performance in IID and non-IID data distribution.
- We conclude that GN gives better performance than BN for non-IID data distribution.
- We proposed MOON-Prox that performs best among all the compared federated learning algorithms.
- We investigate security issues faced by federated learning methods via gradient attack and were able to reconstruct the private data of the client.

**Thanks for your attention**

Method	Model	Norm.	Accuracy
Centralized	ResNet-50	BN	86.91
Centralized	ResNet-50	GN	87.48
FedAVG[IID]	ResNet-50	BN	81.14
FedAVG[IID]	ResNet-50	GN	78.60
FedAVG[nonIID]	ResNet-50	BN	60.00
FedAVG[nonIID]	ResNet-50	GN	60.51
FedGKT[IID]	ResNet-49	BN	55.28
FedGKT[IID]	ResNet-49	GN	53.27
FedGKT[nonIID]	ResNet-49	BN	27.01
FedGKT[nonIID]	ResNet-49	GN	19.79
FedGKT*[IID]	ResNet-49	BN	74.55
FedGKT*[IID]	ResNet-49	GN	74.28
FedGKT*[nonIID]	ResNet-49	BN	47.53
FedGKT*[nonIID]	ResNet-49	GN	48.58