



# Instance-wise Batch Label Restoration via Gradients in Federated Learning

Kailang Ma\*, Yu Sun\*, Jian Cui, Dawei Li, Zhenyu Guan, Jianwei Liu

School of Cyber Science and Technology,  
Beihang University,  
China



# Background & Motivation

- Federated Learning: a distributed learning paradigm, asking participants to exchange **model updates or gradients** instead of raw data.
- **Breaking Privacy! Gradient-Matching strategy** (Gradient Inversion Attack, GIA) allows for reverse-engineer private input from shared gradients. [[Zhu et al' 19](#), [Geiping et al '20...](#)]
  - **Label Restoration remains a challenge.**



Original



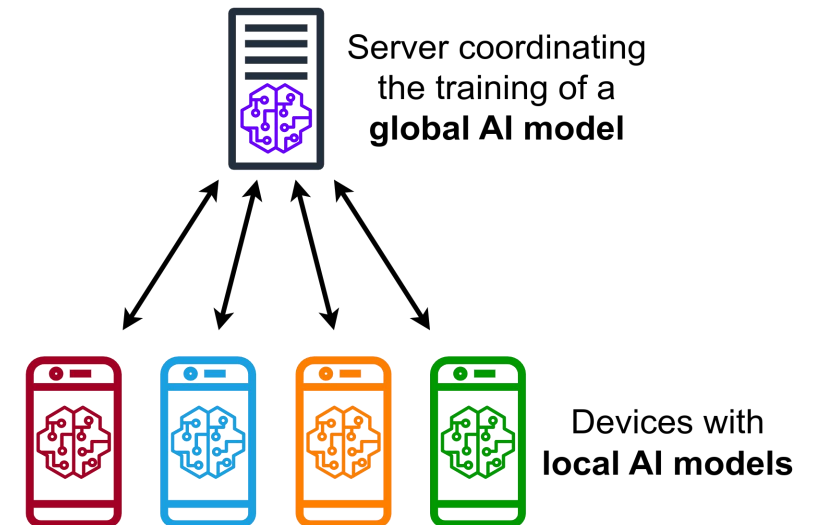
[Zhu et al' 19](#)



[Geiping et al '20](#)



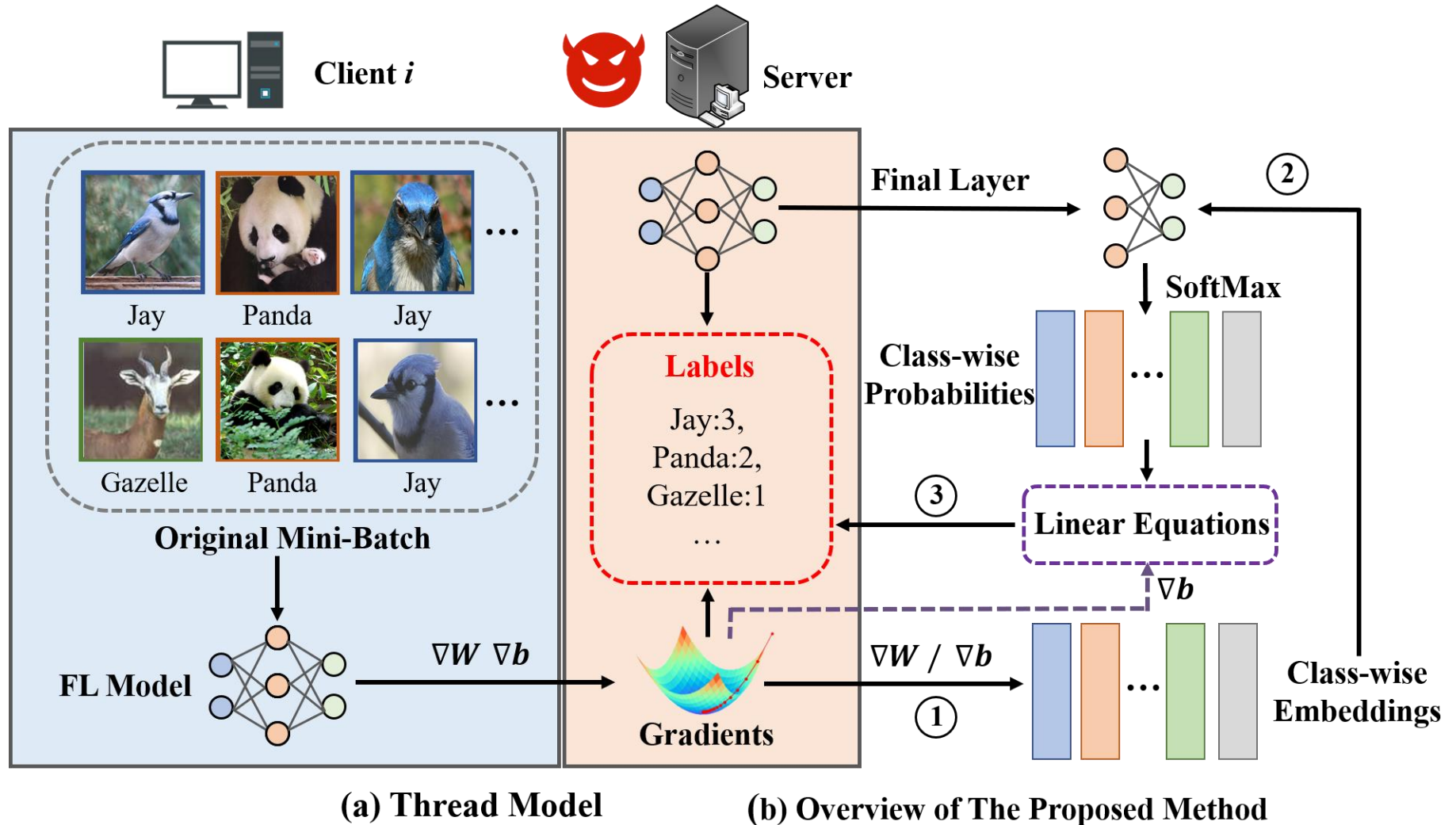
[Yin et al '21](#)



# Prior Works

- The optimization-based methods do not perform well [[Zhu et al' 19](#)].
  - It is possible to extract labels totally by **analytic observation or computation**.
    - **Signs of Gradients** [[Zhao et al' 20](#), for a single sample; [Yin et al '21](#), from batch-averaged gradients]. (Requires Non-negative Activation Functions).
    - **Find a classifier to separate  $q^c$**  (gradient column *w.r.t.* the ground truth label  $c$ ) from  $q^{j \neq c}$  by linear programming [[Dang et al'21](#)].
    - Existing methods remains limited to identify the presence of categories (i.e., **class-wise label restoration**).
- **Not only Category but also Frequency (i.e., instance-wise label restoration)**

# Overview of iLRG

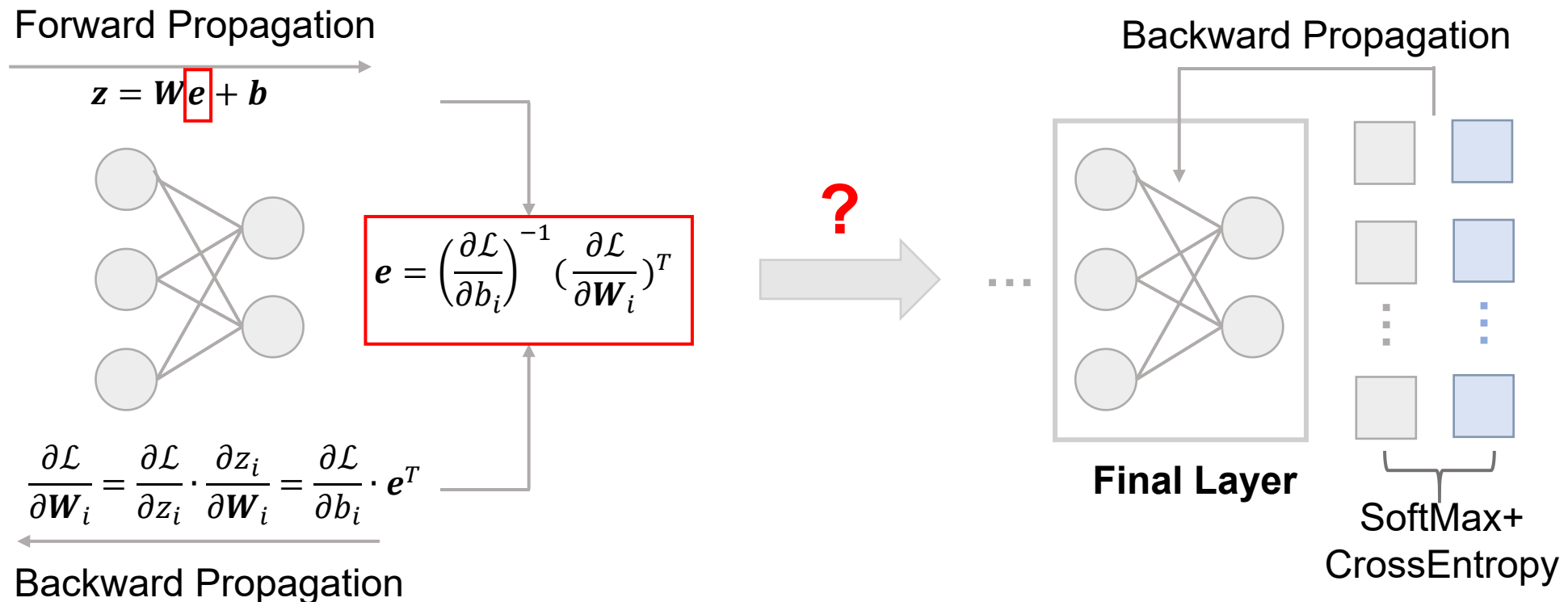


# Class-wise Embedding Reconstruction

- A provable **single embedding reconstruction** at **any fully-connected layer** has been proposed.

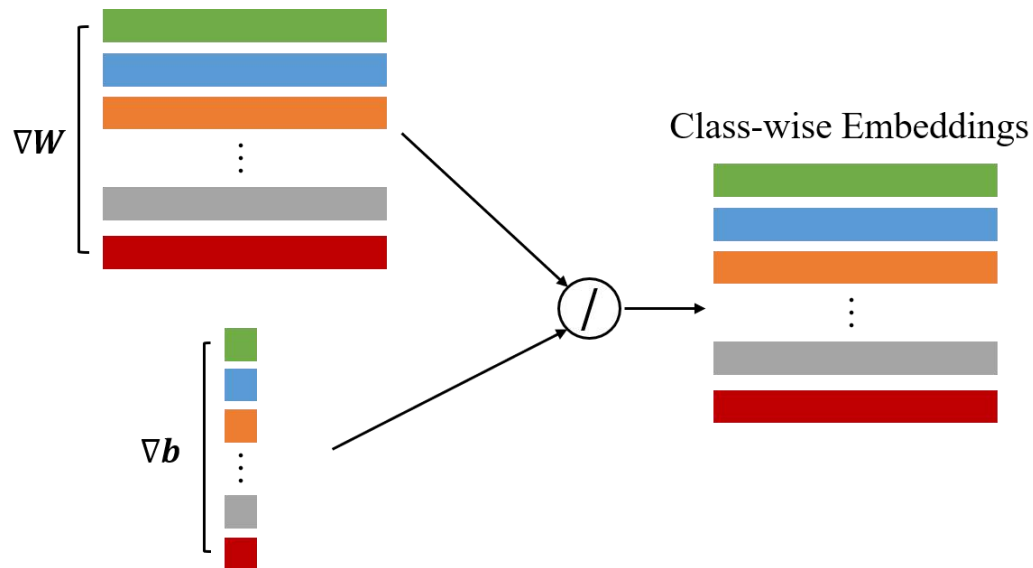
→ Extend to **recover batch embeddings** after average operations?

- **Distinct to Class-Averaged.**



# Class-wise Embedding Reconstruction

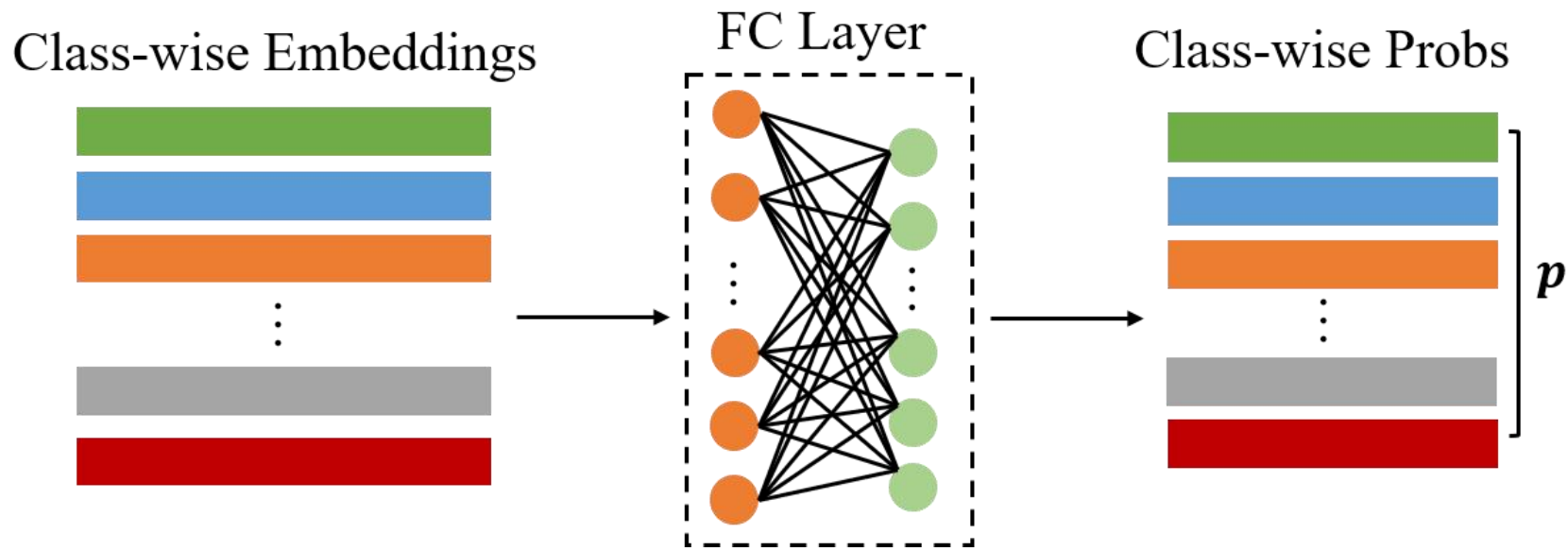
- Two Observations / Approximations (inspired by [\[Sun et al' 21\]](#)):
  - Approx 1 **Intra-class Uniformity and Concentration of Embedding Distribution**
    - We can replace gradients w.r.t.  $i$ -class with the arithmetic mean of this category.
  - Approx 2 **Inter-class Low Entanglement of Gradient Contributions**
    - $i$ -class samples mainly contribute to the  $i$ -th gradient row.



$$\bar{\mathbf{e}}_{\mathbb{B}i} \approx \frac{\partial \mathcal{L}}{\partial b_{i_{\mathbb{B}i}}}^{-1} \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{i_{\mathbb{B}i}}}^T \approx \frac{\partial \mathcal{L}}{\partial b_{i_{\mathbb{B}}}}^{-1} \cdot \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{i_{\mathbb{B}}}}^T$$

# Class-wise Probability Reconstruction

- Feed embeddings into the layer to obtain the subsequent post-softmax probabilities:
  - Approx 3 **Average Probabilities from Average Embeddings**
    - $\bar{p}_{\mathbb{B}_j} \approx \text{SoftMax}(\mathbf{W}\bar{e}_{\mathbb{B}_j} + \mathbf{b})$ .



# Instance-wise Label Restoration

- Construct a system of linear equations:

- Find out that **gradient w.r.t. model output z with SoftMax and CrossEntropy**

- $\nabla z_i = p_i - y_i$  (difference between post-softmax probabilities and one-hot labels)

- Equation composition:  $\nabla b_i$ ,  $\bar{p}_{\mathbb{B}_j}$ ,  $K$ (batch size),  $k_i$ (label frequency) → **Targets!**

$$k_1 + \dots + k_i + \dots + k_C = K$$

$$\sum_j k_j \bar{p}_{i\mathbb{B}_j} - K \nabla b_i = k_i$$

$$\begin{array}{l}
 \nabla z_i \xrightarrow{\sum} \sum_k \nabla z_i^k = \sum_k \nabla b_i^k = K \nabla b_i \\
 \parallel \qquad \parallel \qquad \bar{p}_{\mathbb{B}_1} \begin{bmatrix} \boxed{1} & \dots & \boxed{i} & \dots & \boxed{C} \end{bmatrix} \times \begin{bmatrix} k_1 \\ \vdots \\ k_i \\ \vdots \\ k_C \end{bmatrix} \\
 p_i \xrightarrow{\sum} \sum_k p_i^k = \sum_j k_j \bar{p}_{i\mathbb{B}_j} = \bar{p}_{\mathbb{B}_i} \begin{bmatrix} \boxed{1} & \dots & \boxed{i} & \dots & \boxed{C} \end{bmatrix} \times \begin{bmatrix} k_1 \\ \vdots \\ k_i \\ \vdots \\ k_C \end{bmatrix} \\
 | \qquad | \qquad \bar{p}_{\mathbb{B}_C} \begin{bmatrix} \boxed{1} & \dots & \boxed{i} & \dots & \boxed{C} \end{bmatrix} \times \begin{bmatrix} k_1 \\ \vdots \\ k_i \\ \vdots \\ k_C \end{bmatrix} \\
 y_i \xrightarrow{\sum} \sum_k y_i^k = k_i
 \end{array}$$

**Summation and Decomposition**



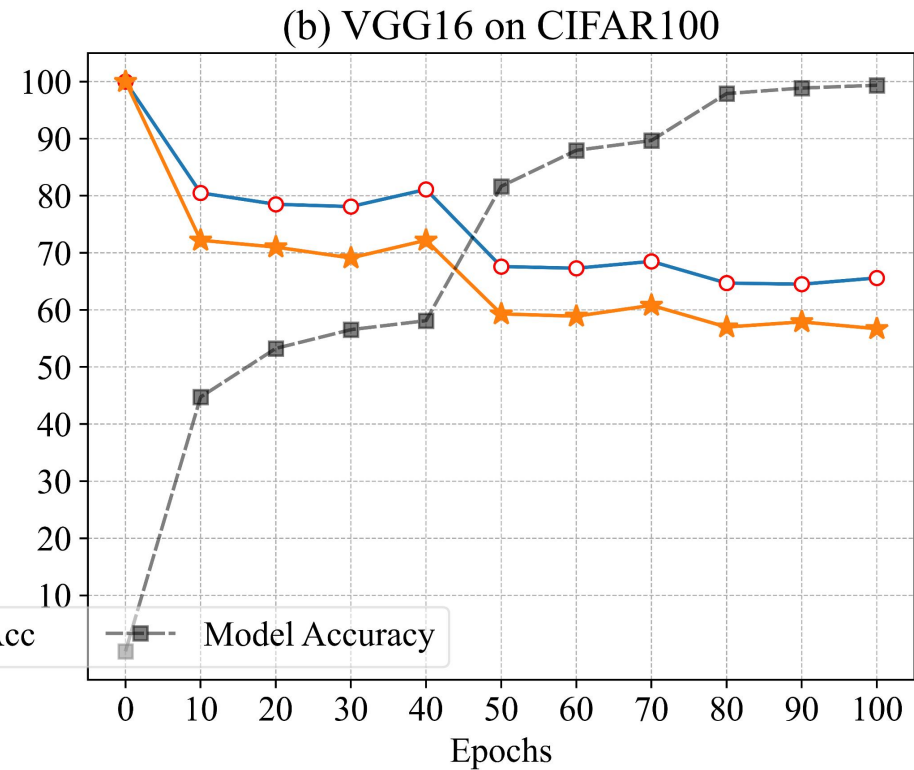
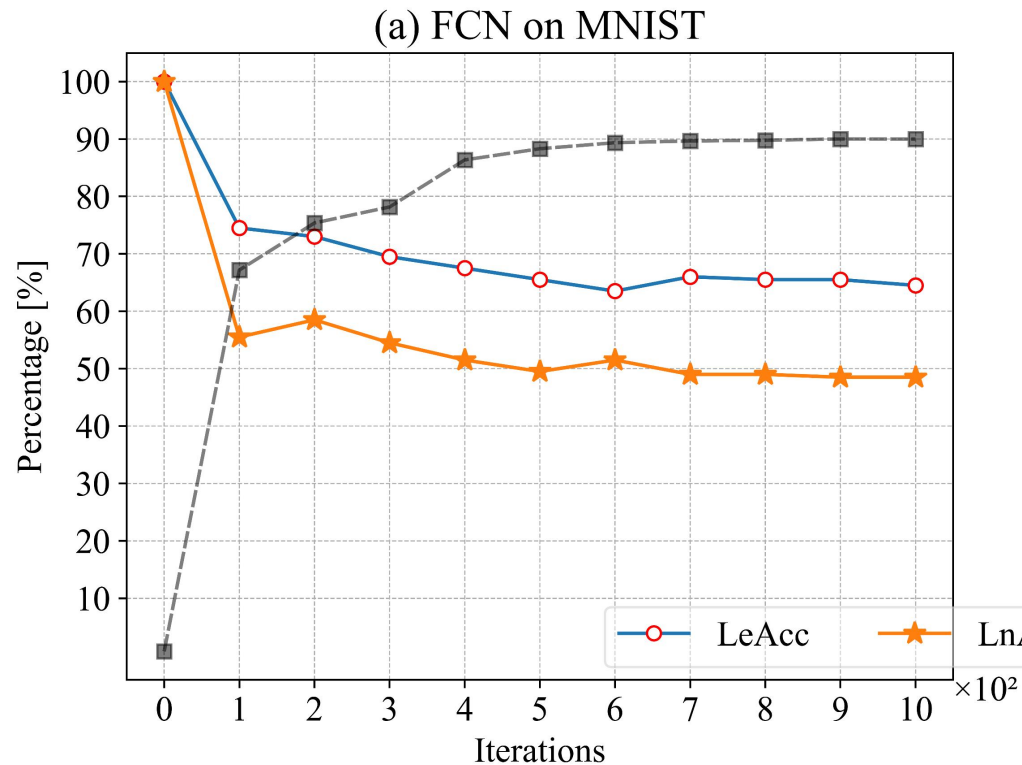
# Results on Image Classification

- We show iLRG **outperforms** prior works across several datasets, architectures and evaluation metrics.
  - Last column shows that recovery of class-wise probabilities are quite **precise** and **corroborate** with label restoration accuracies.

Model	Dataset	iDLG		GI		RLG		Ours		
		LeAcc	LnAcc	LeAcc	LnAcc	LeAcc	LnAcc	LeAcc	LnAcc	CosSim
FCN-3	MNIST	0.514	-	1.000	-	0.932	-	<b>1.000</b>	<b>0.994</b>	<b>0.979</b>
LeNet-5	CIFAR100	1.000	-	1.000	-	1.000	-	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
LeNet-S*	CIFAR100	0.150	-	0.213	-	1.000	-	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
VGG-16	ImageNet	1.000	-	1.000	-	0.981	-	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
ResNet-50	ImageNet	1.000	-	1.000	-	1.000	-	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>

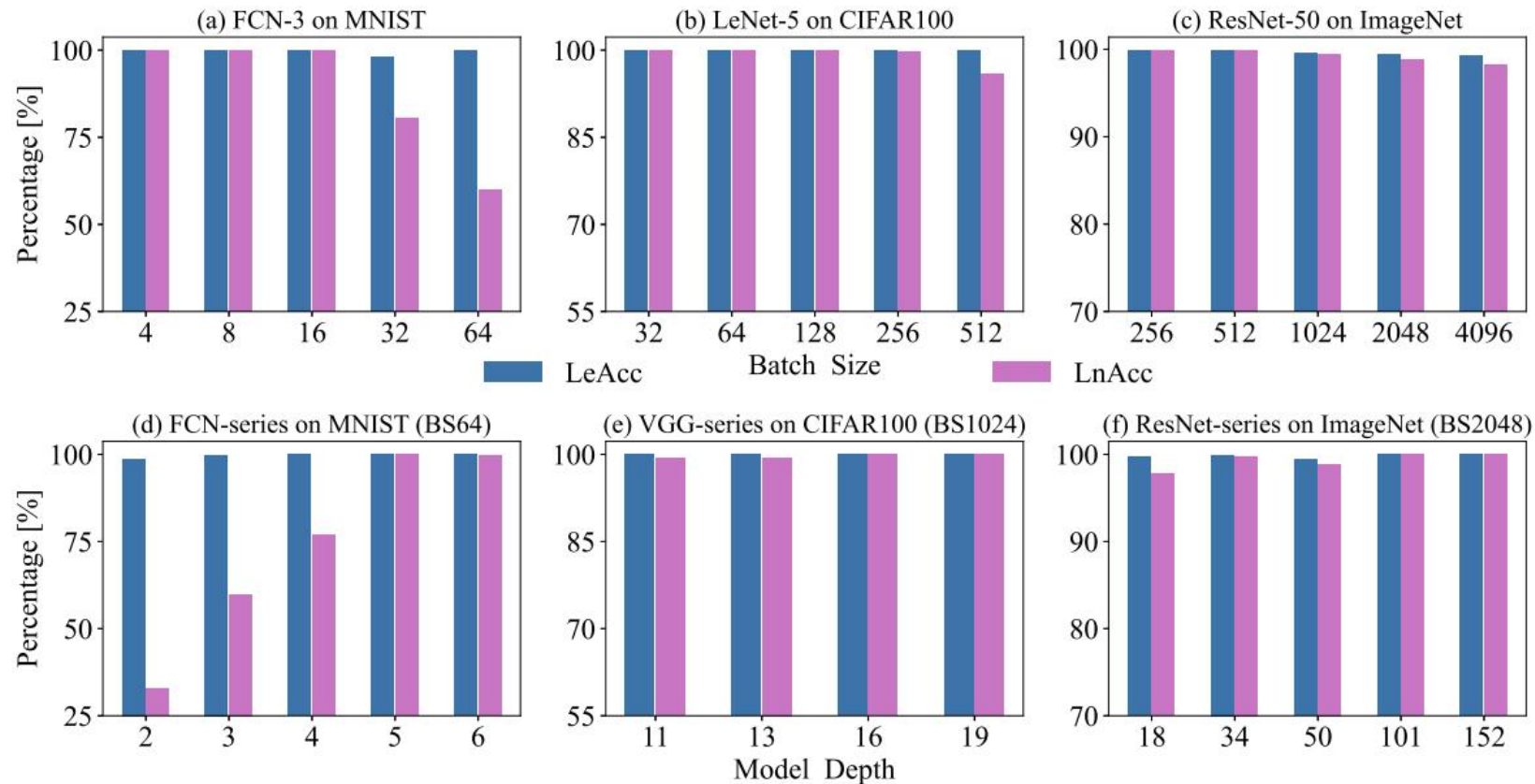
# Results on Different Training Stages

- Attack effects **get worse** as training progresses. (performs extremely well for untrained models)
  - **Recommend attacking in the early stages.**



# Results on Different Batch Sizes and Model Depths

- As we raise the batch size, the attack capability **decreases**.
- The attack performance **increases** as the network deepens.



# Results on Improved Image Recovery with iLRG

- iLRG can be used to specify an optimization objective for each instance, which **leads to an improved version of GIA.**



MNIST - Ground Truth



CIFAR100 - Ground Truth



MNIST - IG, PSNR $\uparrow$ : 12.46, LPIPS $\downarrow$ : 0.4269



CIFAR100 - IG, PSNR $\uparrow$ : 10.02, LPIPS $\downarrow$ : 0.5997



MNIST - Ours, PSNR $\uparrow$ : 18.73, LPIPS $\downarrow$ : 0.1357



CIFAR100 - Ours, PSNR $\uparrow$ : 17.51, LPIPS $\downarrow$ : 0.3439

# Conclusion

- We propose instance-wise Label Restoration from Gradients (iLRG), a method to **reveal instance-wise labels and class-wise embeddings** via shared batch-averaged gradients in FL.
- We conduct a large number of experiments to **demonstrate its effectiveness** and **explore several factors** that may influence it.
- We further **facilitate the existing gradient inversion attacks** by exploiting the recovered labels.



Code: <https://github.com/BUAA-CST/iLRG>