

# Building Gradient Bridges: Label Leakage from Restricted Gradient Sharing in Federated Learning

Rui Zhang

Department of Computing  
The Hong Kong Polytechnic University  
csrzhang1@comp.polyu.edu.hk

Ka-Ho Chow

Department of Computer Science  
The University of Hong Kong  
kachow@cs.hku.hk

Ping Li

Department of Computing  
The Hong Kong Polytechnic University  
p.li@polyu.edu.hk

**Abstract**—The growing concern over data privacy, the benefits of utilizing data from diverse sources for model training, and the proliferation of networked devices with enhanced computational capabilities have all contributed to the rise of federated learning (FL). The clients in FL collaborate to train a global model by uploading gradients computed on their private datasets without collecting raw data. However, a new attack surface has emerged from gradient sharing, where adversaries can restore the label distribution of a victim’s private data by analyzing the obtained gradients. To mitigate this privacy leakage, existing lightweight defenses restrict the sharing of gradients, such as encrypting the final-layer gradients or locally updating the parameters within. In this paper, we introduce a novel attack called Gradient Bridge (GDBR) that recovers the label distribution of training data from the limited gradient information shared in FL. GDBR explores the relationship between the layer-wise gradients, tracks the flow of gradients, and analytically derives the batch training labels. Extensive experiments show that GDBR can accurately recover more than 80% of labels in various FL settings. GDBR highlights the inadequacy of restricted gradient sharing-based defenses and calls for the design of effective defense schemes in FL.

**Index Terms**—Federated Learning, Gradient Inversion, Label Recovery, Privacy Leakage

## I. INTRODUCTION

Federated Learning (FL) is a privacy-preserving technology that enables multiple clients to collaboratively train a machine learning model without sharing their private data [1], [2]. In each communication round, the FL server sends the global model to each client, who uses their private data to compute gradients w.r.t. the model parameters and then transmits these gradients back to the server. The server aggregates the contributions from all clients and updates the global model for the next round of training. This collaborative process is repeated until the global model converges to a satisfactory state.

Despite its privacy-enhancing nature, recent studies [3]–[5] have shown that the shared gradients in FL provide a novel attack surface for privacy leakage. In particular, adversaries can reveal the label distribution of a client’s private training data by inspecting the gradients [6]–[9]. This label leakage can lead to more severe privacy threats, such as data reconstruction attacks [10], [11] that recover the training data of the victim client, membership inference attacks [12] that reveal whether a training instance belongs to the victim client’s dataset, and

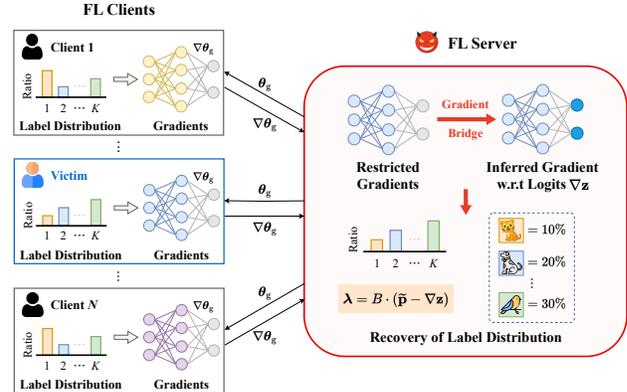


Fig. 1: Illustration of our Gradient Bridge (GDBR) attack. The gray parts in gradients denote the unshared information, while the colored parts depict the shared gradients. GDBR first infers the gradient w.r.t. output logits (in sky blue) from the obtained gradients and then recovers the label distribution of the victim client from the derived equation (in orange).

attribute inference attacks [13] that infer the values of sensitive attributes of the victim client’s training data.

To mitigate the above label leakage, the majority of existing defenses focus on encrypting the gradients with homomorphic encryption [14], [15] or protecting the gradients with differential privacy [16], [17]. However, these defenses either incur huge computational overhead or degrade the performance of the training model. To balance the trade-off between privacy and utility, recent studies have proposed dedicated lightweight defense strategies to protect critical gradients. These strategies include homomorphically encrypting the gradients in the last fully connected (FC) layer [18], removing the bias parameters from the FC layers [19], and locally updating the FC layers during FL training process [20]. From the perspective of raw label information contained in the gradients, these lightweight defenses aim to restrict gradient sharing, especially within the final FC layer, which is the most vulnerable layer for leakage of label distributions. Although these defense strategies appear to provide privacy without compromising utility, we argue that they offer a false sense of security.

To this end, we propose a novel label recovery attack, named Gradient Bridge (GDBR), to restore the label distribution from

\*Ka-Ho Chow and Ping Li are the corresponding authors.

the limited gradient information in FL. The key insight behind GDBR is to construct a bridge between the obtained gradients and the training labels. By exploring the relationship between the layer-wise gradients and the model parameters, we first trace the flow of gradients in the bottom layers of the model. Specifically, the gradient w.r.t. the input features of each layer can be recursively deduced from the starting layer at the bottom. Then, we derive the gradient w.r.t. the logits in a batch-averaged form, which contains the label information. Finally, we present the formulation of label restoration by associating the derived gradient with the observable variables. We conduct extensive experiments on various image datasets and classification models to validate the effectiveness and robustness of our GDBR attack. The experiment results show that more than 80% of the labels in the batch training data can be accurately recovered, regardless of the model architectures, dataset complexities, batch sizes, and class distributions. Even if the shared gradients are pruned or perturbed, GDBR can still recover the label distribution with acceptable degradation in performance. Fig. 1 illustrates our proposed GDBR attack. Our main contributions are summarized as follows:

- We propose a novel label recovery attack, named Gradient Bridge (GDBR), designed to recover the label distribution from the restricted gradient information in FL.
- We build a bridge between the shared gradients and the target labels by tracing the flow of gradients in the bottom layers, and derive the equations for label recovery.
- Extensive experiments show that GDBR is effective and robust in restoring the label distribution of private training data. Moreover, GDBR can still recover a certain amount of labels even when additional defenses are applied.

## II. RELATED WORK

In this section, we introduce the related works on gradient inversion attacks, analytical label recovery attacks, and representative lightweight defense strategies in FL.

### A. Gradient Inversion Attacks

Gradient inversion attacks [5] are designed to recover private training data by exploiting the transmitted gradients in FL. Zhu et al. [10] propose the DLG attack, which iteratively optimizes dummy data and labels to minimize the discrepancy between the generated gradients and the ground-truth gradients. Geiping et al. [21] enhance this attack by using cosine similarity as the loss function and introducing total variance regularization to improve the quality of the reconstructed images. Yin et al. [22] devise group consistency regularization to better restore the locations of main objects in images. Zhu et al. [23] present a recursive-based attack that recovers the input features of each layer by solving a system of linear equations. Furthermore, recent studies have also investigated inverting the gradients in other collaboration tasks, such as natural language processing [24], [25] and speech recognition [26].

### B. Analytical Label Recovery

In addition to data reconstruction via solving optimization problems, analytical label recovery attacks are also proposed to effectively reveal the private label distribution from gradient sharing. Zhao et al. [27] introduce iDLG, which discovers the correlation between the gradient of the last FC layer and the one-hot label. Yin et al. [22] propose GI to recover the labels by selecting the rows with the minimum values in the gradient of the last FC layer. Dang et al. [6] present RLG to distinguish the target class by separating the columns of the gradient w.r.t. weights in the output layer. Geng et al. [7] propose ZLG to estimate the posterior probabilities and extracted features of the last layer from auxiliary data to restore class-wise labels. Wainakh et al. [8] design LLG by utilizing both the direction and magnitude of the gradient w.r.t. weights in the FC layer to determine the overall impact factor and class-wise offset, and then restore the labels one by one. Ma et al. [28] devise iLRG to derive the embedding features and Softmax probabilities of the output layer and solve linear equations for label recovery.

However, all these attacks rely on the gradients of the final FC layer. Specifically, iDLG, GI and RLG can only recover non-repeating labels of the training data, while iLRG requires the gradient w.r.t. the bias term for feature restoration.

### C. Lightweight Defense Strategies

To preserve label privacy, several dedicated but lightweight defense strategies have been proposed to mitigate the leakage of private labels from uploaded gradients in FL. Sotthiwat et al. [18] propose Partially Encrypted Multi-Party Computation (PE-MPC) to encrypt the gradients in the FC layer closest to the data source. This ensures that the gradients in the output layer are neither accessible nor observable to the adversary. For some attacks that exploit the bias gradient to restore data or labels, such as [14], [28], [29], Scheliga et al. [19] present to remove the bias term from all the FC layers to entirely avoid such kind of leakage. Additionally, Mei et al. [20] propose the personalized FL training method FedVF, which allows each participant to locally update the parameters of the classification layers. FedVF seeks to balance the model performance and privacy protection by preventing the attacker from accessing the critical gradients in these layers.

## III. PRELIMINARY

In this section, we introduce the preliminaries of this work, including the inference of a single one-hot label, the gradients in typical layers, and our threat model in FL.

### A. Inference of Single One-Hot Label

In a classification task with  $C$  classes, the cross-entropy loss function  $\mathcal{L}$  is widely used to measure the distance between the predicted Softmax probabilities  $\mathbf{p} \in \mathbb{R}^C$  and the target one-hot labels  $\mathbf{y} \in \mathbb{R}^C$ . The cross-entropy loss is defined as:

$$\mathcal{L} = - \sum_{i=1}^C y_i \log(p_i) = - \log \frac{\exp(z_c)}{\sum_{j=1}^C \exp(z_j)},$$

where  $c$  is the index of the true class,  $z_i$  is the  $i$ -th element of the output logits  $\mathbf{z}$ , and  $p_i = \frac{\exp(z_i)}{\sum_{j=1}^C \exp(z_j)}$ .

In the backward propagation, the gradient of the loss function  $\mathcal{L}$  w.r.t. the output logits  $z_i$  can be formulated as:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial z_i} &= -\frac{\partial \log \exp(z_c)}{\partial z_i} + \frac{\partial \log \sum_{j=1}^C \exp(z_j)}{\partial z_i} \\ &= -\delta_{ic} + \frac{\exp(z_i)}{\sum_{j=1}^C \exp(z_j)}, \end{aligned}$$

where  $\delta_{ic}$  is the Kronecker delta function, which is equal to 1 if  $i = c$  and 0 otherwise. It is obvious to observe that  $y_i = \delta_{ic}$ , and  $\delta_{ic}$  just denotes the one-hot label. Therefore, the gradient w.r.t. the output logits  $\mathbf{z}$  can be rewritten as:

$$\nabla \mathbf{z} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}} = \mathbf{p} - \mathbf{y}. \quad (1)$$

The attacker can infer the one-hot label  $\mathbf{y}$  by exploiting the gradient  $\nabla \mathbf{z}$  and the estimated probabilities  $\mathbf{p}$ .

### B. Gradients in Typical Layers

Here, we introduce the gradients in typical layers, including the fully connected (FC) layer, the convolutional (Conv) layer, and the ReLU activation layer. For simplicity, we will ignore the bias items in the following discussion.

**FC Layer.** The forward propagation of an FC layer can be expressed as  $\mathbf{z} = \mathbf{W}\mathbf{x}$ , where  $\mathbf{z} \in \mathbb{R}^N$  is the output,  $\mathbf{W} \in \mathbb{R}^{N \times M}$  is the weight matrix, and  $\mathbf{x} \in \mathbb{R}^M$  is the input.

In the backward propagation, the gradient of the loss function  $\mathcal{L}$  w.r.t. the weight matrix  $\mathbf{W}$  is calculated as:

$$\nabla \mathbf{W} = \frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \cdot \frac{\partial \mathbf{z}}{\partial \mathbf{W}} = \nabla \mathbf{z} \mathbf{x}^\top, \quad (2)$$

where  $\nabla \mathbf{z}$  is the gradient w.r.t. the layer output  $\mathbf{z}$ .

Similarly, the gradient of the loss function  $\mathcal{L}$  w.r.t. the input feature  $\mathbf{x}$  can be derived as:

$$\nabla \mathbf{x} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}} \cdot \frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \mathbf{W}^\top \nabla \mathbf{z}. \quad (3)$$

**Conv Layer.** The weights in a Conv layer are also called kernels, which are shared across different locations within the input. The kernel is generally a 4D tensor, denoted as  $\mathbf{W} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times H_k \times W_k}$ , where  $C_{\text{in}}$  and  $C_{\text{out}}$  are the numbers of input and output channels, and  $H_k$  and  $W_k$  are the height and width of the kernel, respectively.

The forward pass can be represented as  $\mathbf{Z} = \mathbf{W} * \mathbf{X}$ , where  $\mathbf{Z} \in \mathbb{R}^{C_{\text{out}} \times H_{\text{out}} \times W_{\text{out}}}$  is the output,  $\mathbf{X} \in \mathbb{R}^{C_{\text{in}} \times H_{\text{in}} \times W_{\text{in}}}$  is the input, and  $*$  denotes the convolution operation. Unfolding the process of convolutional operation, a specific element  $Z_{k,i,j}$  in the output  $\mathbf{Z}$  is calculated as follows:

$$\begin{aligned} Z_{k,i,j} &= \sum_{c=1}^{C_{\text{in}}} \sum_{h=1}^{H_k} \sum_{w=1}^{W_k} W_{k,c,h,w} \cdot X_{c,i+h-1,j+w-1} \\ &= \sum_{c,h,w} W_{k,c,h,w} \cdot X_{c,i+h-1,j+w-1}, \end{aligned}$$

where  $k$  is the index of the output channel, and  $i$  and  $j$  are the indices of the height and width of the output, respectively.

In addition, we use the Einstein summation convention, such as  $\sum_{c,h,w}$ , to simplify the notation.

In the backward pass, the gradient w.r.t. element  $W_{k,c,h,w}$  in the kernel  $\mathbf{W}$  can be expressed as:

$$\begin{aligned} \nabla W_{k,c,h,w} &= \sum_{i=1}^{H_{\text{out}}} \sum_{j=1}^{W_{\text{out}}} \nabla Z_{k,i,j} \cdot X_{c,i+h-1,j+w-1} \\ &= \sum_{i,j} \nabla Z_{k,i,j} \cdot X_{c,i+h-1,j+w-1}. \end{aligned} \quad (4)$$

**ReLU Layer.** ReLU is the most commonly used activation function in deep neural networks. ReLU activation function is defined as  $a = \text{ReLU}(z) = \max(0, z)$ , where  $a$  is the output,  $z$  is the input, and  $\max(\cdot)$  is the maximum operation.

The derivative of the output  $a$  w.r.t. the input  $z$  is given by:

$$\sigma'(z) = \frac{da}{dz} = \begin{cases} 1, & \text{if } z > 0, \\ 0, & \text{otherwise.} \end{cases}$$

In the backward propagation, the gradient w.r.t. an element  $z_k$  in the input  $\mathbf{z}$  can be deduced as:

$$\nabla z_k = \frac{\partial \mathcal{L}}{\partial a_k} \cdot \sigma'(z_k) = \nabla a_k \cdot \begin{cases} 1, & \text{if } z_k > 0, \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

### C. Threat Model

In this paper, we assume an honest-but-curious attacker in FL, who follows the FL protocol but attempts to extract private information from the shared gradients. We mainly focus on the FedSGD [1], which is a typical FL algorithm in other inversion attacks [7], [8], [10], [22]. The attacker aims to recover the private labels of the clients' training data by analyzing the shared gradients. The attacker can launch GDBR in a white-box setting, where the attacker has full knowledge of the model architecture and parameters. Moreover, the attacker can also leverage some auxiliary information, such as testing datasets, to enhance the attack performance. In our attack scenario, the gradients of the final FC layer(s) in the model are not fully shared. The client will only transmit the first layer's gradients in the bottom layer stacking. Without loss of generality, we assume that the attacker will launch the GDBR attack at the early stage of the FL training process.

## IV. METHODOLOGY

In this section, we introduce the proposed Gradient Bridge (GDBR) attack against FL. We first present the overview of the GDBR. Then, we provide the theoretical foundation of GDBR, including the correlation between layer-wise gradients and the derived gradient w.r.t. output logits. Finally, we provide the formulation for label recovery by utilizing the estimated features and Softmax probabilities.

### A. Overview

Our GDBR aims to recover the labels of training data in a victim's local dataset by exploiting the gradient information shared in FL collaboration. Different from previous attacks, the adversary in GDBR cannot obtain the complete gradient information from a client's update. In particular, the gradient

from the last FC layer is always unavailable to the semi-honest server in FL. In GDBR, only one layer of gradient information is required for label recovery, which could be:

- gradients of middle layers in the MLP [30] model;
- gradients of the first or second FC layer in the LeNet [31], AlexNet [32] or VGG [33] model;
- gradients the last Conv layer in the ResNet [34] model, where the average pooling layer is assumed to be replaced by a Conv layer whose output size is  $1 \times 1$ .

Based on the gradient obtained from any of the above layers, the GDBR attack first builds the gradient bridge that connects the given layer to the output logits. Then, GDBR reconstructs the batch-averaged gradient w.r.t. the output logits. Finally, the attack infers the one-hot labels by associating the estimated probabilities with the gradient w.r.t. the logits.

### B. Correlation Between Layer-wise Gradients

In the following, we explore the correlations between the gradients in typical layers, such as FC, Conv and ReLU layers, or the layer stacks. These correlations are essential for GDBR to recover the labels from the gradients in FL.

**Lemma 1.** *In an FC layer, let  $\nabla \mathbf{x}$ ,  $\nabla \mathbf{W}$ , and  $\nabla \mathbf{z}$  represent the gradients w.r.t. the input  $\mathbf{x}$ , the weight  $\mathbf{W}$ , and the output  $\mathbf{z}$ , respectively. Then the following relationships hold:*

$$\nabla \mathbf{x} \mathbf{x}^\top = \mathbf{W}^\top \nabla \mathbf{W} \quad (6)$$

$$\nabla \mathbf{z} \mathbf{z}^\top = \nabla \mathbf{W} \mathbf{W}^\top \quad (7)$$

$$\nabla \mathbf{z} \odot \mathbf{z} = \text{diag}(\nabla \mathbf{W} \mathbf{W}^\top), \quad (8)$$

where  $\odot$  denotes the element-wise product, and  $\text{diag}(\cdot)$  denotes the diagonal elements of a matrix.

*Proof.* If we multiply  $\mathbf{x}^\top$  on both sides of Eq. (3), and then substitute Eq. (2) into it, we can derive:

$$\nabla \mathbf{x} \mathbf{x}^\top = \mathbf{W}^\top \nabla \mathbf{z} \mathbf{x}^\top = \mathbf{W}^\top \nabla \mathbf{W}.$$

Likewise, if we multiply  $\mathbf{W}^\top$  on both sides of Eq. (2), then we can obtain the following equation:

$$\nabla \mathbf{W} \mathbf{W}^\top = \nabla \mathbf{z} \mathbf{x}^\top \mathbf{W}^\top = \nabla \mathbf{z} \mathbf{z}^\top.$$

Since  $\nabla \mathbf{z} \mathbf{z}^\top$  is a symmetric matrix, the diagonal elements of it are equal to the element-wise product of the output  $\mathbf{z}$  and its gradient  $\nabla \mathbf{z}$ . So we also have:

$$\nabla \mathbf{z} \odot \mathbf{z} = \text{diag}(\nabla \mathbf{z} \mathbf{z}^\top) = \text{diag}(\nabla \mathbf{W} \mathbf{W}^\top),$$

where  $\odot$  denotes the element-wise product.  $\square$

**Lemma 2.** *In an FC layer, given the weight matrix  $\mathbf{W}$  and the input gradient  $\nabla \mathbf{x}$ , the output gradient  $\nabla \mathbf{z}$  can be derived as follows:*

$$\nabla \mathbf{z} = (\mathbf{W} \mathbf{W}^\top)^{-1} \mathbf{W} \nabla \mathbf{x}. \quad (9)$$

*Proof.* From the derived Eq. (6), we can individually denote the gradient  $\nabla \mathbf{W}$  by multiplying the inverse of  $\mathbf{W} \mathbf{W}^\top$  and  $\mathbf{W}$  on both sides of the equation, that is:

$$\nabla \mathbf{W} = (\mathbf{W} \mathbf{W}^\top)^{-1} \mathbf{W} \nabla \mathbf{x} \mathbf{x}^\top.$$

Combining with Eq. (2), we can transfer the role of  $\mathbf{x}^\top$  by multiplying  $\mathbf{x}$  and the inverse of  $\mathbf{x}^\top \mathbf{x}$  on both sides of the above equation, and deduce the gradient  $\nabla \mathbf{z}$  as:

$$\begin{aligned} \nabla \mathbf{z} &= (\mathbf{W} \mathbf{W}^\top)^{-1} \mathbf{W} \nabla \mathbf{x} \mathbf{x}^\top \mathbf{x} (\mathbf{x}^\top \mathbf{x})^{-1} \\ &= (\mathbf{W} \mathbf{W}^\top)^{-1} \mathbf{W} \nabla \mathbf{x}. \end{aligned}$$

$\square$

**Note:** Lemma 1 and 2 present the relationships between the gradients w.r.t. the input, output and weight in an FC layer. Gradient  $\nabla \mathbf{z}$  can be inferred by using: (1) gradient  $\nabla \mathbf{x}$ , or (2) gradient  $\nabla \mathbf{W}$  and output  $\mathbf{z}$ .

**Lemma 3.** *In a Conv layer, let  $\nabla \mathbf{W}_k$  and  $\nabla \mathbf{Z}_k$  represent the gradients w.r.t. the  $k$ -th kernel  $\mathbf{W}_k$  and the output  $\mathbf{Z}_k$ , respectively. Then the following relationship holds:*

$$\langle \nabla \mathbf{W}_k, \mathbf{W}_k \rangle_F = \begin{cases} \nabla \mathbf{Z}_k \odot \mathbf{Z}_k, & \text{if } \mathbf{Z}_k \in \mathbb{R}, \\ \langle \nabla \mathbf{Z}_k, \mathbf{Z}_k \rangle_F, & \text{otherwise,} \end{cases} \quad (10)$$

where  $\odot$  denotes the element-wise product, and  $\langle \cdot, \cdot \rangle_F$  denotes the Frobenius inner product that calculates the sum of the element-wise product of two tensors or matrices.

*Proof.* If we multiply  $\mathbf{W}_{k,c,h,w}$  on both sides of Eq. (4), and then sum over  $c$ ,  $h$ , and  $w$ , we can derive:

$$\begin{aligned} &\sum_{c,h,w} \nabla \mathbf{W}_{k,c,h,w} \cdot \mathbf{W}_{k,c,h,w} \\ &= \sum_{c,h,w} \sum_{i,j} \nabla \mathbf{Z}_{k,i,j} \cdot X_{c,h+i-1,w+j-1} \cdot \mathbf{W}_{k,c,h,w} \\ &= \sum_{i,j} \nabla \mathbf{Z}_{k,i,j} \cdot \left( \sum_{c,h,w} \mathbf{W}_{k,c,h,w} \cdot X_{c,i+h-1,j+w-1} \right) \\ &= \sum_{i,j} \nabla \mathbf{Z}_{k,i,j} \cdot \mathbf{Z}_{k,i,j}. \end{aligned}$$

So we can simplify the above equation as  $\langle \nabla \mathbf{W}_k, \mathbf{W}_k \rangle_F = \langle \nabla \mathbf{Z}_k, \mathbf{Z}_k \rangle_F$ . When the shape of the output  $\mathbf{Z}_k$  is  $1 \times 1$ , that is,  $\mathbf{Z}_k$  is a scalar, we can represent the result using the element-wise product  $\odot$ , so that  $\langle \nabla \mathbf{W}_k, \mathbf{W}_k \rangle_F = \nabla \mathbf{Z}_k \odot \mathbf{Z}_k$ .  $\square$

**Note:** Lemma 3 provides the relationship between the gradients w.r.t. the kernel and the output in a Conv layer. When the output size is  $1 \times 1$ , the Frobenius inner product can be replaced by the element-wise product.

**Lemma 4.** *In a ReLU layer, let  $\nabla \mathbf{z}$  and  $\nabla \mathbf{a}$  represent the gradients w.r.t. the input  $\mathbf{z}$  and output  $\mathbf{a}$ , respectively. Then the following relationship holds:*

$$\nabla \mathbf{z} \odot \mathbf{z} = \nabla \mathbf{a} \odot \mathbf{a}, \quad (11)$$

where  $\odot$  denotes the element-wise product, and  $\nabla \mathbf{Z} \odot \mathbf{Z} = \nabla \mathbf{A} \odot \mathbf{A}$  holds for both  $\mathbf{Z}$  and  $\mathbf{A}$  being matrices or tensors.

*Proof.* By multiplying  $z_k$  on both sides of Eq. (5), for each element  $k$ , we can deduce the following relationship:

$$\nabla z_k \cdot z_k = \nabla a_k \cdot \max(0, z_k) = \nabla a_k \cdot a_k.$$

The relationship can be easily extended to the entire vector  $\mathbf{z}$ , and we have  $\nabla \mathbf{z} \odot \mathbf{z} = \nabla \mathbf{a} \odot \mathbf{a}$  for each element in this ReLU layer. It is the same for the matrix or tensor case.  $\square$

**Theorem 1.** *In a single stack of an FC layer followed by a ReLU activation (FC-ReLU), the gradient  $\nabla \mathbf{a}$  can be inferred from the following relationships:*

$$\nabla \mathbf{a} = \begin{cases} (\mathbf{W}\mathbf{W}^\top)^{-1}\mathbf{W}\nabla \mathbf{x}, & \text{if } \nabla \mathbf{x} \text{ is given,} \\ \text{diag}(\nabla \mathbf{W}\mathbf{W}^\top) \odot \mathbf{a}, & \text{if } \nabla \mathbf{W} \text{ is given,} \end{cases} \quad (12)$$

where  $\odot$  denotes the element-wise division, and all elements in  $\mathbf{a}$  are assumed to be non-zero in the second case.

*Proof.* Combining Lemma 2 and Lemma 4, we can establish the connection between the activated output gradient  $\nabla \mathbf{a}$  and the input gradient  $\nabla \mathbf{x}$  in the FC-ReLU stack as follows:

$$\nabla \mathbf{a} = \nabla \mathbf{z} \odot \mathbf{z} \odot \mathbf{a} = (\mathbf{W}\mathbf{W}^\top)^{-1}\mathbf{W}\nabla \mathbf{x} \odot \mathbf{z} \odot \mathbf{a}.$$

Since we have assumed that all elements in  $\mathbf{a}$  are non-zero, it is obvious that  $\mathbf{a} = \mathbf{z}$  holds from the definition of ReLU activation. In this case,  $\mathbf{z} \odot \mathbf{a}$  can be simplified as  $\mathbf{1}$ , where all the elements in  $\mathbf{1}$  are equal to 1. Therefore, we can simplify the above equation, that is:

$$\nabla \mathbf{a} = (\mathbf{W}\mathbf{W}^\top)^{-1}\mathbf{W}\nabla \mathbf{x}.$$

Combining Eq. (8) and Eq. (11), there is another way to infer the gradient  $\nabla \mathbf{a}$  when the gradient  $\nabla \mathbf{W}$  and the output  $\mathbf{a}$  are given. The derivation of  $\nabla \mathbf{a}$  can be expressed as:

$$\nabla \mathbf{a} = \text{diag}(\nabla \mathbf{W}\mathbf{W}^\top) \odot \mathbf{a},$$

where all elements in  $\mathbf{a}$  are assumed to be non-zero.  $\square$

**Note:** Theorem 1 illustrates two ways to derive the input gradient  $\nabla \mathbf{a}$  in a single FC-ReLU stack. The gradient  $\nabla \mathbf{a}$  can be either deduced from (1) the gradient  $\nabla \mathbf{x}$ , or (2) the gradient  $\nabla \mathbf{W}$  and the activation  $\mathbf{a}$ .

**Theorem 2.** *In a single stack of a Conv layer followed by a ReLU activation (Conv-ReLU), the output gradient  $\nabla \mathbf{A} \in \mathbb{R}^N$  can be inferred from the following equations:*

$$\begin{cases} \nabla \mathbf{A}_k = \langle \nabla \mathbf{W}_k, \mathbf{W}_k \rangle_F \odot \mathbf{A}_k, \\ \nabla \mathbf{A} = [\nabla \mathbf{A}_1, \nabla \mathbf{A}_2, \dots, \nabla \mathbf{A}_K]^\top, \end{cases} \quad (13)$$

where  $\odot$  denotes the element-wise division, and each activated output  $\mathbf{A}_k$  is assumed to be non-zero.

*Proof.* According to Lemma 3 and Lemma 4, we have  $\nabla \mathbf{Z} \odot \mathbf{Z} = \nabla \mathbf{A} \odot \mathbf{A}$  and  $\nabla \mathbf{Z} \odot \mathbf{Z} = \langle \nabla \mathbf{W}, \mathbf{W} \rangle_F$  in this Conv-ReLU stack, where  $\mathbf{Z}$  and  $\mathbf{A}$  are the outputs of the Conv layer and the ReLU layer, respectively. We mainly consider that the shape of  $\mathbf{A}_k$  is  $1 \times 1$ , that is,  $\mathbf{A}_k \in \mathbb{R}$  and  $\nabla \mathbf{A} \in \mathbb{R}^N$ , where  $N$  is the number of kernels in the Conv layer. Thus, for the  $k$ -th kernel, we can derive the gradient  $\nabla \mathbf{A}_k$  as:

$$\nabla \mathbf{A}_k = \langle \nabla \mathbf{W}_k, \mathbf{W}_k \rangle_F \odot \mathbf{A}_k,$$

where  $\mathbf{A}_k$  is assumed to be non-zero. Finally, we can deduce the gradient  $\nabla \mathbf{A}$  by concatenating all the  $\nabla \mathbf{A}_k$  together.  $\square$

**Note:** When  $\mathbf{A} \in \mathbb{R}^N$ , Theorem 2 presents the gradient  $\nabla \mathbf{A}$  in a single Conv-ReLU stack, which can be inferred by using the gradient  $\nabla \mathbf{W}$  and the activation  $\mathbf{A}$ .

In deep neural networks, the bottom layers of a classification model are usually composed of Conv-ReLU-FC-ReLU stacks or FC-ReLU stacks. We consider a more general structure that  $L$  layers of such stacks are cascaded in the bottom layers. The first stack is an FC-ReLU or Conv-ReLU, and the last stack is an FC layer without the ReLU activation. The other stacks are all FC-ReLU. We use superscripts  $[l]$  to denote the layer index, where  $l \in \{1, 2, \dots, L\}$ . This structure is common in the LeNet, AlexNet, VGG, and ResNet models. In summary, the layers that we analyze include:

- 1 FC-ReLU or Conv-ReLU stack ( $l = 1$ );
- $(L - 2)$  FC-ReLU stacks ( $l \in \{2, 3, \dots, L - 1\}$ );
- 1 FC output layer without ReLU activation ( $l = L$ ).

Combining the previous Theorem 1 and Theorem 2, we can infer the gradient  $\nabla \mathbf{A}_k^{[1]}$  in the first stack as follows:

$$\nabla \mathbf{a}_k^{[1]} = \begin{cases} \text{diag}(\nabla \mathbf{W}_k^{[1]} \mathbf{W}_k^{[1]\top}) \odot \mathbf{a}_k^{[1]}, & \text{if FC,} \\ \langle \nabla \mathbf{W}_k^{[1]}, \mathbf{W}_k^{[1]} \rangle_F \odot \mathbf{a}_k^{[1]}, & \text{if Conv,} \end{cases} \quad (14)$$

where  $\nabla \mathbf{W}_k^{[1]}$  can be directly obtained from the shared gradients in FL. In the design of GDBR, all we need is the gradient  $\nabla \mathbf{W}_k^{[1]}$  in the first stack of the bottom layers.

For the  $l$ -th layer,  $l \in \{2, 3, \dots, L - 1\}$ , we can recursively deduce the gradients  $\nabla \mathbf{a}^{[l]}$  and  $\nabla \mathbf{x}^{[l+1]}$  from Theorem 1 as:

$$\begin{cases} \nabla \mathbf{x}^{[2]} = \nabla \mathbf{a}^{[1]} = [\nabla \mathbf{a}_1^{[1]}, \nabla \mathbf{a}_2^{[1]}, \dots, \nabla \mathbf{a}_K^{[1]}]^\top \\ \nabla \mathbf{x}^{[3]} = \nabla \mathbf{a}^{[2]} = (\mathbf{W}^{[2]} \mathbf{W}^{[2]\top})^{-1} \mathbf{W}^{[2]} \nabla \mathbf{x}^{[2]} \\ \vdots \\ \nabla \mathbf{x}^{[l+1]} = \nabla \mathbf{a}^{[l]} = (\mathbf{W}^{[l]} \mathbf{W}^{[l]\top})^{-1} \mathbf{W}^{[l]} \nabla \mathbf{x}^{[l]}. \end{cases} \quad (15)$$

Finally, according to Lemma 2, we can derive the gradient  $\nabla \mathbf{z}^{[L]}$  in the last FC layer as the following equation:

$$\nabla \mathbf{z}^{[L]} = (\mathbf{W}^{[L]} \mathbf{W}^{[L]\top})^{-1} \mathbf{W}^{[L]} \nabla \mathbf{x}^{[L]}. \quad (16)$$

We name the above derivations as the *gradient bridge* that connects the gradient w.r.t. the first FC or Conv layer to the gradient w.r.t. the output logits. In this scenario, only the gradient  $\nabla \mathbf{W}^{[1]}$  of the first stack in the bottom layers and the global model  $\theta_g$  are required for inferring the gradient  $\nabla \mathbf{z}^{[L]}$  of the output logits. In FL collaboration, the model weights  $\mathbf{W}^{[l]}$  in each layer are known to the server, and the activated features  $\mathbf{a}^{[1]}$  can also be estimated by fitting some auxiliary data into the model. Our finding of the gradient bridge breaks the limitation of the restricted gradient sharing in FL, which is the foundation of the GDBR attack.

**Note:** Given the gradient  $\nabla \mathbf{W}^{[1]}$  of the first FC or Conv in the bottom layers, the *gradient bridge* can be built to infer the gradient  $\nabla \mathbf{z}^{[L]}$  w.r.t. the output logits.

### C. Derivation of Batch-Averaged Gradients

**Assumption 1.** We assume that the activation features  $\mathbf{a}^{[1]}$  in the first stack of the bottom layers exhibit similar Gaussian distributions across different dimensions for various training samples. Specifically, the features  $\mathbf{a}^{[1(n)}$  for the  $n$ -th sample in a batch of  $B$  can be approximated by some auxiliary data, denoted as  $\tilde{\mathbf{a}}^{[1]}$ . Consequently, we have:

$$\tilde{\mathbf{a}}^{[1]} \approx \mathbf{a}^{[1(n)}, \quad \forall n \in \{1, 2, \dots, B\}.$$

Combining Eq. (14) and Assumption 1, we can deduce the batch-averaged gradient  $\overline{\nabla \mathbf{a}}^{[1]}$  by utilizing  $\tilde{\mathbf{a}}^{[1]}$  to approximate the features  $\mathbf{a}^{[1(n)}$ . So the averaged gradient  $\overline{\nabla \mathbf{a}}^{[1]}$  is given by the following equations:

$$\begin{aligned} \overline{\nabla \mathbf{a}}_k^{[1]} &= \frac{1}{B} \sum_{n=1}^B \nabla \mathbf{a}_k^{[1(n)} \\ &\approx \begin{cases} \left( \frac{1}{B} \sum_n \text{diag}(\nabla \mathbf{W}_k^{[1(n)} \mathbf{W}_k^{[1\top]}) \right) \odot \tilde{\mathbf{a}}_k^{[1]} \\ \left( \frac{1}{B} \sum_n \langle \nabla \mathbf{W}_k^{[1(n)}, \mathbf{W}_k^{[1]} \rangle_{\text{F}} \right) \odot \tilde{\mathbf{a}}_k^{[1]} \end{cases} \\ &= \begin{cases} \text{diag}(\overline{\nabla \mathbf{W}}_k^{[1]} \mathbf{W}_k^{[1\top]}) \odot \tilde{\mathbf{a}}_k^{[1]}, & \text{if FC,} \\ \langle \overline{\nabla \mathbf{W}}_k^{[1]}, \mathbf{W}_k^{[1]} \rangle_{\text{F}} \odot \tilde{\mathbf{a}}_k^{[1]}, & \text{if Conv.} \end{cases} \end{aligned}$$

For the  $l$ -th layer ( $l \in \{2, 3, \dots, L-1\}$ ), the batch-averaged gradient  $\overline{\nabla \mathbf{a}}^{[l]}$  can be obtained as follows:

$$\begin{aligned} \overline{\nabla \mathbf{a}}^{[l]} &= \frac{1}{B} \sum_{n=1}^B \nabla \mathbf{a}^{[l(n)} \\ &= (\mathbf{W}^{[l]} \mathbf{W}^{[l\top]})^{-1} \mathbf{W}^{[l]} \left( \frac{1}{B} \sum_{n=1}^B \nabla \mathbf{x}^{[l(n)} \right) \\ &= (\mathbf{W}^{[l]} \mathbf{W}^{[l\top]})^{-1} \mathbf{W}^{[l]} \overline{\nabla \mathbf{x}}^{[l]}. \end{aligned}$$

Similarly, the batch-averaged gradient w.r.t. the logits  $\mathbf{z}^{[L]}$  in the output layer can be derived in the same way as above.

In summary, for the batch-averaged setting, the gradients for each stack of the bottom layers can be recursively inferred using the following equations:

$$\overline{\nabla \mathbf{a}}_k^{[1]} \approx \begin{cases} \text{diag}(\overline{\nabla \mathbf{W}}_k^{[1]} \mathbf{W}_k^{[1\top]}) \odot \tilde{\mathbf{a}}_k^{[1]}, & \text{if FC,} \\ \langle \overline{\nabla \mathbf{W}}_k^{[1]}, \mathbf{W}_k^{[1]} \rangle_{\text{F}} \odot \tilde{\mathbf{a}}_k^{[1]}, & \text{if Conv.} \end{cases} \quad (17)$$

$$\begin{cases} \overline{\nabla \mathbf{x}}^{[2]} = \overline{\nabla \mathbf{a}}^{[1]} = [\overline{\nabla \mathbf{a}}_1^{[1]}, \overline{\nabla \mathbf{a}}_2^{[1]}, \dots, \overline{\nabla \mathbf{a}}_{K'}^{[1]}]^\top \\ \overline{\nabla \mathbf{x}}^{[3]} = \overline{\nabla \mathbf{a}}^{[2]} = (\mathbf{W}^{[2]} \mathbf{W}^{[2\top]})^{-1} \mathbf{W}^{[2]} \overline{\nabla \mathbf{x}}^{[2]} \\ \vdots \\ \overline{\nabla \mathbf{x}}^{[l+1]} = \overline{\nabla \mathbf{a}}^{[l]} = (\mathbf{W}^{[l]} \mathbf{W}^{[l\top]})^{-1} \mathbf{W}^{[l]} \overline{\nabla \mathbf{x}}^{[l]}. \end{cases} \quad (18)$$

$$\overline{\nabla \mathbf{z}}^{[L]} = (\mathbf{W}^{[L]} \mathbf{W}^{[L\top]})^{-1} \mathbf{W}^{[L]} \overline{\nabla \mathbf{x}}^{[L]}. \quad (19)$$

### D. Label Recovery from Inferred Gradients

**Assumption 2.** We assume that the Softmax probabilities  $\mathbf{p}$  of the model output exhibit similar Gaussian distributions across different classes for various training samples. Specifically, the probabilities  $\mathbf{p}^{(n)}$  for the  $n$ -th sample in a batch of size  $B$  can be approximated by some auxiliary data, denoted as  $\tilde{\mathbf{p}}$ . Consequently, we have:

$$\tilde{\mathbf{p}} \approx \mathbf{p}^{(n)}, \quad \forall n \in \{1, 2, \dots, B\}.$$

Combining Eq. (1) in Preliminary, the batch-averaged gradient  $\overline{\nabla \mathbf{z}}^{[L]}$  can be represented as:

$$\overline{\nabla \mathbf{z}}^{[L]} = \frac{1}{B} \sum_{n=1}^B \nabla \mathbf{z}^{[L(n)} = \frac{1}{B} \sum_{n=1}^B (\mathbf{p}^{(n)} - \mathbf{y}^{(n)}). \quad (20)$$

In the batch-averaged setting, we use  $\lambda$  to denote the sum of the one-hot labels for the batch of  $B$  samples. Each element in  $\lambda$  is an integer, which represents the number of samples that have the corresponding class label in the batch. By substituting  $\lambda$  into Eq. (20) and use  $\tilde{\mathbf{p}}$  in Assumption 2 to approximate the probabilities  $\mathbf{p}^{(n)}$ , we can derive:

$$\lambda = \sum_{i=1}^B \mathbf{y}^{(i)} = B \cdot (\tilde{\mathbf{p}} - \overline{\nabla \mathbf{z}}^{[L]}). \quad (21)$$

Finally, we can infer the number of samples for each class in the batch by solving  $\lambda$  in Eq. (21). Gradient  $\overline{\nabla \mathbf{z}}^{[L]}$  is the output of the gradient bridge, which has been provided in Eq. (19).

## V. EXPERIMENTS

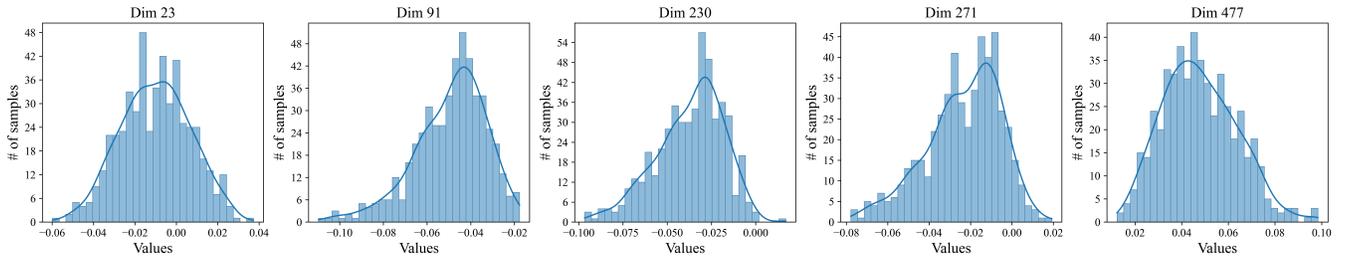
In this section, we conduct extensive experiments on various datasets and models to evaluate the effectiveness of GDBR in recovering labels. We first validate the reasonability of the assumptions in Section IV. Then, we compare the performance with prior baseline methods and analyze the influence of different factors on GDBR, such as batch sizes, different layers, class distributions, etc. Finally, we present the performance of GDBR against two defense mechanisms, which are gradient pruning and noise perturbation.

### A. Setups

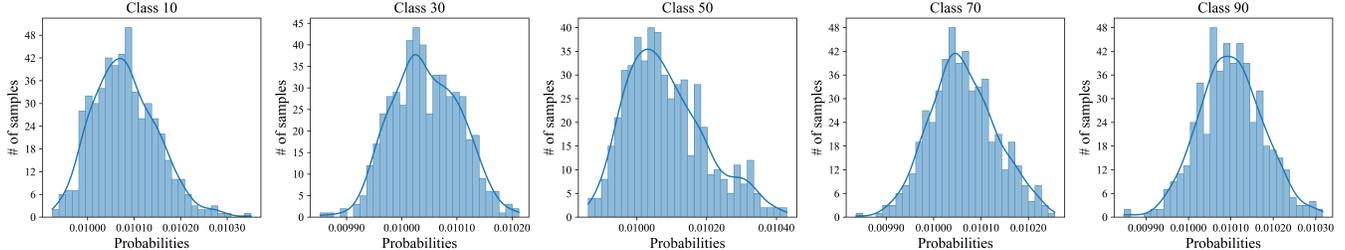
**Datasets and Models.** We evaluate the GDBR attack on five benchmark datasets and five popular neural network models. Here are the details of the datasets and models.

- MNIST [31] contains 60k training and 10k testing gray-scale digit images of size  $28 \times 28$  with 10 classes.
- SVHN [35] contains 73k training and 26k testing color digit images of size  $32 \times 32$  with 10 classes.
- CIFAR-10/100 [36] contain 50k training and 10k testing color images of size  $32 \times 32$  with 10/100 classes.
- ImageNet-1K [37] contains 1.28M training and 50k validation color images of size  $224 \times 224$  with 1,000 classes.

We employ five neural network models with different architectures, including MLP (Multi-Layer Perceptron) [30], LeNet [31], AlexNet [32], VGG series [33], and ResNet series [34]. In particular, the MLP model consists of 6 FC layers with



(a) Input feature distributions of the FC layer in ResNet18 ( $B = 512$ ). The shown dimensions are randomly selected, and the values in each dimension approximately follow a Gaussian distribution whose mean can be estimated from similar training samples.



(b) Softmax probability distributions of ResNet18 model ( $B = 512$ ). Five 10-fold classes are chosen for illustration, and the probabilities of different classes exhibit similar Gaussian distributions, with mean values close to 0.01, across various training samples.

Fig. 2: The distributions of input features and output probabilities of the FC layer in ResNet18, which is trained on the CIFAR-100 dataset. The x-axis represents the values of features or probabilities, and the y-axis represents the number of samples.

[2048, 1024, 512, 256, 128, 64] hidden units. All models are implemented using the PyTorch framework [38].

**Implementation Details.** We consider an FL system with  $N$  participants, where one of them is designated as the victim. Each client trains the model using its private dataset with the FedSGD algorithm [1]. During each training iteration, a client randomly selects  $B$  samples from its dataset to create a mini-batch, with the default batch size set to 64.

For the ResNet series models, we replace the average pooling layer with a convolutional layer. We use PyTorch’s default initialization for the layers preceding the average pooling. For the bottom layers, we initialize the weights using a uniform distribution between 0.01 and 0.2, and we ignore the bias terms for simplicity. By default, we use the uniform distribution to initialize the weights of the bottom layers, and we provide the gradient of the penultimate layer to GDBR.

To estimate the input features of the penultimate layer and the output probabilities of the model, we utilize the validation or testing dataset as an auxiliary dataset. We sample a total of 1,000 samples from the auxiliary data, distributing them evenly across all the classes. To eliminate experimental randomness, we repeat each experiment 20 times and report the average result. All the experiments are conducted using PyTorch 2.0.0 and CUDA 12.2 on a workstation equipped with an Intel i9-10900K CPU @ 3.70GHz, 64GB of RAM, and an NVIDIA GeForce RTX 4090 (24GB) GPU.

**Baselines and Evaluation Metrics.** We compare our proposed GDBR with two state-of-the-art methods, ZLG [7] and LLG [8]. ZLG leverages the relationship between the labels and the gradient  $\nabla \mathbf{W}^{[L]}$  in the final FC layer to recover labels,

while LLG exploits the direction and magnitude of  $\nabla \mathbf{W}^{[L]}$  to restore the labels of the batch samples. However, gradient  $\nabla \mathbf{W}^{[L]}$  is not directly accessible in our attack scenario. For a fair comparison, we provide part of the ground-truth gradient  $\nabla \mathbf{W}^{[L]}$  to both ZLG and LLG for label recovery. Moreover, we provide the same subset of auxiliary data to GDBR and the baselines to estimate the features of the penultimate layer and the output probabilities of the model.

To quantify the performance of label restoration, we adopt two metrics: *Instance-level Accuracy (InsAcc)* and *Class-level Accuracy (ClsAcc)* [9]. InsAcc evaluates the accuracy of recovering the ground-truth labels of individual data samples, while ClsAcc assesses the accuracy of recovering the ground-truth classes to which the data samples belong. All the results are reported as percentages (%).

### B. Verification of Assumptions

We first verify the rationalization of Assumption 1 and 2 by conducting experiments on the ResNet18 model trained on the CIFAR-100 dataset. We select five random dimensions within input features and five 10-fold classes to visualize the distributions of the features and the probabilities in the final FC layer. The findings, shown in Fig. 2, indicate that at the early stage of training, the extracted features and the output probabilities of different samples exhibit similar distributions. This observation supports the use of auxiliary data to estimate the hidden features and the output probabilities of the model, which are essential for GDBR to build the gradient bridge.

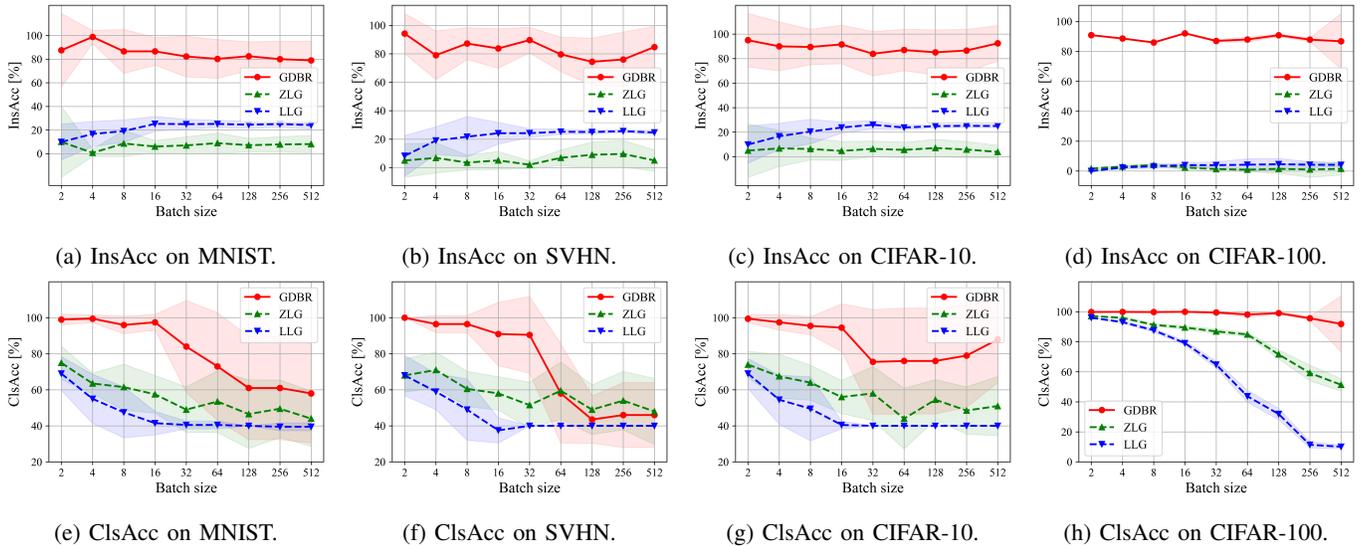


Fig. 3: Comparison of GDBR with baselines on InsAcc and ClsAcc across different datasets and batch sizes. The experiments are performed on four model-dataset pairs: LeNet on MNIST, AlexNet on SVHN, VGG11 on CIFAR-10, and ResNet18 on CIFAR-100, with batch sizes ranging from 2 to 512. The batch data is randomly sampled from a subset of training classes.

### C. Comparison with Baselines

Next, we compare GDBR with the baselines across various combinations of datasets and models. The model-dataset pairs are as follows: LeNet on MNIST, AlexNet on SVHN, VGG11 on CIFAR-10, and ResNet18 on CIFAR-100. We set the batch sizes ranging from 2 to 512, with the batch data randomly sampled from a subset of all classes in the training datasets. To provide the weight gradient in the last layer to ZLG and LLG, we shuffle the column elements of the matrix. The results, displayed in Fig. 3, show that GDBR outperforms ZLG and LLG on both InsAcc and ClsAcc metrics in almost all cases. Even with access to part of the ground-truth gradient, ZLG and LLG still struggle to recover the labels of the target batch. In contrast, GDBR can effectively recover the labels with high accuracy, demonstrating the effectiveness of our method.

### D. Analysis of Different Factors

In this subsection, we investigate the influence of various factors on the performance of GDBR. These factors include the gradients from different layers, the class distribution of the batch data, the use of auxiliary data or dummy data, and the gradient simulation and model initialization modes.

1) *Effect of Gradients from Different Layers*: We investigate the influence of gradients from different layers provided to GDBR. The experiments are conducted on a 6-layer MLP model using the SVHN and CIFAR-10 datasets, with the batch data randomly sampled from the training datasets. As shown in Fig. 4, the gradient of the first layer (6th from bottom) has the worst results. This is because the features of the first layer are more scattered and directly affected by the training data, making it harder to estimate and build an accurate bridge to the output logits. Additionally, due to the more complex nature of CIFAR-10 compared to SVHN, GDBR’s performance on

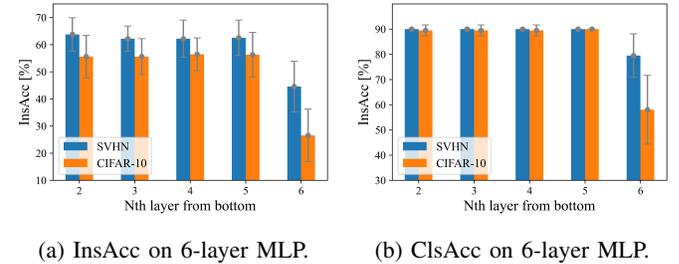


Fig. 4: Comparison of utilized gradients from different layers in a 6-layer MLP model, trained on MNIST and CIFAR-10.

CIFAR-10 is inferior to that on SVHN. The gradients from other layers exhibit similar performance, indicating the robustness of GDBR when applied to different layers.

2) *Effect of Batch Class Distribution*: We next examine the impact of different batch class distributions on the performance of GDBR. We evaluate GDBR with five different batch class distributions: *random*, *uniform*, *single*, *subclassed*, and *imbalanced*, and report the InsAcc. The experiments are conducted on different datasets and models, and the results are shown in Table I. The findings indicate that GDBR performs well in all settings. Although the performance is slightly affected by the data distribution, GDBR consistently achieves high accuracy in recovering the labels of the target batch data. Notably, batch data containing only one class yields better results than mixed classes, likely because the variables are more concentrated and easier to estimate. These analyses underscore the robustness of GDBR against different data distributions.

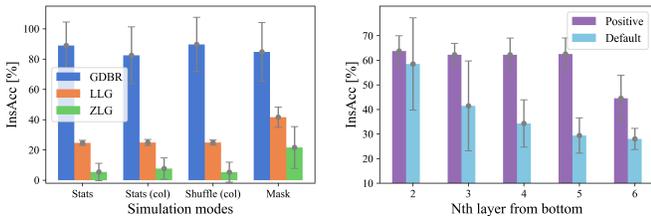
3) *Effect of Auxiliary Data vs. Dummy Data*: We further compare the effect of using auxiliary data versus dummy data in the GDBR. As mentioned in the Implementation Details, we

TABLE I: Comparison of different class distributions within the batch training data across various datasets and models.

Dataset	Model	InsAcc [%]				
		R	U	S	SC	IM
MNIST	LeNet	80.7	83.0	86.3	85.0	93.5
SVHN	AlexNet	85.6	81.6	81.7	85.2	90.1
CIFAR-10	VGG11	84.1	81.2	89.6	90.0	97.0
CIFAR-100	ResNet18	89.8	88.5	94.3	87.9	90.4
ImageNet	ResNet50	95.2	97.0	98.8	90.0	91.0

\*R: random, U: uniform, S: single, SC: subclassed, IM: imbalanced.

use the validation or testing dataset as auxiliary data, selecting 1,000 samples evenly distributed across all classes. For dummy data, we generate the data from a standard Gaussian distribution. The evaluation results are shown in Table II. For simple datasets like MNIST, GDBR can achieve even higher InsAcc with dummy data. However, for the other color image datasets, GDBR performs better with auxiliary data. The auxiliary data provides more accurate information to estimate the extracted features and the output probabilities, helping GDBR recover the labels more accurately. Nonetheless, when auxiliary data is unavailable, GDBR can still achieve comparable performance by exploiting dummy data. We also find that the initialized training model lacks feature extraction capabilities, resulting in similar embeddings for different samples in the feature space.



(a) Gradient simulation modes. (b) Model initialization modes.

Fig. 5: Comparison of gradient simulation modes for baselines and model initialization modes. The training dataset is SVHN.

#### 4) Effect of Gradient Simulation and Model Initialization:

Then, we investigate the impact of gradient simulation modes and model initialization modes on the effectiveness of GDBR. The simulation methods are proposed to provide a portion of the ground-truth gradient to ZLG and LLG. Specifically, we consider the following simulation methods.

- *Stats*: Using the statistics of the ground-truth gradient to generate similar distributions for simulation.
- *Stats (col)*: Using the column-wise statistics of the gradient to generate similar distributions for each column.
- *Shuffle (col)*: Shuffling the column-wise elements of the ground-truth gradient for simulation.
- *Mask*: Masking 50% of the gradient and replacing with mean values of the masked elements for simulation.

As shown in Fig. 5a, although the simulation of *Mask* has slightly better performance, GDBR outperforms the baselines

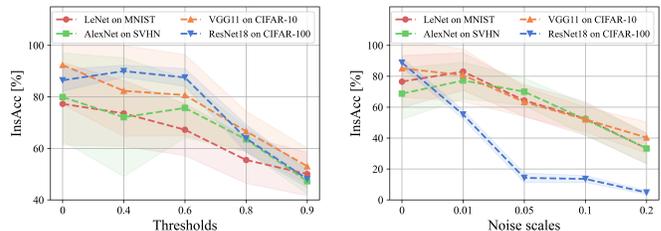
TABLE II: Comparison of label recovery using auxiliary data vs. dummy data across various datasets and models.

Dataset	Model	InsAcc [%]		ClsAcc [%]	
		Aux	Dummy	Aux	Dummy
MNIST	LeNet	81. ± 5.	83. ± 6.	98. ± 4.	95. ± 5.
SVHN	AlexNet	86. ± 5.	80. ± 4.	99. ± 3.	98. ± 4.
CIFAR-10	VGG11	84. ± 6.	85. ± 7.	97. ± 5.	99. ± 2.
CIFAR-100	ResNet18	90. ± 4.	85. ± 3.	98. ± 2.	97. ± 2.
ImageNet	ResNet50	95. ± 2.	92. ± 2.	99. ± 1.	98. ± 2.

\*Aux: auxiliary data, Dummy: dummy data.

with any of the gradient simulation methods.

We also compare the influence of different model initialization modes on GDBR. As introduced in the Implementation Details, we use a uniform distribution between 0.01 and 0.2 to initialize the weights of the bottom layers. This modification ensures that all the features for estimation are positive. The default initialization of PyTorch adopts Kaiming uniform [39], which uses values between the negative and positive bounds. In addition, we replace zero values in the features with the mean of the non-zero elements. We use a 6-layer MLP model trained on SVHN to compare the performance of GDBR with these two initialization modes, utilizing gradients from different layers. The results, shown in Fig. 5b, indicate that GDBR achieves higher accuracy with the positive initialization. However, the performance gap is acceptable, suggesting that GDBR is robust to different initialization modes.



(a) Gradient pruning ([0, 0.9]). (b) Noise perturbation ([0, 0.2]).

Fig. 6: Comparison of gradient pruning and noise perturbation.

#### E. Performance against Defense Mechanisms

Finally, we evaluate the performance of GDBR against two defense mechanisms: gradient pruning and noise perturbation. We compare the performance of GDBR with different gradient pruning thresholds and noise perturbation scales, as shown in Fig. 6. The results demonstrate that the performance of GDBR deteriorates significantly under a high pruning threshold (e.g.,  $\geq 0.9$ ) or a high noise scale (e.g.,  $\geq 0.2$ ). Notably, the Resnet18 model trained on CIFAR-100 is more sensitive to noise than the other models. This increased sensitivity is caused by the accumulation of noise in deducing the gradient w.r.t.  $\mathbf{a}_k$  in the penultimate layer. According to  $\nabla \mathbf{a}_k = \langle \nabla \mathbf{W}_k, \mathbf{W}_k \rangle_F \odot \mathbf{a}_k$ , the noise in  $\nabla \mathbf{W}_k$  is amplified by the Frobenius inner product. Overall, GDBR can still recover the labels with high accuracy

under moderate defense mechanisms, highlighting the necessity of developing more robust defense strategies.

## VI. CONCLUSION

In this work, we propose GDBR, a label recovery attack that reconstructs the label distribution of the victim client's private training data from the limited gradient information shared in FL. By exploring the correlation between the gradients and the model parameters, GDBR builds a bridge between the shared gradients and the target labels and presents the formulation for label restoration. Extensive experiments demonstrate that GDBR can accurately recover at least 80% of the class-wise batch labels in different FL settings. However, GDBR mainly focuses on the early training stage of FL collaboration. In the future, we plan to investigate the effectiveness of GDBR in the trained models and explore effective defenses against GDBR.

## REFERENCES

- [1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [2] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingerman, V. Ivanov, C. Kiddon, J. Konečný, S. Mazzocchi, B. McMahan *et al.*, "Towards federated learning at scale: System design," in *Proceedings of Machine Learning and Systems (MLSys)*, 2019, pp. 374–388.
- [3] Z. Wang, M. Song, Z. Zhang, Y. Song, Q. Wang, and H. Qi, "Beyond inferring class representatives: User-level privacy leakage from federated learning," in *IEEE INFOCOM 2019-IEEE conference on computer communications*. IEEE, 2019, pp. 2512–2520.
- [4] W. Wei, L. Liu, M. Loper, K.-H. Chow, M. E. Gursoy, S. Truex, and Y. Wu, "A framework for evaluating gradient leakage attacks in federated learning," in *ESORICS 2020: 25th European Symposium on Research in Computer Security, ESORICS 2020, Guildford, UK, September 14–18, 2020*. Springer, 2020, pp. 545–566.
- [5] R. Zhang, S. Guo, J. Wang, X. Xie, and D. Tao, "A survey on gradient inversion: Attacks, defenses and future directions," in *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence*, 2023, pp. 5678–685.
- [6] T. Dang, O. Thakkar, S. Ramaswamy, R. Mathews, P. Chin, and F. Beaufays, "Revealing and protecting labels in distributed training," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [7] J. Geng, Y. Mou, F. Li, Q. Li, O. Beyan, S. Decker, and C. Rong, "Towards general deep leakage in federated learning," *arXiv preprint arXiv:2110.09074*, 2021.
- [8] A. Wainakh, F. Ventola, T. Müßig, J. Keim, C. G. Cordero, E. Zimmer, T. Grube, K. Kersting, and M. Mühlhäuser, "User-level label leakage from gradients in federated learning," *Proceedings on Privacy Enhancing Technologies*, vol. 2, pp. 227–244, 2022.
- [9] R. Zhang, S. Guo, and P. Li, "Posterior probability-based label recovery attack in federated learning," in *Privacy Regulation and Protection in Machine Learning*, 2024.
- [10] L. Zhu, Z. Liu, and S. Han, "Deep leakage from gradients," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [11] A. Hatamizadeh, H. Yin, H. R. Roth, W. Li, J. Kautz, D. Xu, and P. Molchanov, "Gradvit: Gradient inversion of vision transformers," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 021–10 030.
- [12] X. Wang, L. Wu, and Z. Guan, "Graddiff: Gradient-based membership inference attacks against federated distillation with differential comparison," *Information Sciences*, vol. 658, p. 120068, 2024.
- [13] S. Mehnaz, S. V. Dibbo, R. De Viti, E. Kabir, B. B. Brandenburg, S. Mangard, N. Li, E. Bertino, M. Backes, E. De Cristofaro *et al.*, "Are your sensitive attributes private? novel model inversion attribute inference attacks on classification models," in *31st USENIX Security Symposium (USENIX Security 22)*, 2022, pp. 4579–4596.
- [14] Y. Aono, T. Hayashi, L. Wang, S. Moriai *et al.*, "Privacy-preserving deep learning via additively homomorphic encryption," *IEEE transactions on information forensics and security*, vol. 13, no. 5, pp. 1333–1345, 2017.
- [15] C. Zhang, S. Li, J. Xia, W. Wang, F. Yan, and Y. Liu, "{BatchCrypt}: Efficient homomorphic encryption for {Cross-Silo} federated learning," in *2020 USENIX annual technical conference (USENIX ATC 20)*, 2020, pp. 493–506.
- [16] Y. Huang, S. Gupta, Z. Song, K. Li, and S. Arora, "Evaluating gradient inversion attacks and defenses in federated learning," *Advances in neural information processing systems*, vol. 34, pp. 7232–7241, 2021.
- [17] W. Wei, L. Liu, J. Zhou, K.-H. Chow, and Y. Wu, "Securing distributed sgd against gradient leakage threats," *IEEE Transactions on Parallel and Distributed Systems*, vol. 34, no. 7, pp. 2040–2054, 2023.
- [18] E. Sothiawat, L. Zhen, Z. Li, and C. Zhang, "Partially encrypted multi-party computation for federated learning," in *2021 IEEE/ACM 21st International Symposium on Cluster, Cloud and Internet Computing (CCGrid)*. IEEE, 2021, pp. 828–835.
- [19] D. Scheliga, P. Mäder, and M. Seeland, "Combining variational modeling with partial gradient perturbation to prevent deep gradient leakage," *arXiv preprint arXiv:2208.04767*, 2022.
- [20] Y. Mei, B. Guo, D. Xiao, and W. Wu, "Fedvf: Personalized federated learning based on layer-wise parameter updates with variable frequency," in *2021 IEEE International Performance, Computing, and Communications Conference (IPCCC)*. IEEE, 2021, pp. 1–9.
- [21] J. Geiping, H. Bauermeister, H. Dröge, and M. Moeller, "Inverting gradients – how easy is it to break privacy in federated learning?" in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020, pp. 16 937–16 947.
- [22] H. Yin, A. Mallya, A. Vahdat, J. M. Alvarez, J. Kautz, and P. Molchanov, "See through gradients: Image batch recovery via gradinversion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 16 337–16 346.
- [23] J. Zhu and M. Blaschko, "R-gap: Recursive gradient attack on privacy," in *International Conference on Learning Representations (ICLR)*, 2021.
- [24] S. Gupta, Y. Huang, Z. Zhong, T. Gao, K. Li, and D. Chen, "Recovering private text in federated learning of language models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 8130–8143, 2022.
- [25] M. Balunovic, D. Dimitrov, N. Jovanović, and M. Vechev, "Lamp: Extracting text from gradients with language model priors," *Advances in Neural Information Processing Systems*, vol. 35, pp. 7641–7654, 2022.
- [26] T. Dang, O. Thakkar, S. Ramaswamy, R. Mathews, P. Chin, and F. Beaufays, "A method to reveal speaker identity in distributed asr training, and how to counter it," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 4338–4342.
- [27] B. Zhao, K. R. Mopuri, and H. Bilen, "idlg: Improved deep leakage from gradients," *arXiv preprint arXiv:2001.02610*, 2020.
- [28] K. Ma, Y. Sun, J. Cui, D. Li, Z. Guan, and J. Liu, "Instance-wise batch label restoration via gradients in federated learning," in *The Eleventh International Conference on Learning Representations*, 2023.
- [29] L. H. Fowl, J. Geiping, W. Czaja, M. Goldblum, and T. Goldstein, "Robbing the fed: Directly obtaining private data in federated learning with modified models," in *International Conference on Learning Representations*, 2021.
- [30] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning internal representations by error propagation, parallel distributed processing, explorations in the microstructure of cognition, ed. de rumelhart and j. mcllelland. vol. 1. 1986," *Biometrika*, vol. 71, no. 599-607, p. 6, 1986.
- [31] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [32] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, 2012.
- [33] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [34] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [35] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, A. Y. Ng *et al.*, "Reading digits in natural images with unsupervised feature learning," in *NIPS workshop on deep learning and unsupervised feature learning*, vol. 2011, no. 2. Granada, 2011, p. 4.
- [36] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," *Technical report, University of Toronto*, 2009.

- [37] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [38] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems*, vol. 32, 2019.
- [39] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1026–1034.