# Concealing Sensitive Samples against Gradient Leakage in Federated Learning

**Jing Wu[1], Munawar Hayat[2], Mingyi Zhou[3], Mehrtash Harandi[1]**

[1]Department of Electrical and Computer Systems Engineering, Monash University
[2]Department of Data Science and AI, Monash University
[3]Department of Software Systems and Cybersecurity, Monash University
{jing.wu1, munawar.hayat, mingyi.zhou, mehrtash.harandi}@monash.edu

## Abstract

Federated Learning (FL) is a distributed learning paradigm that enhances users' privacy by eliminating the need for clients to share raw, private data with the server. Despite the success, recent studies expose the vulnerability of FL to model inversion attacks, where adversaries reconstruct users' private data via eavesdropping on the shared gradient information. We hypothesize that a key factor in the success of such attacks is the low entanglement among gradients per data within the batch during stochastic optimization. This creates a vulnerability that an adversary can exploit to reconstruct the sensitive data. Building upon this insight, we present a simple, yet effective defense strategy that obfuscates the gradients of the sensitive data with concealed samples. To achieve this, we propose synthesizing concealed samples to mimic the sensitive data at the gradient level while ensuring their visual dissimilarity from the actual sensitive data. Compared to the previous art, our empirical evaluations suggest that the proposed technique provides the strongest protection while simultaneously maintaining the FL performance. Code is located at https://github.com/JingWu321/DCS-2.

## 1 Introduction

Consider an Artificial Intelligence (AI) service that aids in disease diagnosis. Multiple hospitals train a model for this service in collaboration. Publishing such a service could benefit a large number of doctors and patients, but it is critical to ensure that private medical data is secure and the utility of the service is normal. Federated Learning (FL) (McMahan et al. 2017a) is an essential technology for such critical applications where the confidentiality of private data is important. FL provides a distributed learning paradigm that enables multiple clients (*e.g.*, hospitals, businesses, or even mobile devices) to train a unified model jointly under the orchestration of a central server. A key advantage of FL lies in its promise of privacy for participating clients. With data decentralized and users' information kept solely with the client, only model updates (*e.g.*, gradients) are transmitted to the central server. Since the model's updates are specifically tailored to the learning task, they may create a false sense of security for FL clients, leading them to believe that the shared updates contain no information on their private training data (Kairouz et al. 2021).

Recent *model inversion attacks* (Zhu, Liu, and Han 2019; Geiping et al. 2020; Balunović et al. 2022; Fowl et al. 2022; Li et al. 2022) have shown that the users' private data can be reconstructed from the gradients shared during the learning process. This alarming finding has led to the exploration of various defense schemes to mitigate privacy leakage. Zhu *et al.* (Zhu, Liu, and Han 2019) employed a strategy that adds noise to gradients, guided by Differential Privacy (DP) (Dwork et al. 2006; Abadi et al. 2016; Song, Chaudhuri, and Sarwate 2013; McMahan et al. 2017b), a concept originally designed to constrain information disclosure. They also utilized gradient compression (Lin et al. 2017), which prunes gradients below a threshold magnitude, as a protective measure. Latest techniques have further advanced the field, with developments such as Automatic Transformation Search (ATS) (Gao et al. 2021) (augmenting data to hide sensitive information), PRivacy EnhanCing mODulE (PRECODE) (Scheliga, Mäder, and Seeland 2022) (use of bottleneck to hide the sensitive data), and Soteria (Sun et al. 2021) (pruning gradients in a single layer).

However, as defense techniques improve, **attacks evolve** as well. New findings, as highlighted by Balunović et al. (2022) and Li et al. (2022), indicate that modern defenses may be ineffective against more sophisticated attacks. For example, Balunović et al. (2022) show that an adversary can disregard the gradients pruned by Soteria and still reconstruct inputs, even without knowledge of the specific layers where pruning is applied. The vulnerability also extends to other defenses; data can be readily reconstructed in the initial communication rounds against the defense ATS (Balunović et al. 2022). In the case of the defense PRECODE, the mere presence of a single non-zero entry in the bias term can enable perfect reconstruction by adversaries (Balunović et al. 2022).

Most current defenses seek to protect all data equally, even if this results in a poor privacy-performance trade-off. In this work, we argue for a more realistic and practical setup where the focus should be given to the sensitive data (*e.g.*, personal data revealing racial or ethnic origin, political opinions, and religious beliefs as mentioned in European Union's General Data Protection Regulation (Voigt and Von dem Bussche 2017)). Consider a malignant skin lesion recognition system as an example. Skin images with tattoos that contain personal information demand extra attention than

images without such information. As such, preserving the former's privacy should be the algorithms' priority.

Exploring the underlying mechanism of model inversion attacks, we hypothesize that these attacks capitalize on the characteristic of relatively low entanglement among the gradients of data points during stochastic optimization. Building upon this understanding, we introduce a defense strategy that obfuscates the gradients of sensitive data using concealed samples. Formally, our goal is to ensure that an adversary is unable to reconstruct sensitive data while simultaneously preserving the performance of the FL system. To achieve this, we propose an algorithm that can adaptively synthesize concealed samples in lieu of sensitive data. We design the concealed points to have high gradient similarity with the sensitive data but visually disparate. For this purpose, our proposed defense has two main characteristics; **1) Enhancing the privacy of sensitive data.** Even though the gradients from the concealed data are similar to those of the sensitive data, inverting these gradients results in data points that are visually very different from the sensitive data. By obfuscating the gradients of the sensitive data with those of the concealed data, the reconstruction of sensitive information becomes confounded, which in turn leads to enhancing the privacy of sensitive data in FL. **2) Maintaining the FL performance.** The introduction of concealed data could potentially disrupt the learning process as it alters the gradient information. Our algorithm mitigates this by ensuring that the shared gradients, after the introduction of concealed data, align with the gradients of the original training samples, including sensitive data. This alignment is achieved through a gradient projection-based approach, preserving the learning capability of the FL system. Unlike existing defenses, our approach proposes a practical solution to enhance privacy in FL. It presents a significant challenge for an adversary to reconstruct the user-defined sensitive samples, all without sacrificing the overall performance of the FL system.

Our main contributions can be summarized as follows:

- We show that model inversion attacks predominantly exploit the characteristic of relatively low entanglement among gradients of samples during stochastic optimization. Based on this finding, we propose to adaptively synthesize concealed samples that obfuscate the gradients of sensitive data.

- The proposed approach crafts concealed samples that are adaptively learned to enhance privacy for sensitive data while simultaneously avoiding performance degradation.

- We thoroughly evaluate and compare our algorithm against various baselines (*e.g.*, injecting noise to the gradients as in the previous works (Sun et al. 2021; Gao et al. 2021; Zhu, Liu, and Han 2019)), and empirically observe that our algorithm consistently outperforms the current state-of-the-art defense methods.

## 2 Related work

**Model Inversion Attacks.** Several model inversion attacks breach FL privacy by reconstructing the clients' data *e.g.*, (Zhu and Blaschko 2020; Fan et al. 2020; Zhu, Liu, and

Han 2019; Yin et al. 2021; Jin et al. 2021; Jeon et al. 2021; Li et al. 2022; Takahashi, Liu, and Liu 2023; Nguyen et al. 2023). Deep Leakage from Gradients (DLG) (Zhu, Liu, and Han 2019) and its variants (Zhao, Mopuri, and Bilen 2020) employ an optimization-based technique to reconstruct private data from the given gradient updates. While the original algorithm (Zhu, Liu, and Han 2019) works best if the number of training samples in each batch is small, subsequent works (Geiping et al. 2020; Wei et al. 2020; Mo et al. 2021; Jeon et al. 2021; Yin et al. 2021) including Gradient Similarity (GS) (Geiping et al. 2020) and GradInversion attack (Yin et al. 2021) are able to reconstruct high-resolution images with larger batch sizes by incorporating stronger image priors. Jin et al. (2021) introduce catastrophic data leakage (CAFE) in vertical federated learning (VFL), showing improved data recovery quality in VFL. Balunović et al. (2022) formalize the gradient leakage problem within the Bayesian framework and demonstrate that the existing optimization-based attacks could be approximated as the optimal adversary with different assumptions on the input and gradients (*ie.*, the prior knowledge about the input and conditional probability of the gradient given the input). They further show that most existing defenses are not quite effective against stronger attacks once appropriate priors (*e.g.*, using generative adversarial networks (Li et al. 2022)) are incorporated to reconstruct data.

While aforementioned optimization-based model inversion attacks assume the server is honest-but-curious (Goldreich 2009), recent works (Fowl et al. 2022; Boenisch et al. 2021) introduce model modification attacks by a malicious server. Boenisch et al. (2021) apply trap weights to initialize the model with the goal of activating parts of its parameters, enabling perfect reconstruction within milliseconds. Similarly, Fowl et al. (2022) proposes the insertion of a tailored imprint module into the network structure. The imprinting module will store information exclusively about a specific subset of data points during the updates, and as a result, data can be recovered precisely and quickly, even when aggregated over large batches.

**Privacy Preserving defenses.** Several approaches propose defense against model inversion attacks that breach users' privacy in FL. We can broadly categorize the existing defenses against model inversion attacks into four categories: gradient compression (Lin et al. 2017; Sun et al. 2021) and perturbation (Dwork et al. 2006; Abadi et al. 2016; Song, Chaudhuri, and Sarwate 2013), data encryption (Gao et al. 2021; Huang et al. 2020), architectural modifications (Scheliga, Mäder, and Seeland 2022), and secure aggregation via changing the communication and training protocol (Bonawitz et al. 2017; Mohassel and Zhang 2017; Lee et al. 2021; Wei et al. 2021) (not considered here). Zhu, Liu, and Han (2019) show that gradient compression can help, while Sun et al. (2021) propose Soteria, suggesting gradient pruning in a single layer as a defense strategy. Zhu, Liu, and Han (2019) also explore adding Gaussian or Laplacian noise guided by DP (Dwork et al. 2006; Abadi et al. 2016; Song, Chaudhuri, and Sarwate 2013; McMahan et al. 2017b) to prevent data being reconstructed. ATS relies on

heavy data augmentation on training images to hide sensitive information, while InstaHide (Huang et al. 2020, 2021) encrypts the private data with data from public datasets. Scheliga, Mäder, and Seeland (2022) introduce PRECODE, which inserts a bottleneck to hide the users' data. Despite these significant efforts to develop defense schemes against FL attacks, recent works highlight the vulnerabilities of existing defenses. For example, several studies show that DP requires a large number of participants in the training process to converge (Zhu, Liu, and Han 2019; Gao et al. 2021; Sun et al. 2021). Balunović et al. (2022) show that an adversary can get an almost perfect reconstruction after dropping the gradients pruned by Soteria. Balunović *et al.* (Balunović et al. 2022) also suggests that it is easy to reconstruct the data using the GS attack in the initial communication rounds against ATS, while Carlini et al. (2020) shows that the private data can be recovered from the encodings of InstaHide (Huang et al. 2020, 2021). For PRECODE, Balunović et al. (2022) demonstrate that an adversary can completely reconstruct the data with at least one non-zero entry in the bias. Further, strong defenses like Soteria can still be bypassed by the Generative Gradient Leakage (GGL) attack method (Li et al. 2022).

# 3  Methodology

In this section, we outline our proposed defense against model inversion attacks. We begin by introducing a basic FL framework, followed by an explanation of a simple reconstruction formulation that illustrates how model inversion attacks operate with shared gradient information. Subsequently, we describe how our proposed approach counters these attacks. Throughout the paper, we denote scalars by lowercase symbols, vectors by bold lowercase symbols, and matrices by bold uppercase symbols (*e.g.*, $a$, $\boldsymbol{a}$, and $\boldsymbol{A}$).

## 3.1  Federated learning

Let $f_{\boldsymbol{\theta}} : \mathcal{X} \to \mathcal{Y}$ be a model with parameters $\boldsymbol{\theta}$, classifying inputs $\boldsymbol{x} \in \mathcal{X}$ to labels $\boldsymbol{y}$ in the label space $\mathcal{Y}$. In FL, we assume that there are $C$ clients and a central server. The data $\mathcal{D}_c$ resides with the client $c$, and the server receives the gradient updates from the clients to update the model parameters $\boldsymbol{\theta}$ as

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{(\boldsymbol{X},\boldsymbol{Y})\sim\mathcal{D}_c}[\mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{X}),\boldsymbol{Y};\boldsymbol{\theta})]. \quad (1)$$

In the $t$-th training round, each client $c$ will compute the gradients $\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{X}),\boldsymbol{Y})$ over local training data and send it to the server. The server then updates the model parameters $\boldsymbol{\theta}^t$ using gradients from the selected $\tilde{C}$ clients:

$$\boldsymbol{\theta}^t = \boldsymbol{\theta}^{t-1} - \frac{\eta}{\tilde{C}} \sum_{c=1}^{\tilde{C}} \nabla_{\boldsymbol{\theta}^{t-1}}\mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{X}),\boldsymbol{Y};\boldsymbol{\theta}^{t-1}), \quad (2)$$

where $\eta$ is the learning rate. The server propagates back the updated parameters $\boldsymbol{\theta}^t$ to each client, repeating the process until convergence. Even though the private training data never leaves the local clients, in the following, we show how an adversary can still reconstruct the data based on the shared gradients $\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{X}),\boldsymbol{Y})$ from client $c$ in the $t$-th communication round.

**Remark.** *If we assume that each client has its own objective, the FL problem can be formulated as*

$$\min_{\boldsymbol{\theta}} \mathbb{E}_{(\boldsymbol{X},\boldsymbol{Y})\sim\mathcal{D}_c}[\mathcal{L}_c(f_{\boldsymbol{\theta}}(\boldsymbol{X}),\boldsymbol{Y};\boldsymbol{\theta})].$$

*Our solution is generic and can also be used to address this scenario.*

## 3.2  Privacy Leakage

**Individual data point leakage.** Without loss of generality, we consider the case of a network having only one fully connected layer, for which the forward pass is given by $\mathbb{R}^m \ni \boldsymbol{y} = \boldsymbol{W}^{\top}\boldsymbol{x} + \boldsymbol{b}$, where $\boldsymbol{W} \in \mathbb{R}^{n\times m}$ is the weight and $\boldsymbol{b} \in \mathbb{R}^m$ is the bias. Let $\mathcal{L}$ denote the objective to update the parameters, then the adversary reconstructs the input $\boldsymbol{x} \in \mathbb{R}^n$ by computing the gradients of the objective w.r.t. the weight and the bias:

$$\nabla_{\boldsymbol{W}}\mathcal{L} = [\frac{\partial\mathcal{L}}{\partial y_1}\frac{\partial y_1}{\partial\boldsymbol{W}_{:1}},\cdots,\frac{\partial\mathcal{L}}{\partial y_m}\frac{\partial y_m}{\partial\boldsymbol{W}_{:m}}],$$
$$\nabla_{\boldsymbol{b}}\mathcal{L} = [\frac{\partial\mathcal{L}}{\partial y_1},\cdots,\frac{\partial\mathcal{L}}{\partial y_m}]. \quad (3)$$

Note that $\frac{\partial y_l}{\partial\boldsymbol{W}_{:l}} = \boldsymbol{x}$ for $1 \le l \le m$. Thus, we can perfectly reconstruct the input from the gradient information as $\boldsymbol{x}^* = \nabla_{\boldsymbol{W}_{:l}}\mathcal{L}/\nabla_{\boldsymbol{b}_l}\mathcal{L} = (\frac{\partial\mathcal{L}}{\partial y_l}\frac{\partial y_l}{\partial\boldsymbol{W}_{:l}})/\frac{\partial\mathcal{L}}{\partial y_l} = \boldsymbol{x}$, provided that at least one element of the gradient of the loss with respect to the bias is non-zero (*ie.*, $\frac{\partial\mathcal{L}}{\partial y_l} \neq 0, 1 \le l \le m$).

**Multiple data points leakage.** Let $\boldsymbol{x}_j$, $j \in [1,B], B > 1$ denotes samples of a mini-batch of size $B$. The gradient of the mini-batch is:

$$\nabla_{\boldsymbol{W}}\mathcal{L} = \frac{1}{B}\sum_{j=1}^{B}[\frac{\partial\mathcal{L}}{\partial y_{1,j}}\frac{\partial y_{1,j}}{\partial\boldsymbol{W}_{:1}},\cdots,\frac{\partial\mathcal{L}}{\partial y_{m,j}}\frac{\partial y_{m,j}}{\partial\boldsymbol{W}_{:m}}],$$
$$\nabla_{\boldsymbol{b}}\mathcal{L} = \frac{1}{B}\sum_{j=1}^{B}[\frac{\partial\mathcal{L}}{\partial y_{1,j}},\cdots,\frac{\partial\mathcal{L}}{\partial y_{m,j}}], \quad (4)$$

which encapsulates a linear combination of all data points $\boldsymbol{x}_j$ in the mini-batch. Sun et al. (2021) observe that for data coming from different classes, the corresponding data representations tend to be embedded in different rows/columns of gradients. Suppose that within the mini-batch, only $\boldsymbol{x}_1$ belongs to class $y_c$ ($1 \le c \le m$), then the column $c$ of the gradient in Eq. (4) will have

$$\frac{\sum_{j=1}^{B}\frac{\partial\mathcal{L}}{\partial y_{c,j}}\frac{\partial y_{c,j}}{\partial\boldsymbol{W}_{:c}}}{\sum_{j=1}^{B}\frac{\partial\mathcal{L}}{\partial y_{c,j}}} \approx \frac{\frac{\partial\mathcal{L}}{\partial y_{c,1}}\frac{\partial y_{c,1}}{\partial\boldsymbol{W}_{:c}}}{\frac{\partial\mathcal{L}}{\partial y_{c,1}}} = \boldsymbol{x}_1. \quad (5)$$

Due to this property, *ie.* relatively low entanglement among gradients per data points within a batch, the adversary can reconstruct the data in practice.

Boenisch *et al.* (Boenisch et al. 2021) also observe that for a ReLU network, over-parameterization can cause all but one training data in a mini-batch to have zero gradients, allowing the individual data point leakage in the mini-batch and the passive adversaries to obtain perfect reconstruction in various cases.

Optimization-based attacks aim to reconstruct data by minimizing the distance between the gradient of the input and that of the reconstruction. In contrast, model modification attacks utilize specific parameters with the goal of amplifying the leakage of individual data points (Boenisch et al. 2021) within the mini-batch or allowing portions of the gradient to contain information exclusive to a subset of data points (Fowl et al. 2022). It is important to note that neither optimization-based attacks nor model modification attacks can precisely separate the gradient for individual data points. **This limitation in the attack algorithms is a vulnerability that we leverage in our approach to protect the data.**

## 3.3 Defense by Concealing Sensitive Samples (DCS$^2$)

Our objective is to protect sensitive data without modifying any FL settings (*e.g.*, model structure) and the sensitive data themselves, while minimizing the impact of the proposed defense on the model performance. Previously, we discussed that model inversion attacks reconstruct the inputs using the gradient information since the gradient encapsulates sufficient information about data samples to reconstruct them (see Eq. (4)). We note that while theoretically, attacks cannot precisely separate the gradient for each sample, they can be extremely successful in practice. Our key insight is to insert samples (referred to as concealed samples) to imitate the sensitive data on the gradient level while ensuring that these samples are visually dissimilar to the sensitive data. Our goal is to make it difficult or even impossible for the adversary to distinguish the gradient of the synthesized concealed samples from the gradient of the sensitive data.

Without loss of generality, assume that there is only one sensitive data point, denoted by $\boldsymbol{x}_s$. Our task is to construct the concealed sample $\tilde{\boldsymbol{x}}_c$ for this sensitive data to achieve the following goals as part of our defense strategy:

**Goal-1:** To protect sensitive data from model inversion attacks, we would like to maximize the dissimilarity between the concealed sample $\tilde{\boldsymbol{x}}_c$ and the sensitive sample $\boldsymbol{x}_s$, as measured by $\|\tilde{\boldsymbol{x}}_c - \boldsymbol{x}_s\|$. Simultaneously, we seek to minimize the similarity between the gradient of the concealed sample w.r.t. sensitive data. This is quantified by the cosine similarity between the gradient vectors, *ie.*, $\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\tilde{\boldsymbol{x}}_c), \tilde{\boldsymbol{y}}_c)$ and $\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{x}_s), \boldsymbol{y}_s)$, while ensuring that the resulting latent representation is similar to the sensitive latent representation, *ie.*, $\|f_{\boldsymbol{\theta}}(\tilde{\boldsymbol{x}}_c) - f_{\boldsymbol{\theta}}(\boldsymbol{x}_s)\| \leq \epsilon$.

**Goal-2:** To facilitate the server's ability to learn and enhance the FL model, we must ensure that the resulting gradient closely resembles the gradient of the batch without concealed samples. This can be achieved by satisfying $\langle\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\{\boldsymbol{x}_s\} \cup \{\tilde{\boldsymbol{x}}_c\}), \{\boldsymbol{y}_s\} \cup \{\tilde{\boldsymbol{y}}_c\}), \nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{x}_s), \boldsymbol{y}_s)\rangle > 0$.

To accomplish the aforementioned goals, our defense strategy consists of two phases: **1.** *synthesizing the concealed samples* and **2.** *gradient projection*, which we discuss below.

**Synthesizing the concealed samples.** To obtain concealed samples that are visually dissimilar to sensitive data

but whose gradient is similar to the sensitive data, we would like to solve the following optimization problem:

$$\min_{\tilde{\boldsymbol{x}}_c} 1 - \frac{\langle\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\tilde{\boldsymbol{x}}_c), \tilde{\boldsymbol{y}}_c), \nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{x}_s), \boldsymbol{y}_s)\rangle}{\|\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\tilde{\boldsymbol{x}}_c), \tilde{\boldsymbol{y}}_c)\| \|\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{x}_s), \boldsymbol{y}_s)\|} \quad (6)$$

$$\max_{\tilde{\boldsymbol{x}}_c} \|\tilde{\boldsymbol{x}}_c - \boldsymbol{x}_s\| \quad (7)$$

$$\text{s.t.} \|f_{\boldsymbol{\theta}}(\tilde{\boldsymbol{x}}_c) - f_{\boldsymbol{\theta}}(\boldsymbol{x}_s)\| \leq \epsilon. \quad (8)$$

We propose the following objective to achieve this

$$\mathcal{L}_{obj} = (1 - \frac{\langle\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\tilde{\boldsymbol{x}}_c), \tilde{\boldsymbol{y}}_c), \nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{x}_s), \boldsymbol{y}_s)\rangle}{\|\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\tilde{\boldsymbol{x}}_c), \tilde{\boldsymbol{y}}_c)\| \times \|\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{x}_s), \boldsymbol{y}_s)\|})$$
$$+ e^{-\lambda_x\|\tilde{\boldsymbol{x}}_c - \boldsymbol{x}_s\|} + \lambda_z(\frac{\|f_{\boldsymbol{\theta}}(\tilde{\boldsymbol{x}}_c) - f_{\boldsymbol{\theta}}(\boldsymbol{x}_s)\|}{\|f_{\boldsymbol{\theta}}(\boldsymbol{x}_s)\|} - \epsilon), \quad (9)$$

where $\lambda_z$ and $\lambda_x$ are hyperparameters to balance the different terms in the objective, $\epsilon$ controls the latent distance. The first term and third term target achieving Goal-1 by ensuring that the concealed sample is similar to the sensitive data at the gradient level, while the second term learns the concealed sample to be visually dissimilar to the sensitive data.

**Remark.** *The label corresponding to the concealed sample $\tilde{\boldsymbol{x}}_c$ is denoted by $\tilde{\boldsymbol{y}}_c$ in Eq. (9). To obtain $\tilde{\boldsymbol{x}}_c$, we solve an optimization problem, starting from $\boldsymbol{x}_0$, which may be a sample different from $\boldsymbol{x}_s$. In such cases, we assign $\tilde{\boldsymbol{y}}_c$ with the label of $\boldsymbol{x}_0$, ie., $\tilde{\boldsymbol{y}}_c = \boldsymbol{y}_0$. In our experiments, we show that $\tilde{\boldsymbol{x}}_c$ can be randomly initialized, and accordingly, we set $\tilde{\boldsymbol{y}}_c$ at random. Our empirical evaluations in § 4 show that the proposed method works equally well under both conditions.*

**Gradient projection.** Using Eq. (9), we can obtain the concealed sample $\boldsymbol{x}_c$. What we need to do next is to ensure that the gradient of the mini-batch augmented with the concealed sample is aligned with the gradient of the original mini-batch, as this way, the server can improve its model. This will be achieved via the gradient projection, but before delving into details of projection and inspired by the mixup regularization (Zhang et al. 2017), we propose an enhancement. Let $\boldsymbol{g}$ be the gradient of the original mini-batch $\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{x}_s), \boldsymbol{y}_s)$. We obtain the gradient with the concealed sample as

$$\boldsymbol{g}_c \triangleq \nabla_{\boldsymbol{\theta}}\Big\{\mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{x}_s), \boldsymbol{y}_s) + \lambda_g\mathcal{L}(f_{\boldsymbol{\theta}}(\tilde{\boldsymbol{x}}_c), \tilde{\boldsymbol{y}}_c)$$
$$+ (1 - \lambda_g)\mathcal{L}(f_{\boldsymbol{\theta}}(\tilde{\boldsymbol{x}}_c), \boldsymbol{y}_s)\Big\}, \quad (10)$$

where $\lambda_g$ is a hyperparameter. Note that if $\lambda_g = 1$, we indeed attain the gradient of the mini-batch augmented by the concealed sample. However, including the gradient in the form $\nabla_{\boldsymbol{\theta}}\mathcal{L}(f_{\boldsymbol{\theta}}(\tilde{\boldsymbol{x}}_c), \boldsymbol{y}_s)$ is empirically observed to be beneficial. Analysis can be found in § 4.

To align the resulting gradient $\boldsymbol{g}_c$ with the original gradient of the mini-batch $\boldsymbol{g}$, we opt for the technique developed in (Lopez-Paz and Ranzato 2017). This will ensure that the gradient sent to the server will improve the FL model. To this end, we compute the angle between the original gradient vector and the new gradient and check if it satisfies $\langle\boldsymbol{g}, \boldsymbol{g}_c\rangle \geq 0$. If the constraints is satisfied, the new gradient $\boldsymbol{g}_c$ behaves similarly to that of obtained from the mini-batch

**Algorithm 1:** Defense by Concealing Sensitive Samples (DCS$^2$)

1: **procedure** GRADIENT OBFUSCATION
2:     initialize the start point for constructing
            the concealed data $\tilde{\boldsymbol{x}}_c \leftarrow \boldsymbol{x}_0, \tilde{\boldsymbol{y}}_c \leftarrow \boldsymbol{y}_0$;
3:     get the concealed sample $\tilde{\boldsymbol{x}}_c \leftarrow Eq.$ (9);
4:     compute the new gradient $\boldsymbol{g}_c \leftarrow Eq.$ (10);
5: **procedure** GRADIENT PROJECTION
6:     get the gradient from the original batch
            $\boldsymbol{g} \leftarrow \nabla_{\boldsymbol{\theta}} \mathcal{L}(f_{\boldsymbol{\theta}}(\boldsymbol{x}_s), \boldsymbol{y}_s)$;
7:     **if** $\langle \boldsymbol{g}, \boldsymbol{g}_c \rangle < 0$ **then**
8:         get the solution $v^* \leftarrow Eq.$ (12);
9:         project the new gradient to the closest gradient
                $\hat{\boldsymbol{g}}_c = \boldsymbol{g}v^* + \boldsymbol{g}_c$.

$\boldsymbol{x}_s$; otherwise, we project the new gradient $\boldsymbol{g}_c$ to the closest gradient $\hat{\boldsymbol{g}}_c$ according to:

$$\underset{\hat{\boldsymbol{g}}_c}{\arg\min} \quad \frac{1}{2} \|\boldsymbol{g}_c - \hat{\boldsymbol{g}}_c\|_2^2,$$
$$s.t. \quad \langle \boldsymbol{g}, \hat{\boldsymbol{g}}_c \rangle \geq 0. \tag{11}$$

To efficiently solve Eq. (11), we employ the Quadratic Programming (QP) with inequality constraints:

$$\underset{v}{\arg\min} \quad \frac{1}{2}\boldsymbol{g}^\top \boldsymbol{g}v + \boldsymbol{g}_c^\top \boldsymbol{g}v,$$
$$\text{s.t.} \quad v \geq 0. \tag{12}$$

The projected gradient $\hat{\boldsymbol{g}}_c$ is given from the solution $v^*$ in Eq. (12) as $\hat{\boldsymbol{g}}_c = \boldsymbol{g}v^* + \boldsymbol{g}_c$. The complete pseudocode for the algorithm is provided in Algorithm 1.

# 4 Experiments

In this section, we first describe our evaluation settings, followed by a comparison of our defense with existing defenses against model inversion attacks in FL, to answer the following research questions (RQs):

**RQ1:** Can the proposed method DCS$^2$ effectively protect sensitive data against model inversion attacks in FL?

**RQ2:** Is the proposed method DCS$^2$ capable of maintaining FL performance while providing protection?

**RQ3:** How does the proposed method DCS$^2$ compare with existing defenses?

**RQ4:** How does the proposed method DCS$^2$ perform when defending against adaptive attacks?

**RQ5:** How does the proposed method DCS$^2$ perform when the starting point for generating concealing samples varies?

Additional details, including the values of hyperparameters, are available in the supplementary material.

## 4.1 Experimental setup

**Attack methods.** We evaluate defenses against classical and state-of-the-art (SOTA) attacks in FL: the improved version of the classical Deep Leakage from Gradients (Zhu, Liu, and Han 2019) called *GS attack* (Geiping et al. 2020) that introduces image prior and uses cosine similarity as a distance metric to enhance reconstruction, and SOTA attack *GGL attack* that uses a Generative Adversarial Network

(GAN) to learn prior knowledge from public datasets. We also include the recently proposed SOTA model modification attack *ie. Imprint attack* (Fowl et al. 2022). Furthermore, we provide an evaluation when the *adaptive attack* has strong prior knowledge about the private training data.

**Defense baselines.** Following recent works (Sun et al. 2021; Gao et al. 2021), we compare our approach with defenses including *DP-Gaussian* (adding Gaussian noise to gradients, following the implementation in (Sun et al. 2021; Gao et al. 2021)), and *Prune* (Gradient Compression) (Lin et al. 2017). We further compare against the recently proposed defense *Soteria* (Sun et al. 2021), which perturbs the representations. In the supplementary material, we also provide a comparison with defenses that alter the sensitive data *e.g.*, ATS (Gao et al. 2021).

**Datasets.** We consider four datasets, namely MNIST (Le-Cun et al. 1998) with image resolution $28 \times 28$, CIFAR10 (Krizhevsky, Hinton et al. 2009) with image resolution $32 \times 32$, CelebFaces Attributes (CelebA) Dataset (Liu et al. 2015) with image resolution rescaled to $32 \times 32$ for a fair evaluation on GGL attack and TinyImageNet (Le and Yang 2015) with image resolution rescaled to $224 \times 224$.

**Models.** Being consistent with existing literature, we consider three model architectures i.e., LeNet (LeCun et al. 1998) for MNIST, ConvNet (with the same structure as in Soteria (Sun et al. 2021)) for CIFAR10 and CelebA, ResNet18 (He et al. 2016) for TinyImageNet.

**Metrics.** To quantify the quality of reconstructed images and compare them with the original sensitive data, we use peak signal-to-noise ratio (PSNR) as used in the work (Balunović et al. 2022), and structural similarity index measure (SSIM) (Wang et al. 2004). Besides, we use the learned perceptual image patch similarity (LPIPS) metric (Zhang et al. 2018) for experiments on TinyImageNet. When measuring PSNR and SSIM, lower values indicate better performances. When it comes to LPIPS, a higher number indicates a better performance. We report classification accuracy values on the respective test sets and the protected data to measure the FL performance.

**FL setting.** We apply FedAvg (McMahan et al. 2017a) during training to report our results. We set the local epoch $E$ as 1 (easier for attacks) and batch size $B$ as 64. We have 10 clients in total; each client only has 200 samples on MNIST, 2000 samples on CIFAR10, 500 samples on CelebA, and 2000 samples on TinyImageNet.

## 4.2 Privacy-performance trade-off

We consider 100% of the training data in the target client as sensitive samples. The optimal conditions for an adversary to invert gradients are a batch size of one, a low image resolution, and an untrained target network.

**Results on MNIST and CIFAR10.** We first evaluate defenses against the GS attack on the MNIST and CIFAR10 datasets using models with randomly initialized weights. Results on Tab. 1 indicate that, compared with existing defenses, our proposed approach provides a better defense

| | MNIST | | | | CIFAR10 | | | |
|---|---|---|---|---|---|---|---|---|
| Defense | PSNR↓ | SSIM↓ | Acc↑ (Sensitive Data) | Acc↑ (Test set) | PSNR↓ | SSIM↓ | Acc↑ (Sensitive Data) | Acc↑ (Test set) |
| None | $59.20_{\pm2.71}$ | $1.00_{\pm4.87}$ | $86.98_{\pm0.00}$ | $87.16_{\pm0.01}$ | $20.41_{\pm3.15}$ | $0.73_{\pm0.09}$ | $90.35_{\pm0.04}$ | $80.41_{\pm0.01}$ |
| DP-Gaussian | $35.38_{\pm2.44}$ | $0.83_{\pm0.07}$ | $85.94_{\pm0.00}$ | $86.91_{\pm0.01}$ | $12.34_{\pm1.34}$ | $0.28_{\pm0.06}$ | $77.19_{\pm0.18}$ | $79.65_{\pm0.04}$ |
| Prune | $14.13_{\pm2.29}$ | $0.37_{\pm0.06}$ | $85.94_{\pm0.00}$ | $86.91_{\pm0.00}$ | $11.26_{\pm1.75}$ | $0.22_{\pm0.06}$ | $77.80_{\pm0.32}$ | $79.51_{\pm0.08}$ |
| Soteria | $9.67_{\pm1.09}$ | $0.30_{\pm0.07}$ | $86.98_{\pm0.00}$ | $86.94_{\pm0.00}$ | $11.48_{\pm1.42}$ | $0.29_{\pm0.06}$ | $\mathbf{84.70_{\pm0.32}}$ | $79.76_{\pm0.04}$ |
| DCS$^2$ (Ours) | $\mathbf{7.84_{\pm2.56}}$ | $\mathbf{0.17_{\pm0.09}}$ | $\mathbf{86.98_{\pm0.00}}$ | $\mathbf{86.98_{\pm0.01}}$ | $\mathbf{8.04_{\pm1.10}}$ | $\mathbf{0.15_{\pm0.05}}$ | $80.39_{\pm0.07}$ | $\mathbf{79.79_{\pm0.03}}$ |

Table 1: Defenses against GS attack on MNIST and CIFAR10. Values are averaged. For DP-Gaussian, we follow the implementation in studies (Sun et al. 2021; Gao et al. 2021).

| | CelebA | | | | TinyImageNet | | | |
|---|---|---|---|---|---|---|---|---|
| Defense | PSNR↓ | SSIM↓ | Acc↑ (Sensitive Data) | Acc↑ (Test set) | SSIM↓ | LPIPS↑ | Acc↑ (Sensitive Data) | Acc↑ (Test set) |
| None | $19.92_{\pm2.18}$ | $0.75_{\pm0.07}$ | $100.0_{\pm0.00}$ | $93.79_{\pm0.07}$ | $1.00_{\pm0.00}$ | $0.00_{\pm0.00}$ | $73.94_{\pm1.21}$ | $66.41_{\pm0.02}$ |
| DP-Gaussian | $13.95_{\pm1.52}$ | $0.44_{\pm0.08}$ | $90.51_{\pm0.47}$ | $93.19_{\pm0.04}$ | $1.00_{\pm0.00}$ | $0.00_{\pm0.00}$ | $53.28_{\pm0.78}$ | $65.65_{\pm0.07}$ |
| Prune | $9.57_{\pm2.66}$ | $0.24_{\pm0.12}$ | $91.41_{\pm1.10}$ | $93.25_{\pm0.06}$ | $0.91_{\pm0.12}$ | $0.16_{\pm0.20}$ | $52.77_{\pm0.07}$ | $\mathbf{65.73_{\pm0.20}}$ |
| Soteria | $8.89_{\pm2.63}$ | $0.24_{\pm0.11}$ | $100.0_{\pm0.00}$ | $93.86_{\pm0.01}$ | $1.00_{\pm0.00}$ | $0.00_{\pm0.00}$ | $41.84_{\pm1.14}$ | $52.06_{\pm1.47}$ |
| DCS$^2$ (Ours) | $\mathbf{8.24_{\pm2.71}}$ | $\mathbf{0.17_{\pm0.12}}$ | $100.0_{\pm0.00}$ | $\mathbf{94.31_{\pm0.01}}$ | $\mathbf{0.79_{\pm0.22}}$ | $\mathbf{0.22_{\pm0.23}}$ | $59.88_{\pm0.71}$ | $65.68_{\pm0.05}$ |

Table 2: Defenses against GGL attack on CelebA and Imprint attack on TinyImageNet. Values are averaged.

| $\lambda_g$ | SSIM↓ | LPIPS↑ | Acc↑ (Sensitive Data) | Acc↑ (Test set) |
|---|---|---|---|---|
| 0.5 | $0.80_{\pm0.20}$ | $0.22_{\pm0.21}$ | $\mathbf{60.33_{\pm0.71}}$ | $\mathbf{65.76_{\pm0.04}}$ |
| 0.7 | $0.79_{\pm0.22}$ | $0.22_{\pm0.23}$ | $59.88_{\pm0.71}$ | $65.68_{\pm0.05}$ |
| 1.0 | $\mathbf{0.78_{\pm0.22}}$ | $\mathbf{0.23_{\pm0.23}}$ | $58.54_{\pm0.46}$ | $65.24_{\pm0.21}$ |

Table 3: DCS$^2$ with different $\lambda_g$ on TinyImageNet.

| Defense | PSNR↓ | SSIM↓ |
|---|---|---|
| None | $59.22_{\pm2.71}$ | $1.00_{\pm4.77}$ |
| DCS$^2$ | $7.87_{\pm2.44}$ | $0.18_{\pm0.09}$ |



Figure 1 & Table 4: Defend against adaptive attacks.

against the GS attack. Specifically, on MNIST, the defense baselines reduce the PSNR from 59.20 to ∼ 10, while our defense can reduce the PSNR to around 8. On CIFAR10, our method reduces the SSIM to 0.17 when other defenses only reduce it to around 0.3. In terms of the FL performance, as shown in Tab. 1, our proposed defense method DCS$^2$ largely retains the performance compared with other defenses. Specifically, on MNIST, when most defense baseline drops the performance by about 1% on the sensitive data, our defense maintains the performance.

**Results on CelebA and TinyImageNet.** Further, we compare different defenses for more complex datasets, with larger capacity networks, on CelebA and TinyImageNet, to defend against stronger attacks. We use randomly initialized weights and use the attribute gender as the target label in CelebA to perform binary classification. A pre-trained ResNet18 was applied for TinyImageNet. As shown in Table 2, our defense provides the best protection while competitively maintaining the original FL performance. Specifically, on CelebA, defending against the GGL attack, our method provides the best protection, and the FL performance is even improved while defenses DP-Gaussian and Prune drop by around 0.5% on the test set. On TinyImageNet, when defending against the Imprint attack, the defense Soteria cannot know where the adversary would insert the imprint module, so it cannot withstand the Imprint attack. While most defenses cannot provide protection, our defense method increases the LPIPS from 0.00 to 0.22.

Fig. 2 shows the example of reconstructions from different attacks with defenses on different datasets. The attacks could still recover some parts of the sensitive data with other defenses, while they fail with our proposed defense method. The training process on various datasets with different defenses is illustrated in Fig. 4. Training with these defenses typically results in convergence. However, in the case of Soteria on TinyImageNet, approximately 90% of the representations are perturbed, resulting in a convergence failure.

Tab. 3 presents the results for DCS$^2$ on TinyImageNet under varying values of $\lambda_g$. As $\lambda_g$ increases, the protection for sensitive data improves. However, this leads to a reduction in the performance of the FL system.

### 4.3 Comparison against adaptive attacks.

We compared the proposed defense method DCS$^2$ against two SOTA attacks: Imprint and GGL. Imprint modifies the architecture, and GGL uses a GAN to learn prior knowledge from public datasets. As per Gao et al. (2021), both these attacks are adaptive since the adversary "starts the reconstruction from an image with certain semantic information" or "designs attack techniques instead of optimizing the distance between the real and dummy gradients". Results in Tabs. 1 and 2 indicate that our defense provides the best protection with minimal drop in accuracy. For example, on TinyImageNet, our defense reduces the SSIM score from 1.0 to 0.79. In comparison, the defense Prune decreases it to approximately 0.9 and other defenses prove inadequate against this attack. The accuracy of the FL system using our defense
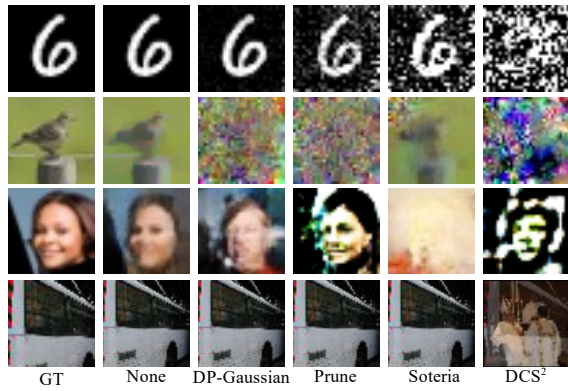
Figure 2: Example of reconstructions for Tabs. 1 and 2. From top to bottom are reconstructions from GS attacks on MNIST, those from GS attacks on CIFAR10, those from GGL attacks on CelebA, and those from Imprint attacks on TinyImageNet, respectively. (Best viewed in color)

| MNIST | Noise | MixUP | PSNR↓ | SSIM↓ | Acc↑ (Sensitive Data) | Acc↑ (Test set) |
|---|---|---|---|---|---|---|
| ✓ | ✗ | ✓ | $7.97_{\pm 2.48}$ | $0.18_{\pm 0.08}$ | $86.56_{\pm 0.57}$ | $86.99_{\pm 0.01}$ |
| ✓ | ✗ | ✗ | $7.84_{\pm 2.56}$ | $0.17_{\pm 0.09}$ | $86.98_{\pm 0.00}$ | $86.98_{\pm 0.01}$ |
| ✗ | ✓ | ✓ | $7.69_{\pm 2.38}$ | $0.18_{\pm 0.08}$ | $\mathbf{86.98_{\pm 0.00}}$ | $86.99_{\pm 0.00}$ |
| ✗ | ✓ | ✗ | $\mathbf{7.40_{\pm 2.22}}$ | $\mathbf{0.16_{\pm 0.08}}$ | $85.94_{\pm 0.00}$ | $86.94_{\pm 0.00}$ |

Table 5: Different start points on MNIST.

on the sensitive data decreases by about 14%, whereas other defenses drop exceeding 20%.

Further, we design another strong attack where the adversary has strong prior knowledge and initializes the GS attack with the average image for each class. Results are shown in Tab. 4, our proposed method can still provide good protection against such an attack with prior knowledge about the sensitive data. Fig. 1 shows an example of the reconstructions from this attack. The GS attack would initialize the dummy input with the AvgImg (average image) shown in Fig. 1. The average image already explicitly includes information about the sensitive data, while our defense method could still protect the data against this adaptive attack.

## 4.4 Effect of Starting Points on Generating Concealed Samples

We further evaluate our defense by choosing different initial starting points to craft the concealed samples. Tab. 5 show the performance with different start points. 'MixUP' means that $\tilde{x}_c$ is initialized with $0.7x_0 + 0.3x_s$. Tab. 6 show the results when the start points are from CIFAR10, which has different distribution than the target task dataset CelebA. As shown in Tabs. 5 and 6, even starting from random noise and different domains, our defense method could still provide protection and retain the model's performance.

| Defense | PSNR↓ | SSIM↓ | Acc↑ |
|---|---|---|---|
| None | $19.92_{\pm 2.18}$ | $0.75_{\pm 0.07}$ | $93.79_{\pm 0.07}$ |
| DCS$^2$ | $8.68_{\pm 2.78}$ | $0.18_{\pm 0.12}$ | $94.13_{\pm 0.03}$ |



Figure 3 & Table 6: Start points from CIFAR10 for CelebA.



Figure 4: Training process. (Best viewed in color)

| Defense | PSNR↓ | SSIM↓ | Acc↑ (Sensitive Data) | Acc↑ (Test set) |
|---|---|---|---|---|
| Prune | $14.13_{\pm 2.29}$ | $0.37_{\pm 0.06}$ | $85.94_{\pm 0.00}$ | $86.91_{\pm 0.00}$ |
| DCS$^2$ | $7.84_{\pm 2.56}$ | $0.17_{\pm 0.09}$ | $\mathbf{86.98_{\pm 0.00}}$ | $\mathbf{86.98_{\pm 0.01}}$ |
| Prune&DCS$^2$ | $\mathbf{6.08_{\pm 1.60}}$ | $\mathbf{0.12_{\pm 0.06}}$ | $86.15_{\pm 0.47}$ | $86.92_{\pm 0.01}$ |

Table 7: Combination of defenses.

## 4.5 Combination with existing defenses

An illustration of combining DCS$^2$ with the defense 'Prune' is presented in Tab. 7. In this scenario, the enhancement of protection for private training data is notable. While the performance experiences a slight decrease compared to the standalone proposed defense method, it still surpasses the performance of the defense 'Prune' alone.

## 5 Limitations

While our empirical evaluations show that our proposed defense is effective in enhancing privacy and retaining FL performance, it requires additional computation to craft concealed samples (refer to the supplementary material for details on compute complexity). Future directions to improve concealed samples-based defense include finding the best starting points and reducing the time to craft the concealed samples. We hope our defense can provide a new perspective for defending against model inversion attacks in FL.

## 6 Conclusion

In this work, we proposed an effective defense algorithm against model inversion attacks in FL. Our approach crafts concealed samples that imitate the sensitive data, but can obfuscate their gradients, thus making it challenging for an adversary to reconstruct sensitive data from the shared gradients. To enhance the privacy of the sensitive data, the concealed samples are adaptively learned to be visually very dissimilar to the sensitive samples, while their gradients are aligned with the original samples to avoid FL performance drop. Our evaluations on four benchmark datasets showed that, compared with other defenses, our approach offers the best protection against model inversion attacks while simultaneously retaining or even improving the FL performance.

# 7 Acknowledgments

# References

Abadi, M.; Chu, A.; Goodfellow, I.; McMahan, H. B.; Mironov, I.; Talwar, K.; and Zhang, L. 2016. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, 308–318.

Balunović, M.; Dimitrov, D. I.; Staab, R.; and Vechev, M. 2022. Bayesian Framework for Gradient Leakage. *ICLR*.

Boenisch, F.; Dziedzic, A.; Schuster, R.; Shamsabadi, A. S.; Shumailov, I.; and Papernot, N. 2021. When the Curious Abandon Honesty: Federated Learning Is Not Private. *arXiv preprint arXiv:2112.02918*.

Bonawitz, K.; Ivanov, V.; Kreuter, B.; Marcedone, A.; McMahan, H. B.; Patel, S.; Ramage, D.; Segal, A.; and Seth, K. 2017. Practical secure aggregation for privacy-preserving machine learning. In *proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, 1175–1191.

Carlini, N.; Deng, S.; Garg, S.; Jha, S.; Mahloujifar, S.; Mahmoody, M.; Song, S.; Thakurta, A.; and Tramer, F. 2020. Is Private Learning Possible with Instance Encoding? *arXiv preprint arXiv:2011.05315*.

Dwork, C.; McSherry, F.; Nissim, K.; and Smith, A. 2006. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, 265–284. Springer.

Fan, L.; Ng, K. W.; Ju, C.; Zhang, T.; Liu, C.; Chan, C. S.; and Yang, Q. 2020. Rethinking privacy preserving deep learning: How to evaluate and thwart privacy attacks. In *Federated Learning*, 32–50. Springer.

Fowl, L.; Geiping, J.; Czaja, W.; Goldblum, M.; and Goldstein, T. 2022. Robbing the Fed: Directly Obtaining Private Data in Federated Learning with Modified Models. *ICLR*.

Gao, W.; Guo, S.; Zhang, T.; Qiu, H.; Wen, Y.; and Liu, Y. 2021. Privacy-preserving collaborative learning with automatic transformation search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 114–123.

Geiping, J.; Bauermeister, H.; Dröge, H.; and Moeller, M. 2020. Inverting gradients-how easy is it to break privacy in federated learning? *Advances in Neural Information Processing Systems*, 33: 16937–16947.

Goldreich, O. 2009. *Foundations of cryptography: volume 2, basic applications*. Cambridge university press.

Griewank, A.; and Walther, A. 2008. *Evaluating derivatives: principles and techniques of algorithmic differentiation*. SIAM.

Gustafson, L.; Rolland, C.; Ravi, N.; Duval, Q.; Adcock, A.; Fu, C.-Y.; Hall, M.; and Ross, C. 2023. FACET: Fairness in Computer Vision Evaluation Benchmark. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 20370–20382.

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.

Huang, Y.; Gupta, S.; Song, Z.; Li, K.; and Arora, S. 2021. Evaluating gradient inversion attacks and defenses in federated learning. *Advances in Neural Information Processing Systems*, 34.

Huang, Y.; Song, Z.; Li, K.; and Arora, S. 2020. Instahide: Instance-hiding schemes for private distributed learning. In *International Conference on Machine Learning*, 4507–4518. PMLR.

Jeon, J.; Lee, K.; Oh, S.; Ok, J.; et al. 2021. Gradient inversion with generative image prior. *Advances in Neural Information Processing Systems*, 34: 29898–29908.

Jin, X.; Chen, P.-Y.; Hsu, C.-Y.; Yu, C.-M.; and Chen, T. 2021. Catastrophic Data Leakage in Vertical Federated Learning. *Advances in Neural Information Processing Systems*, 34.

Kairouz, P.; McMahan, H. B.; Avent, B.; Bellet, A.; Bennis, M.; Bhagoji, A. N.; Bonawitz, K.; Charles, Z.; Cormode, G.; Cummings, R.; et al. 2021. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2): 1–210.

Krizhevsky, A.; Hinton, G.; et al. 2009. Learning multiple layers of features from tiny images.

Le, Y.; and Yang, X. 2015. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7): 3.

LeCun, Y.; Bottou, L.; Bengio, Y.; and Haffner, P. 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11): 2278–2324.

Lee, H.; Kim, J.; Ahn, S.; Hussain, R.; Cho, S.; and Son, J. 2021. Digestive neural networks: A novel defense strategy against inference attacks in federated learning. *computers & security*, 109: 102378.

Li, Z.; Zhang, J.; Liu, L.; and Liu, J. 2022. Auditing Privacy Defenses in Federated Learning via Generative Gradient Leakage. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10132–10142.

Lin, Y.; Han, S.; Mao, H.; Wang, Y.; and Dally, W. J. 2017. Deep gradient compression: Reducing the communication bandwidth for distributed training. *arXiv preprint arXiv:1712.01887*.

Liu, Z.; Luo, P.; Wang, X.; and Tang, X. 2015. Deep Learning Face Attributes in the Wild. In *Proceedings of International Conference on Computer Vision (ICCV)*.

Lopez-Paz, D.; and Ranzato, M. 2017. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30.

McMahan, B.; Moore, E.; Ramage, D.; Hampson, S.; and y Arcas, B. A. 2017a. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, 1273–1282. PMLR.

McMahan, H. B.; Ramage, D.; Talwar, K.; and Zhang, L. 2017b. Learning differentially private recurrent language models. *arXiv preprint arXiv:1710.06963*.

Mo, F.; Borovykh, A.; Malekzadeh, M.; Haddadi, H.; and Demetriou, S. 2021. Quantifying information leakage from gradients. *CoRR, abs/2105.13929*.

Mohassel, P.; and Zhang, Y. 2017. Secureml: A system for scalable privacy-preserving machine learning. In *2017 IEEE symposium on security and privacy (SP)*, 19–38. IEEE.

Nguyen, N.-B.; Chandrasegaran, K.; Abdollahzadeh, M.; and Cheung, N.-M. 2023. Re-thinking Model Inversion Attacks Against Deep Neural Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16384–16393.

Scheliga, D.; Mäder, P.; and Seeland, M. 2022. PRECODE-A Generic Model Extension to Prevent Deep Gradient Leakage. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 1849–1858.

Song, S.; Chaudhuri, K.; and Sarwate, A. D. 2013. Stochastic gradient descent with differentially private updates. In *2013 IEEE global conference on signal and information processing*, 245–248. IEEE.

Sun, J.; Li, A.; Wang, B.; Yang, H.; Li, H.; and Chen, Y. 2021. Soteria: Provable defense against privacy leakage in federated learning from representation perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9311–9319.

Takahashi, H.; Liu, J.; and Liu, Y. 2023. Breaching FedMD: Image Recovery via Paired-Logits Inversion Attack. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12198–12207.

Voigt, P.; and Von dem Bussche, A. 2017. The eu general data protection regulation (gdpr). *A Practical Guide, 1st Ed., Cham: Springer International Publishing*, 10(3152676): 10–5555.

Wang, Z.; Bovik, A. C.; Sheikh, H. R.; and Simoncelli, E. P. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4): 600–612.

Wei, W.; Liu, L.; Loper, M.; Chow, K.-H.; Gursoy, M. E.; Truex, S.; and Wu, Y. 2020. A framework for evaluating client privacy leakages in federated learning. In *European Symposium on Research in Computer Security*, 545–566. Springer.

Wei, W.; Liu, L.; Wut, Y.; Su, G.; and Iyengar, A. 2021. Gradient-leakage resilient federated learning. In *2021 IEEE 41st International Conference on Distributed Computing Systems (ICDCS)*, 797–807. IEEE.

Yin, H.; Mallya, A.; Vahdat, A.; Alvarez, J. M.; Kautz, J.; and Molchanov, P. 2021. See through gradients: Image batch recovery via gradinversion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16337–16346.

Zhang, H.; Cisse, M.; Dauphin, Y. N.; and Lopez-Paz, D. 2017. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*.

Zhang, R.; Isola, P.; Efros, A. A.; Shechtman, E.; and Wang, O. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *CVPR*.

Zhao, B.; Mopuri, K. R.; and Bilen, H. 2020. idlg: Improved deep leakage from gradients. *arXiv preprint arXiv:2001.02610*.

Zhu, J.; and Blaschko, M. 2020. R-gap: Recursive gradient attack on privacy. *arXiv preprint arXiv:2010.07733*.

Zhu, L.; Liu, Z.; and Han, S. 2019. Deep leakage from gradients. *Advances in Neural Information Processing Systems*, 32.

# 8 Appendix

## 8.1 Compute Complexity.

Let's consider a model with $d$ parameters, and we'll denote the time complexity for forward propagation as $h(d)$. As per the Baur-Strassen theorem (Griewank and Walther 2008), the time complexity of a single step in backpropagation will be at most $5h(d)$. Our approach introduces an additional time complexity of $6f(d)$ due to the incorporation of concealed samples. Assuming that $x \in \mathbb{R}^n$ and $f(x) \in \mathbb{R}^m$, we can express the overall computational cost of the objective function in Eq. (9) as $\mathcal{O}(d^2 + n^2 + m^2)$. If the perturbed gradient with concealed samples behaves dissimilar to that of the original gradient, then the computational cost of the gradient projection is around $\mathcal{O}(d^3)$. On TinyImageNet, we observed that one round of updates with our method on an NVIDIA RTX A100 GPU takes approximately 172 seconds. In comparison, an update without any defenses requires around 102 seconds.

## 8.2 Model Architectures

Details of the models used in this study are shown in Tab. 8. The activation layers of the model for the MNIST dataset are Sigmoid, and for CIFAR10 and CelebA, TinyImageNet datasets are ReLU.

| MNIST | CIFAR10/CelebA | TinyImageNet |
|---|---|---|
| $5 \times 5$ Conv, 12 | $5 \times 5$ Conv, 32 | $7 \times 7$ Conv, 64 |
| $5 \times 5$ Conv, 12 | $\{5 \times 5$ Conv, $64\} \times 2$ | $3 \times 3$ MaxPool |
| $5 \times 5$ Conv, 12 | $\{5 \times 5$ Conv, $128\} \times 3$ | $\left\{\begin{array}{l}3 \times 3 \text{ Conv, } 64 \\ 3 \times 3 \text{ Conv, } 64\end{array}\right\} \times 2$ |
| $5 \times 5$ Conv, 12 | $3 \times 3$ MaxPool | $\left\{\begin{array}{l}3 \times 3 \text{ Conv, } 128 \\ 3 \times 3 \text{ Conv, } 128\end{array}\right\} \times 2$ |
| FC-10 | $\{5 \times 5$ Conv, $128\} \times 3$ | $\left\{\begin{array}{l}3 \times 3 \text{ Conv, } 256 \\ 3 \times 3 \text{ Conv, } 256\end{array}\right\} \times 2$ |
| | $3 \times 3$ MaxPool | $\left\{\begin{array}{l}3 \times 3 \text{ Conv, } 512 \\ 3 \times 3 \text{ Conv, } 512\end{array}\right\} \times 2$ |
| | FC-10 (CIFAR10) / FC-2 (CelebA) | $7 \times 7$ AveragePool |
| | | FC-200 |

Table 8: Model architectures for different datasets.

## 8.3 Parameters and Details.

We build on the repository using the official implementation of the GS, GGL, and Imprint attack methods. For the defenses Soteria, Prune, and DP-Gaussian, we build on the repository from the study (Sun et al. 2021). For the defense ATS, we build upon the repository from the study (Balunović et al. 2022). For training, we apply SGD optimizer and set the learning rate for updating the local models $\eta = 0.01$ with exponential decay. The pruning rate from the defense Prune is 0.9 for CIFAR10 and 0.7 for others. The variance of noise distribution from the defense DP-Gaussian is 0.01 for MNIST, 0.5 for TinyImageNet, and 0.001 for others. The pruning rate from the defense Soteria is 0.2 for MNIST, 0.5 for CIFAR10, 0.7 for CelebA, and 0.9 for TinyImageNet. Our method set $\epsilon = 0.01$, $\lambda_g = 0.7$, and the number of iteration as 1000 for all datasets, $\lambda_x = 0.01$ and $\lambda_z = 0.01$ for MNIST and CIFAR10, $\lambda_x = 0.01$ and $\lambda_z = 0.1$ for CelebA, $\lambda_x = 0.001$ and $\lambda_z = 0.01$ for

TinyImageNet. The weights of penalty terms in Eq. (9) are chosen to balance them or kept fixed across all experiments (*e.g.*, $\epsilon = 0.01$ for all experiments.) In our experiments, we opt for the entirety of the training data in the target client to be sensitive and generate one synthetic sample per sensitive datum (ratio 1:1). SSIM is 0.92 with an accuracy of 87% and 0.86 with an accuracy of 86.91% for ratios 10:1 and 5:1, respectively. While the protection becomes better, the model's performance will drop.
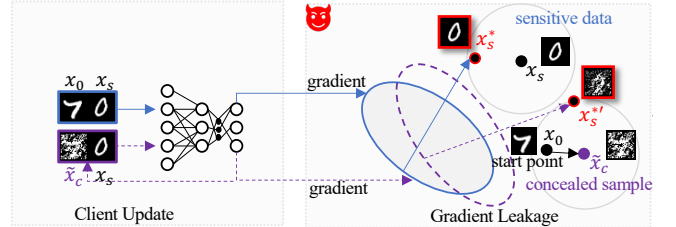


Figure 5: Illustration of gradient leakage and the proposed defense method. The adversary can obtain a perfect reconstruction $x_s^*$ for the sensitive data $x_s$ when the gradient is over the input $x_0$ and $x_s$, but it fails and get the reconstruction $x_s^{*'}$ when the gradient is over our concealed sample $\tilde{x}_c$ and $x_s$.

Consider a sensitive data point $x_s$ (see Fig. 5 for an illustration), we aim to craft the concealed sample $\tilde{x}_c$ which makes $\nabla_\theta \mathcal{L}(f_\theta(\tilde{x}_c), \tilde{y}_c)$ approaching $\nabla_\theta \mathcal{L}(f_\theta(x_s), y_s)$. This strategy obfuscates the sensitive samples, as the reconstruction by the adversary through the gradient will contain information from the concealed sample, which we aim to be visually different from the sensitive data.

## 8.4 Extra Results

| Defense | PSNR↓ | SSIM↓ | Acc↑ (Sensitive Data) | Acc↑ (Test set) |
|---|---|---|---|---|
| ATS | $19.68_{\pm4.60}$ | $0.59_{\pm0.14}$ | $67.99_{\pm0.57}$ | $75.70_{\pm0.02}$ |
| DCS$^2$ (Ours) | $\mathbf{8.04}_{\pm1.10}$ | $\mathbf{0.15}_{\pm0.05}$ | $\mathbf{80.39}_{\pm0.07}$ | $\mathbf{79.79}_{\pm0.03}$ |

Table 9: Defenses against GS attack on CIFAR10. ATS applies data augmentations and here it needs to run for extra 50 rounds to converge.

| Defense | PSNR↓ | SSIM↓ | Acc↑ (Sensitive Data) | Acc↑ (Test set) |
|---|---|---|---|---|
| None | $57.50_{\pm1.95}$ | $1.00_{\pm0.00}$ | $89.84_{\pm0.00}$ | $76.60_{\pm0.00}$ |
| DP-Gaussian | $34.73_{\pm0.79}$ | $0.83_{\pm0.04}$ | $83.59_{\pm0.00}$ | $75.75_{\pm0.01}$ |
| Prune | $14.49_{\pm1.81}$ | $0.39_{\pm0.05}$ | $83.59_{\pm0.00}$ | $75.75_{\pm0.00}$ |
| Soteria | $7.28_{\pm0.60}$ | $0.23_{\pm0.04}$ | $85.42_{\pm0.00}$ | $\mathbf{75.92}_{\pm0.01}$ |
| DCS$^2$ (Ours) | $\mathbf{7.27}_{\pm1.77}$ | $\mathbf{0.16}_{\pm0.06}$ | $\mathbf{89.32}_{\pm0.00}$ | $75.19_{\pm0.02}$ |

Table 10: Against GS attack on MNIST on Non-IID setting. Non-IID data indeed presents additional challenges to training.
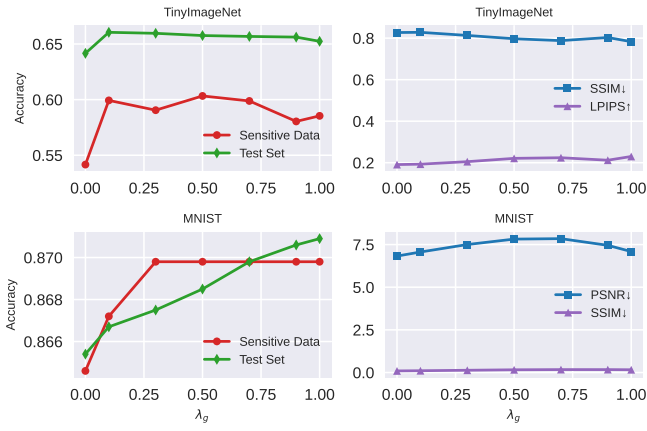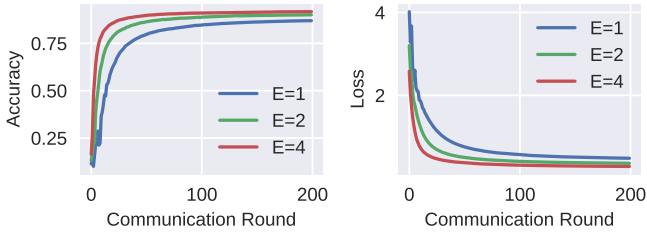
Figure 6: DCS$^2$ with different $\lambda_g$.



Figure 7: FL training process of DCS$^2$ on MNIST with different number of local epochs.

| | None | DP-Gaussian | Prune | Soteria | DCS$^2$ |
|---|---|---|---|---|---|
| PSNR↓ | 59.22±2.71 | 35.28±2.50 | 14.23±2.23 | 9.94±1.10 | **7.87±2.44** |
| SSIM↓ | 1.00±4.77 | 0.82±0.07 | 0.37±0.06 | 0.31±0.07 | **0.18±0.09** |

Table 11: Defend against adaptive attacks on MNIST.



Figure 8: (a) Left to Right: reconstructions by the Imprint attack **w/o** and **w/** DCS$^2$ on FACET. (b) Top to Bottom: GT and reconstructions by the GS attack on MNIST.

We evaluated the Imprint attack with the recently introduced dataset FACET (Gustafson et al. 2023) (Fig. 8 (a)). With our defense applied, the attack failed to perfectly reconstruct the data. The accuracy of the model on the test set **w/o** and **w/** our defense is around 77.21% and 76.62%, respectively.

Note that neither optimization-based attacks nor model modification attacks can precisely separate the gradient for individual data points. In Fig. 8 (b), for instance, there are four '6' within the batch; the gradients w.r.t. these images cannot be fully separated, and the GS attack fails to reconstruct the third image. Besides, similar to any defense mechanism, it is possible to identify the concealed samples. For example, an adversary could potentially modify the model to reconstruct the concealed samples and use a filtering mechanism to identify them. Our model might become vulnerable if the adversary has extensive knowledge about the data and model as exemplified by the GGL attack (use partial training data to learn prior knowledge), or has the ability to modify model architecture as in the Imprint attack. Fig.2 and Fig. 8 (a) provide examples of GGL and Imprint attacks, where the attack managed to reconstruct facial outlines and vague information. It is important to note that model modification methods are generally identifiable and can be countered by vigilant clients.