# WHEN THE CURIOUS ABANDON HONESTY: FEDERATED LEARNING IS NOT PRIVATE

**Franziska Boenisch**[*]
Fraunhofer AISEC
franziska.boenisch@aisec.fraunhofer.de

**Adam Dziedzic**[†]
University of Toronto and Vector Institute
adam.dziedzic@utoronto.ca

**Roei Schuster**[†]
Vector Institute
roei@vectorinstitute.ai

**Ali Shahin Shamsabadi**[†]
Vector Institute and The Alan Turing Institute
a.shahinshamsabadi@turing.ac.uk

**Ilia Shumailov**[†]
Vector Institute
ilia.shumailov@cl.cam.ac.uk

**Nicolas Papernot**
University of Toronto and Vector Institute
nicolas.papernot@utoronto.ca

## ABSTRACT

In federated learning (FL), data does not leave personal devices when they are jointly training a machine learning model. Instead, these devices share gradients, parameters, or other model updates, with a central party (e.g., a company) coordinating the training. Because data never "leaves" personal devices, FL is presented as privacy-preserving. Yet, recently it was shown that this protection is but a thin facade, as even a passive attacker observing gradients can reconstruct data of individual users contributing to the protocol.

In this paper, we argue that prior work *still* largely underestimates the vulnerability of FL. This is because prior efforts exclusively consider passive attackers that are honest-but-curious. Instead, we introduce an active and dishonest attacker acting as the central party, who is able to modify the shared model's weights before users compute model gradients. We call the modified weights *trap weights*. Our active attacker is able to recover user data *perfectly*. Recovery comes with near zero costs: the attack requires no complex optimization objectives. Instead, our attacker exploits inherent data leakage from model gradients and simply amplifies this effect by maliciously altering the weights of the shared model through the trap weights.

These specificities enable our attack to scale to models trained with large mini-batches of data. Where attackers from prior work require hours to recover a single data point, our method needs milliseconds to capture the full mini-batch of data from both fully-connected and convolutional deep neural networks.

Finally, we consider mitigations. We observe that current implementations of differential privacy (DP) in FL are flawed, as they explicitly trust the central party with the crucial task of adding DP noise, and thus provide no protection against a malicious central party. We qualitatively show that, regardless of this flaw, adding enough DP noise to prevent our attack would significantly degrade FL's performance and training time. We consider other defenses and explain why they are similarly inadequate. A significant redesign of FL is required for it to provide any meaningful form of data privacy to users.

## 1 Introduction

With machine learning (ML) being increasingly applied to sensitive data in critical use-cases such as health care [23,41], smart metering [18,47], or the internet of things [29,37], there is a growing need for privacy-preserving training schemes that do not leak sensitive information. Federated learning (FL) is a widely popular approach designed to preserve training-data privacy via a distributed learning protocol [34]. In FL, private user data can be utilized for jointly training an ML model without the data ever leaving the users' device. Instead, the device computes and sends model updates to a central party which aggregates them to produce a shared model. *Assuming* the model updates do not reveal the user data, FL would, thereby, preserve a notion of privacy.

This assumption has been repeatedly contested by prior work. It has been shown how the model updates sent to the central party not only leak training data membership [36] (*i.e.* allow the attacker to tell if a given data point was used in training)

| | |
|---|---|
| Original | |
| Extracted | |
| Original | ive read a few of the reviews and im kinda sad that a lot of the story seems [UNK] ... |
| Extracted | ive read a few of the reviews and im kinda sad that a lot of the story seems [UNK] ... |

Figure 1: Original data and data points extracted from model gradients with our attack. Extraction is **perfect** *i.e.* reconstruction error is zero.

---

[*]Work done while the author was interning at Vector Institute.
[†]Equal contributions, authors ordered alphabetically.

but also properties of the training data [15, 36]. Inspecting model updates allows attackers to even partially *reconstruct* users' training data [17, 48, 49, 51, 52]. Ultimately, FL in its naive implementation offers little to no *guarantees* regarding potential leakage of user data to other users or to the central party.

In this work, we go a step beyond prior attacks by first showing that the gradients sent to the central party include full, memorized training data points. We then proceed to show that a malicious central party can significantly amplify this leakage by simply adversarially setting the model's weights with our *trap weights* method, prior to dispatching them to users. By doing this, we demonstrate that prior work has dramatically underestimated the vulnerability of FL to attacks on privacy, since its threat model—an honest-but-curious attacker who can observe the natural course of FL training but not corrupt its inputs—has not considered the scenario where the central party (*e.g.*, a company) is actively dishonest.

By adversarially initializing the model with our trap weights, the central party can ensure they are able to *perfectly* reconstruct a significant portion of the users' training data, as depicted in Figure 1. We show that this even holds when the gradients are computed over large training data mini-batches, a scenario in which previous attacks usually fail to obtain high-fidelity reconstructions [46]. Since in FL, the central party holds full control over the shared model that is sent out to users, our attack integrates naturally in the FL protocol.

We furthermore show that, while previous reconstruction attacks rely on solving complex optimization problems that require tens of thousands of back-propagation iterations over the model, our attack extracts individual training inputs by simply projecting the appropriate portions of the users' gradients onto the input domain. We thereby reduce the computational time needed to restore individual training data points from several hours to a few milliseconds.

Next, we address and evaluate potential mitigations. For example, a gold standard to reason about privacy is Differential Privacy (DP), which formalizes privacy by bounding the private information an algorithm can possibly reveal about data it analyzes [12]. This contrasts with the privacy notion inherently achieved in FL, namely keeping the original data secret, which is often referred to as *confidentiality*. One could thus hope that employing a differentially-private optimizer, e.g., Differential Private Stochastic Gradient Descent (DPSGD) [4], when training with FL would prevent our data reconstruction attack. Unfortunately, our evaluation shows this is not the case. First, because in popular differentially-private FL implementations, the addition of the calibrated noise that is required to obtain DP guarantees is performed by the central party and not the users themselves. Hence, clearly, for a dishonest central party, no privacy guarantees can be provided to the users. Second, because, as we qualitatively show, DP must add so much noise to mitigate our attack, that it ends up significantly degrading model performance and training time.
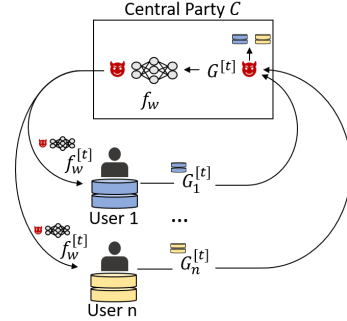


Figure 2: Our attack (😈) targets two points in the FL protocol: At each iteration $t$, the central party actively manipulates the weights $\mathcal{W}$ of the shared model $f_{\mathcal{W}}$ before the model is sent out to the users. This causes the gradients $G_i^{[t]}$ of each user $i$ to contain individual training data points which the central party can then extract before calculating the averaged gradients $G^{[t]}$ and applying them to $f_{\mathcal{W}}$.

We thus conclude that previous robustness assessments of FL with a curious and honest central party have been overly optimistic. In summary, we make the following contributions:

- We observe for the first time that in neural networks starting with a fully-connected layer, even gradients of very large training data mini-batches contain *individual* training data points. In other words, in FL, mini-batch training data points are often directly sent from users to the central party, such that even an honest-but-curious central party has access to them.

- We show that a dishonest and active central party can amplify the leakage of individual training data points and extend it to other model architectures by adversarially initializing the weight of the shared model.

- In this setting, we perform data reconstruction on image and text data. Our attack is able to perform an extremely computationally efficient extraction of individual training data points in only a single-step computation over the received model gradients with the attack setup depicted in Figure 2. For complex image datasets such as ImageNet [10], the attack yields perfect reconstruction of more than $50\%$ of the training data points, even for large training data mini-batches that contain as many as $100$ data points. For textual tasks such as IMDB sentiment analysis [33], it perfectly extracts more than $65\%$ of the data points for mini-batches with $100$ data points.

- We demonstrate that (1) current implementations of FL with DP provide virtually no protection against a non-honest central party because they trust the central party to add DP noise, and (2) enough DP noise to mitigate our attack would result in degraded performance. We then discuss the necessary major redesign of FL to provide meaningful privacy guarantees to users without the need for trust in the central party.

## 2 Background: Neural Networks, Federated Learning, and Differential Privacy

**Neural Networks**. Let $f_{\mathcal{W}} : \mathbb{R}^m \to \{1, \cdots, k\}$ be a $k$-class classifier defined as a set of $l$ layers parameterized by trainable *weights* $\mathcal{W}$. Each layer consists of a linear operation paired with a non-linear activation function (e.g. ReLU). In this work, we consider two popular layer types: fully-connected and convolutional layers.

The goal of the model $f_{\mathcal{W}}$ is to map an input $x_i \in X$ to its desired ground-truth $y_i \in Y$. Therefore, the model weights $\mathcal{W}$ are adapted in a training process, most commonly with the mini-batch *Stochastic Gradient Descent* (SGD). To adjust the initial $\mathcal{W}$, mini-batch SGD repeats the following sequence of steps: (1) sample a mini-batch of size $B$ from the training data $\{(X, Y)_b\}_{b=1}^B$, (2) take a forward pass through the model to obtain its predictions on the mini-batch, (3) compute the difference between predictions and ground-truth labels, called the *loss* $\mathcal{L}$, (4) compute the gradient of $\mathcal{L}$ w.r.t. the weights, called the *weight gradient* $G$, and update the weights accordingly.

To bootstrap mini-batch SGD, weights need to be initialized by sampling from a random distribution; popular distributions include the zero-mean Gaussian [19], *Xavier* [20] or *He* [25] distributions. The choice of distribution has a large effect on learning success [11]. In fact, when weights are maliciously initialized, the final model's utility might be degraded [22].

**Federated Learning**. FL [34] is a communication protocol for training a shared ML model $f_{\mathcal{W}}(\cdot)$ on decentralized data $\{(X_i, Y_i)\}_{i=1}^{\mathsf{N}}$ owned by $\mathsf{N}$ different users $\{u_i\}_{i=1}^{\mathsf{N}}$. Since collecting and managing all the data centrally might be costly, time consuming, and stand in conflict with the confidentiality of these respective users' data, FL enables each user to keep their data locally. A central party coordinates the training of the shared model by iteratively aggregating gradients computed locally by users.

More formally, let $t \in \{1, \cdots T\}$ be the current iteration of the FL protocol. At iteration $t = 0$, the model $f(\cdot)$ is initialized (at random) by the central party denoted as $C$. Let $f_{\mathcal{W}}^{[t]}(\cdot)$ be the model with its weights $\mathcal{W}^{[t]}$ at iteration $t$. At every iteration $t$, M out of the N (M $\ll$ N) users are selected to contribute to the learning. Then, each of the selected M users $u_i$ obtains $f_{\mathcal{W}}^{[t]}(\cdot)$ from $C$ and calculates the gradients $G_i^{[t]}$ for $f_{\mathcal{W}}^{[t]}(\cdot)$ based on one mini-batch $b$ sampled from their local dataset $(X_i, Y_i)_b$. In other words, the user computes the gradient $G_i^{[t]} = \nabla_{\mathcal{W}} \mathcal{L}((X_i, Y_i)_b)$. Each $u_i$ uploads their gradients to $C$, who then averages all of these gradients to update the shared model's parameters:

$$G^{[t]} = \frac{1}{\mathsf{M}} \sum_{i=1}^{\mathsf{M}} G_i^{[t]}, \quad \mathcal{W}^{[t+1]} = \mathcal{W}^{[t]} - \eta G^{[t]}. \quad (1)$$

FL, thereby, represents a decentralization of the mini-batch SGD (i.e. distributed training from mini-batches of user data).

| Method | Attacker | | Rep. | | | Label-Free | B | | Opt./Train.-Free |
|---|---|---|---|---|---|---|---|---|---|
| | U | S | C | U | ID | | 1 | ≥1 | |
| | | | | *stronger* → | | | | *stronger* → | |
| DMU-GAN [26] | ✓ | | ✓ | | | | | ✓ | |
| mGAN-AI [48] | | ✓ | | ✓ | | | | ✓ | |
| DLG [52] | | ✓ | | | ✓ | ✓ | ✓ | | |
| iDLG [51] | | ✓ | | | ✓ | ✓ | ✓ | | |
| GradInv [49] | | ✓ | | | ✓ | ✓ | | ✓ | |
| trap weights [Ours] | | ✓ | | | ✓ | ✓ | | ✓ | ✓ |

Table 1: Comparison of data reconstruction attacks. KEY– U: User, S: Server, C: Class, ID: Individual Data Points, B: Mini-batch size, Opt.: Optimization, Train.: Training, Rep.: Representative.

**Differential Privacy**. User gradients may leak information about their training data [17, 48, 49, 51, 52]. To bound this potential leakage of private information, a gold standard for reasoning about privacy guarantees is the framework of differential privacy (DP) [12]. To implement FL with DP guarantees, mini-batch SGD can be replaced with DPSGD [4]. DPSGD bounds each data point's impact on the learning process by computing per data point gradients, clipping them individually, and adding Gaussian noise to the gradient average for the mini-batch before it is applied to update the model. To perform a model update with DP, the central party in FL takes the clipped gradients received from users and adds noise to the gradients before it aggregates and applies them to the shared model. We refer the reader to Appendix A for the formal definition of DP and the DPSGD algorithm.

## 3 Existing Data Reconstruction Attacks

This section introduces prior work on data reconstruction attacks and discusses the limitations of attacks that are based on iterative optimization.

### 3.1 Prior Work

Phong *et al.* [40] were the first to show how gradients leak information that can be used to recover training data from single neurons or linear layers. Recent work [17, 40, 48, 49, 51, 52] proposed that the central party or users involved in FL training launch data reconstruction attacks based on either training a Generative Adversarial Network [21] (GAN) or solving a second order optimization problem. Table 1 summarizes these data reconstruction attacks (described below) and compares them with our attack.

**Class-wise Representation Reconstruction Attacks**. Hitaj *et al.* [26] were the first to propose a GAN-based data reconstruction attack, called DMU-GAN. The attacker must know the dataset's classes, and the reconstructed data points are generic representations of class-wise properties rather than individual user data points or classes. Wang *et al.* [48] suggested mGAN-AI, which extends DMU-GAN's reconstruction attack to per-user class-wise representations, but still does not extract individual data points. Additionally, both methods require access to data from the same distribution as the users' data. [44] observes that class-wise representations are embedded in model

updates even without the need to reconstruct them using a GAN, and suggest defenses.

**Optimization-based Instance Reconstruction Attacks**. Several attacks aim to reconstruct individual user data points while also relaxing the assumption that data labels are available to the attacker. Zhu *et al.* [52] proposed Deep Leakage from Gradients (DLG), where a data reconstruction attack is formulated as a joint optimization problem on the labels and input data. iDLG [51] sped up the convergence rate of DLG [52] by analytically computing the labels based on the users' gradients of the last layer. These works, and other optimization-based ones [17], are limited to a setting where mini-batches only contain a single example, i.e., $B = 1$. GradInversion [49] regularizes DLG's objective to improve the extraction fidelity, attaining some success in extraction for mini-batches of size $B > 1$. In Section 8.5 we compare performance of our approach against a state-of-the-art optimization-based attack. Our attack is superior in extracting individual training data even for large mini-batch sizes of $B \geq 100$, and being far more computationally efficient (even for passive adversaries in the honest-but-curious model). **Data Extraction Attacks**. In concurrent work, Fowl *et al.* [13] consider a threat model with an active and dishonest central party, similar to our setup. Unlike our trap weights, their active attack requires a change of the model architecture: a fully-connected model layer needs to be inserted early in the architecture. Since this layer's weights have to contain many weight rows with the exact same weight values, this layer is inherently detectable. [1] Additionally, they do not discuss passive analytical-extraction attacks. Finally, our work generalizes their setup and performs successful extraction at arbitrary model layers, also for textual data.

## 3.2 Limitations of Optimization-Based Attacks

We hereby provide a brief exposition to Zhu *et al.* [52]'s DLG, as a representative case study of an optimization-based attack. Their approach, characteristic of optimization-based data reconstruction attacks, is given in Algorithm 2 in the Appendix. It firstly randomly initializes a "dummy data point and corresponding label" $(\hat{x}, \hat{y})$ and computes the resulting "dummy gradients" as $\hat{G} = \nabla_{\mathcal{W}^t} \mathcal{L}(f_{\mathcal{W}^t}(\hat{x}), \hat{y})$. Then, they iteratively optimize the dummy data to produce gradients that are close to the original gradients $G_i^t$ by solving:

$$\mathbf{x}_i^* = \arg\min_{\hat{\mathbf{x}}} \|G_i^{[t]} - \hat{G}\|^2, \quad y_i^* = \arg\min_{\hat{y}} \|G_i^{[t]} - \hat{G}\|^2. \quad (2)$$

DLG often fails to reconstruct high-fidelity data points and discover the ground-truth labels consistently because of a lack of convergence in the optimization. While other methods offer improvements (e.g. iDLG [51] sped up the convergence by

---

[1]This is inherent to the attack because the method relies on each row computing the exact same function on the data and binning its result by varying only the bias term, such that it becomes likely that a bin contains only one input. Conversely, our trap weights are initialized such that it is likely that an output neuron is only activated for a single input in a mini-batch while avoiding imposing a highly regular structure on the weight matrix.

simplifying the objective in Equation 2 from both data and label reconstruction to only data reconstruction; and GradInversion [49] adds useful regularization), they suffer from the same pathology.

We identify several reasons for this. First, the gradient of the loss is non-injective *i.e.* is not invertible everywhere: different mini-batches may yield nearly identical gradients [43]. This holds whether the user samples mini-batches that contain multiple data points or a single data point only, *i.e.* $B = 1$. Second, optimization-based attacks converge to different minima due to the underlying randomness (see step 1 in Algorithm 2 in the Appendix). These minima correspond to different possible reconstructions of the input that often differ from the original training points [49]. Third, optimization-based attacks are computationally expensive: they either need to train a GAN or solve a second-order gradient optimization problem. Instead, our attack extracts *exact* data points from the gradients without any optimization or GAN training.

## 4 Attacker Model

Our attacker's goal is to extract as many individual training data points from an attacker-chosen subset of the participating users. This attacker's primary vantage point is the central party: it controls and can alter the weights of the shared model that is sent to users, and it can read users' gradient updates that are sent back, as Figure 2 illustrates.

Following prior work, we consider an FL protocol where users calculate the model gradients locally on one (potentially large) mini-batch of their training data and share the resulting gradients directly with the central party. We note that when users have abundant amounts of data, they might perform local gradient calculation and averaging over more than one mini-batch, and we evaluate this scenario in Section 8.3.

Our approach can target various diverse architectures. In this work we focus on attacking neural networks that comprise fully-connected and convolutional layers, and show passive and active attacks on common architectures of these types. In our proposed scheme, attacked models must have at least one fully-connected layer. The bigger the mini-batch size, the higher this layer's dimensions need to be if we want the attack to have consistent success rates. Table 6 in Section 8.3 quantifies this trade-off for FC-NNs. The dimension of the other network layers need to be sufficiently large to carry the input features to the fully-connected layer in a forward pass, which occurs naturally in common architectures (see Appendix B and Section 7).

## 5 Passive Analytical Extraction for FC-NNs

Here, we show how the gradients of an FC-NN directly leak individual training points they are computed on, even to a passive attacker who just observes said gradients. In Section 5.1, we mathematically prove that for a single training data point, *i.e.* a mini-batch size of $B = 1$, perfect extraction from the network gradients is possible. Then, in Section 5.2, we moti-

vate why it is also possible to perfectly extract a small number of individual data points from gradients, even when working with larger mini-batches of size $B > 1$. However, the success of this passive extraction attack drops as the mini-batch sizes increase. This limitation motivates our active adversarial weight initialization attack, which we introduce in Section 6.

## 5.1 Single-Input Gradients Directly Leak Input

It has been shown by Geiping *et al.* [17] that a single input data point $\mathbf{x}$ can be reconstructed from the gradients of any fully-connected layer which is preceded only by fully-connected layers and contains a bias $\mathbf{b}$. This holds if the gradient of the loss w.r.t. the layer's output $\mathbf{y} = \mathrm{ReLU}(W\mathbf{x} + \mathbf{b}) = \max(0, W\mathbf{x} + \mathbf{b})$ contains at least one non-zero entry. For detailed proof of the above see Proposition D.1 in [17]. In particular, when considering the first model layer, reconstructing its input data *directly* corresponds to obtaining the original input data point $\mathbf{x}$. Let $y_i$ denote the output of the $i^{th}$ neuron of the first and fully-connected layer of a model, and let $\mathbf{w}_i^T$ be the corresponding row in the weight matrix and $b_i$ the corresponding component in the bias vector. Assume $\mathbf{w}_i^T\mathbf{x} + b_i > 0$, and therefore, $\mathrm{ReLU}(\mathbf{w}_i^T\mathbf{x} + b_i) = \mathbf{w}_i^T\mathbf{x} + b_i$. The reconstruction of the input $\mathbf{x}$ is done by calculating the gradients of the loss w.r.t. the bias and the weighs as follows:

$$\frac{\partial \mathcal{L}}{\partial b_i} = \frac{\partial \mathcal{L}}{\partial y_i} \frac{\partial y_i}{\partial b_i} = \frac{\partial \mathcal{L}}{\partial y_i} \tag{3}$$

since $\frac{\partial y_i}{\partial b_i} = 1$, where $y_i = \mathbf{w}_i^T\mathbf{x} + b_i$.

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}_i^T} = \frac{\partial \mathcal{L}}{\partial y_i} \frac{\partial y_i}{\partial \mathbf{w}_i^T} = \frac{\partial \mathcal{L}}{\partial b_i} \mathbf{x}^T \tag{4}$$

Thus, if any $\frac{\partial \mathcal{L}}{\partial b_i} \neq 0$, perfect reconstruction is given by:

$$\mathbf{x}^T = \left(\frac{\partial \mathcal{L}}{\partial b_i}\right)^{-1} \frac{\partial \mathcal{L}}{\partial \mathbf{w}_i^T} \tag{5}$$

According to Equation (4), the gradient of the loss w.r.t. the weights directly contains a scaled version of the input data. The exact scaling factor is $\left(\frac{\partial \mathcal{L}}{\partial b_i}\right)$, which is the gradient of the loss w.r.t. the bias. This gradient is computed in the regular backward pass together with the gradient of the weights. Therefore, obtaining the scaling factor by just reading it from the gradients of the bias and inverting it to $\left(\frac{\partial \mathcal{L}}{\partial b_i}\right)^{-1}$ comes at zero costs and the factor can be directly applied to rescale the gradient of the weights and obtain the input data point $\mathbf{x}$, see Equation (5). Intuitively, the reason why there is a rescaled version of the input data in the gradients and why this would be beneficial for learning can be motivated by revisiting the simple perceptron algorithm [14]. When an input is misclassified, the weight update in the perceptron algorithm consists simply in adding this input to the weights, which makes the algorithm learn.

## 5.2 Mini-batch Gradients Directly Leak Some Individual Inputs

It turns out that individual data point leakage is not limited to gradients computed over a mini-batch of size $B = 1$: we observe that gradients computed over larger mini-batches also sometimes leak individual training points. To forge an intuition for this phenomenon, Figure 7 in Appendix D visualizes the gradients of the first fully-connected layer's weight matrix of the FC-NN described in Table 7. We see that we are able to clearly distinguish some of the training data points within the rescaled gradients. This is despite the fact that these gradients were computed over a mini-batch of $B = 100$ inputs sampled from the CIFAR10 dataset.

**Why do some gradients contain individual training data points?** We denote a training data mini-batch by $X = \{x_1, x_2, \cdots, x_B\} \in \mathbb{R}^{(m \times B)}$ with $B > 1$. The gradient of this mini-batch $X$ is equal to the average of all gradients computed for each of the data points $\{x_1, x_2, \cdots, x_B\}$ that make up the mini-batch. Let $y_i$ denote again the output of the $i^{th}$ neuron of the fully-connected layer, and let $\mathbf{w_i}$ and $b_i$ be the corresponding row in the weight matrix and the component in the bias vector, respectively. Then the gradient $\mathbf{G}_{\mathbf{w}_i^T}$ and $G_{b_i}$ of $\mathbf{w_i}$ and $b_i$ can be computed as follows:

$$\mathbf{G}_{\mathbf{w}_i^T} = \frac{1}{B} \sum_{j=1}^{B} \frac{\partial \mathcal{L}}{\partial y_{(i,j)}} \frac{\partial y_{(i,j)}}{\partial \mathbf{w}_i^T}$$

$$G_{b_i} = \frac{1}{B} \sum_{j=1}^{B} \frac{\partial \mathcal{L}}{\partial y_{(i,j)}} \frac{\partial y_{(i,j)}}{\partial b_i} \tag{6}$$

with $y_{(i,j)} = \mathrm{ReLU}(\mathbf{w}_i^T\mathbf{x}_j + b_i)$. These equations illustrate that the gradient $\mathbf{G}_{\mathbf{w}_i^T}$ over the data mini-batch $X$ contains a weighted overlay of all the input data points $\mathbf{x}_j$ from the mini-batch. The weighting, therein, depends on the contribution of each data point to the model loss $\mathcal{L}$.

We observe that, in some cases, all but one training data point $\mathbf{x}^*$ from the data mini-batch have zero gradients. This is due to the $\max$ operation in $\mathrm{ReLU}(\mathbf{w}_i^T\mathbf{x} + b_i) := \max(\mathbf{w}_i^T\mathbf{x} + b_i, 0)$. When $\mathbf{w}_i^T\mathbf{x} + b_i$ is negative, the ReLU outputs zero, which results in zero gradients for the corresponding data point. When the gradients are zero for all data points but for the one data point $\mathbf{x}^*$, the weight gradient $\mathbf{G}_{\mathbf{w}_i^T}$ from Equation (6) becomes $\mathbf{G}_{\mathbf{w}_i^T} = \frac{1}{B} \frac{\partial \mathcal{L}}{\partial y_{(i,*)}} \frac{\partial y_{(i,*)}}{\partial \mathbf{w_i^T}}$ with $y_{(i,*)} = \mathrm{ReLU}(\mathbf{w}_i^T\mathbf{x}^* + b_i)$. This reduces the data extraction from the case of $B > 1$ to the case of $B = 1$, for which we saw in Section 5.1 that the data point $\mathbf{x}^*$ can be perfectly extracted. In other words, $\mathbf{w}_i^T\mathbf{x} + b_i$ being negative for all data points but one results in accidental leakage of that data point—enabling its exact reconstruction by a passive adversary.

5

# 6 Active Adversarial Initialization of the First Fully-Connected Layer

Section 5 illustrates under which conditions model gradients leak data points to a passive attacker capable of observing these gradients. This possibility in turn implies that no privacy guarantees for the user's data can be made in the FL setting—in particular not when the central party is dishonest—because the user would share gradients with the central party.

In the following, we show how an active attacker can amplify previously-accidental leakage during the passive attack by controlling the weights $\mathbf{w}_i$ and biases $b_i$. For example, while a passive attacker can extract roughly 20% of arbitrary data points from a batch size $B = 100$ for 1000 neurons (*i.e.* weight rows in the fully-connected layer) on ImageNet, the active attack can more than double the number of extracted data points to 45%. Next, we show how to make such malicious choices to extract a larger number of individual training data points from model gradients.

## 6.1 If-Else Logics over Relationships Between Data Features

Without loss of generality, we will suppress the bias term in the following considerations. The multiplication of a single weight row $\mathbf{w}_i$ corresponding to the $i^{th}$ neuron at the fully-connected layer with some input data point $\mathbf{x}$ can be expressed as a weighted sum of all of the features in $\mathbf{x}$ as follows

$$\mathbf{y}_i = \mathbf{w}_i^T \mathbf{x} = \sum_{j=1}^{m} w_i^{(j)} x_j. \tag{7}$$

In $\mathbf{w}_i$, let N and P denote the sets of indices that hold the negative and positive weight components, respectively. Given ReLU activation, the $i^{th}$ neuron is only activated on $\mathbf{x}$ if the sum of the features weighted by the negative components is smaller than the sum of the features weighted by the positive components:

$$\sum_{n \in \mathbb{N}} w_i^{(n)} x_n < \sum_{p \in \mathbb{P}} w_i^{(p)} x_p. \tag{8}$$

Therefore, $\mathbf{x}$ will yield non-zero gradients at the $i^{th}$ neuron if and only if this particular relationship, expressed by Equation (8), between the features holds. When the inequality holds only for a single data point in a mini-batch, this data point can be individually extracted from the gradients, as described in Section 5.2. The idea behind our trap weights is to set the components within each weight row corresponding to each neuron of the first fully-connected layer, such that this relationship only holds relatively rarely in inputs, and is therefore likely to only hold for a single data point within a mini-batch.

## 6.2 Adversarial Weight Initialization

Intuitively, our approach adversarially initializes each row of the weight matrix to increase the likelihood that only one data point in a given mini-batch will activate the neuron corresponding to that row. To achieve this, we initialize a randomly

---

**Algorithm 1:** Adversarial Initialization of a Weight Row.

**Input:** Weight row $\mathbf{w_i}$ of length $L$, Gaussian distribution $\mathcal{N}(\mu, \sigma)$ with mean $\mu$ and std $\sigma$, Scaling factor $s < 1$, Discrete uniform distribution $\mathcal{U}(\cdot, \cdot)$

**Output:** Adversarially initialised weight row $\mathbf{w_i}$

1: $\mathbb{N} = \{i | i \sim \mathcal{U}(1, L)\}$ where $|\mathbb{N}| = \frac{1}{2}L$ ▷ Select randomly indices for negative weights
2: $\mathbb{P} \leftarrow \{i \notin \mathbb{N} | i \in [L]\}$ ▷ Select indices for positive weights
3: $\mathbf{z}_- \sim \mathcal{N}(\mu, \sigma) | \mathbf{z}_- \in \mathbb{R}_-^{\frac{1}{2}L}$ ▷ Negative samples
4: $\mathbf{z}_+ = -s \cdot \mathbf{z}_-$ ▷ Positive samples
5: $\mathbf{w_i}[\mathbb{N}] \leftarrow \text{Shuffle}(\mathbf{z}_-)$ ▷ Initialize negative weights
6: $\mathbf{w_i}[\mathbb{P}] \leftarrow \text{Shuffle}(\mathbf{z}_+)$ ▷ Initialize positive weights

---

chosen half of the components of the weight row to negative values, and the other half to the corresponding positive values, by sampling from a Gaussian normal distribution. Negative weights are sampled with slightly larger absolute values than the positive weights to decrease the likelihood of multiple training points activating the corresponding neuron. See Algorithm 1 for a formalization of our initialization.

We use a relatively large standard deviation, such as $\sigma = 0.5$, so sampled points span widely across the weight range. This causes some features to be multiplied with relatively large positive or large negative values, and thereby increases the variability of the sums in Equation (8) across data points, ensuring that the inequality holds only for a few of them—in the best case only one—which allows for perfect extraction.

We use the scaling factor $s$ to specify how much larger the absolute values of the negative weight components should be than the positive values. This determines how "aggressively" our activation causes weighted inputs to individual neurons to be negative, thereby to be filtered out by the ReLU function and to have zero gradients for most input data points. The ideal value of $s$ when it comes to attack effectiveness is dataset dependent; it could be fine-tuned by an attacker through access to one mini-batch of data, or over time by evaluating the attack success on the gradients received from the users. However, to set $s$, the attacker does not rely on knowledge about the actual values of the input data features and how they relate. In fact, the only assumptions we need to make is that our attacker knows the *dimensionality* of the data, *i.e.* the number of data features, and that features are scaled in the range $[0, 1]$, which is a standard pre-processing step.

Our adversarial initialization causes the ReLU function for many neurons at the fully-connected layer to activate only for one input data point per mini-batch. Due to the randomness in the initialization of each weight row corresponding to a neuron, different neurons are likely to be activated by different input data points. Thereby, the gradients of different weight rows allow for the extraction of different individual data points. However, there are also cases where the same data point can activate different neurons. We demonstrate the success of this method in Section 8.3 by showing that our trap weights increase the proportion of neurons that only activate on one random individual data point in a mini-batch by more

than factor 10, and thus we are able to extract more than double the number of individual training data points. *E.g.* our trap weights cause $51.4\%$ of active neurons out of 1000 to by activated by individual data points from the ImageNet dataset while random model weights with a Gaussian normal initialization with $\sigma = 0.5$ only yield $4.4\%$. This allows for an individual extraction of $45.7\%$ of the data points in a mini-batch of size $B = 100$ for our trap weights versus $21.8\%$ with random model weights.

We expect that an attacker with more precise background knowledge on the distribution of the training data will be able to craft even more targeted weight initializations for more effective individual data extraction. Such background knowledge may include the distribution of feature values and their relation to each other for the data to be extracted.

# 7 Generalization of Data Extraction Attack to Convolutional Neural Networks

So far, both the passive and active attacks we described are tailored to extracting data from the gradients computed to update a fully-connected layer. However, modern neural network architectures often rely on convolutional layers to model image and text data alike. It is difficult to directly apply our attack strategy to these convolutional layers because they rely on the weight sharing principle: to decrease the effective number of parameters that need to be trained, the same weight values are applied to multiple locations of the image to extract patterns regardless of their location in the image. In this section, we thus propose a second instantiation of our adversarial weight initialization strategy that generalizes extraction attacks to convolutional neural networks (CNNs).

Our solution reduces networks with convolutional layers to the setting we previously considered with fully-connected neural networks. To do so, we observe that a CNN typically composes a few convolutional layers with fully-connected layers. We thus initialize the weights of the convolutional layers such that they transmit the model input *unaltered* up to the fully-connected layers of the model architecture.[2] We then extract the input from gradients computed for the first fully-connected layer that follows convolutional layers, using the attack strategy described in Section 6.

There are two important requirements for our approach to transmitting, or forwarding, model inputs through convolutional layers. The first is to make sure that no feature of the input data is lost. This requires having at least as many parameters at every convolutional layer as the number of input features. The second is to make sure that different features do not get overlaid. We explain how to ensure this next.

---

[2]The same idea can be instantiated to transmit unaltered inputs over fully-connected model layers. This allows for input extraction at other model layers than the first one—which is the layer we worked with in our attack thus far. See Appendix B for details on input transmission over fully-connected layers.

**Two Dimensional Input.** In general, preserving input size over a convolutional layer can be achieved through an adequate combination of padding, stride, and filter sizes. Specifically, we use stride *one* and an adequate zero-padding to preserve the size of the layer input. In order to transmit the input features, we create a filter with uneven dimensions $(w, h)$, where $w = h$, and we initialize it with *zero* everywhere apart from the element in the middle which we set to *one*. For a two dimensional input (*e.g.* a grey-scale image), the described filter perfectly transmits the information to the next layer and creates a feature map that exactly replicates the input. See Figure 3a for this adversarially initialized filter.

**Three Dimensional Input.** Some input data to CNNs is distributed over several input channels, such as color images, that consist of three channels. At every layer, we, therefore need three adversarially initialized convolutional filters to "transmit a copy" of the input channels. A standard architecture can have many more filters per layer, which can, in the case of our attack be randomly initialized since they will be ignored by the attacker. Assume now that the original input features at the current layer $l_i$ are distributed over $a_{l_i}$ of the total $b_{l_{i-1}}$ many feature maps. For example, in the first model layer, $a_{l_i} = b_{l_{i-1}}$ corresponds to the number of color channels required to encode the image. In subsequent layers, the remaining $b_{l_{i-1}} - a_{l_i}$ many feature maps contain random noise, introduced by random filters that do not transmit the input features (*e.g.* Filter 1 in Figure 3b). We denote the indices of the feature maps where the input features are located by $\vec{\alpha_{l_i}}$. We then need $a_{l_i}$ many filters, initialized as described above to transmit the information to the next layer. The filters differ from each other only by the placement of the matrix that contains the *one* element. This placement must correspond to different indices in $\vec{\alpha_{l_i}}$. See Figure 3b for a visualization of this setting. Note that the placement of the feature-transmitting filters at layer $l_i$ will determine the indices $\vec{\alpha_{l_{i+1}}}$ of the feature maps that are input to the next layer.

In the last convolutional layer before the fully-connected layer that we want to transmit the input to, the filters containing noise should be initialized such that they yield negative input to the ReLU function. Thereby, the output of the last convolutional layer becomes zero everywhere apart from the feature maps produced by the filters transmitting input data features. The flattened output then serves as input to the fully-connected layer, and reconstruction can be conducted as described in Section 6.

# 8 Experimental Evaluation

In this section, we validate that our adversarial weight initialization attack allows a central party to reconstruct individual training data points from gradients shared by users. We use three different image datasets, namely MNIST [32], CIFAR10 [31], and ImageNet [10] and the text-based IMDB [33] dataset for sentiment analysis. Because our approach is applicable to FC-NNs and CNNs, we test it against both of these architectures. We instantiate our attack against an FC-NN for the MNIST dataset, and against a CNN for CIFAR10 and Im-

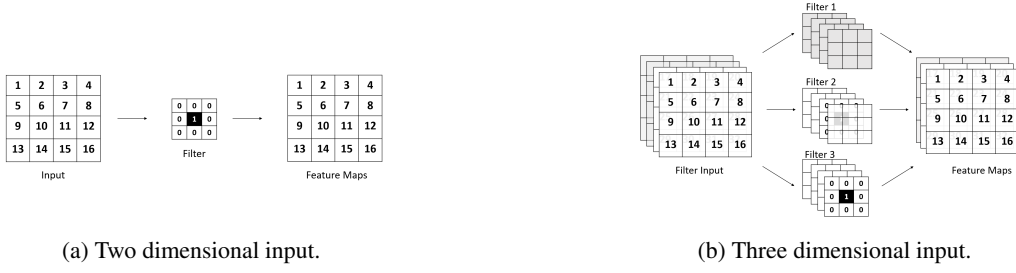(a) Two dimensional input.

(b) Three dimensional input.

Figure 3: Adversarially initialized convolutional filters that transmit their input to the next layer. The numbers indicated in the input and feature map represent the features, the numbers in the filter represent the weight initialization. Grey layers indicate random weight values while white layers indicate zero weights.

ageNet. For the IMDB dataset, we use a model whose input is 250-token sentences, and consists of an *embedding layer*, which maps each token in a 10,000-word vocabulary to a 250-dimensional floating-point vectors, and inputs these to a fully-connected layer. The specifics of our model architectures for image and text data are described in Table 7 and Table 8 in the Appendix, respectively. We implemented our adversarial initialization, trap weights, and the experiments in TensorFlow [5] version 2.4. The code will be open sourced after the peer review process.

**Attack Instantiation.** Because the central party has access to the gradients of *all* model layers uploaded by users, it is able to choose which layer to instantiate the attack on.

For the FC-NN, we adversarially initialize the first layer with our trap weights and extract training data points from its gradients. Through our initialization that allows for transmission of input data over fully-connected layers (see Appendix B), we could also forward the input data to a later layer and perform extraction there.[3]

In the CNNs, we first initialize the convolutional layers to transmit the input data to the first fully-connected layer of the architecture. Then, we adversarially initialize this layer's weights with our trap weights for extraction.

For the text classifier, we initialize the weights of the embedding layer with a random uniform distribution (min=0., max=1.) to create the inputs for the fully-connected layer. We then adversarially initialize this fully-connected-layer's weights with our trap weights to perform extraction of the embeddings there. Finally, after extracting the input's embeddings from then fully-connected layer (like they would extract an image), we map them back to the corresponding text. To reconstruct the original text tokens from a sequence of extracted embeddings, the attacker creates a lookup dictionary, mapping its initialized embeddings back to their corresponding tokens (this is the inverse mapping to the embedding layer). To avoid vector-comparisons for each lookup, we use hash values for vector embeddings as keys.

---

[3]Such a forwarding to a later fully-connected layer might be useful when this layer consists of more neurons than the first fully-connected layer: The experimental results in Table 6 show that the more neurons, the more data points can be extracted individually.

## 8.1 Extraction Success Metrics

We introduce two novel metrics to measure the success of individual data point extraction from model gradients.

**Extraction-Precision.** The first metric, which we call *extraction-precision*, captures the percentage of gradient rows at the given layer's weight matrix from which we can extract any input data point individually. This metric enables us to quantify how well the adversarial weight initialization manages to generate weights that cause activation for exactly one single data point. *Extraction-Precision* can be calculated as follows:

$$\text{P} = \frac{G_0}{N}, \tag{9}$$

with $N$ being the total number of neurons, and $G_0$ denoting the number of gradient rows from which we can extract a data point individually and with an $\ell_2$-distance of zero to any of the input data points. Note that the total number of gradient rows $G$ is equivalent to the number of neurons $N$ at that layer.

However, the *extraction-precision* metric alone would not be expressive enough since a high *extraction-precision* could be achieved despite the exact same individual training input being reconstructed from all gradient rows. Therefore, we defined a second metric that we call *extraction-recall*.

**Extraction-Recall.** The *extraction-recall* measures the percentage of input data points that can be perfectly extracted from any gradient row. We define it by

$$\text{R} = \frac{B_0}{B}, \tag{10}$$

where $B$ is the number of data points in the given mini-batch and $B_0$ is the number of these data points that we can extract with an $\ell_2$-error of zero from the rescaled gradients.

## 8.2 Evaluating the Passive Attack

Recall from Section 5 that extraction of training data from gradients is possible even when model weights are initialized randomly. We evaluate this passive attack to obtain a baseline for our adversarial weight initialization strategies. To evaluate the passive attack, we measure the extraction success of individual training data points from the gradients of randomly initialized models.

| Weights Initializer | MNIST | | CIFAR10 | | ImageNet | |
|---|---|---|---|---|---|---|
| | P | R | P | R | P | R |
| Xavier Normal | .004 | .037 | .048 | .203 | .046 | .213 |
| Xavier Uniform | .005 | .048 | .053 | .229 | .040 | .201 |
| Gaussian ($\sigma$=0.01) | .005 | .048 | .051 | .226 | .041 | .203 |
| Gaussian ($\sigma$=0.1) | .005 | .049 | .053 | .238 | .043 | .209 |
| Gaussian ($\sigma$=0.5) | .006 | .050 | .058 | .255 | .044 | **.218** |
| Gaussian ($\sigma$=1) | .006 | .059 | .058 | .256 | .045 | **.218** |
| Gaussian ($\sigma$=2) | .007 | **.061** | .058 | **.259** | .047 | .217 |

Table 2: Impact of random initialization functions on the *extraction-precision* (P) and *extraction-recall* (R) of individual training data points from the model gradients. The displayed numbers refer to a mini-batch of 100 data points and 1000 neurons for extraction in the first model layer (see model architecture in Table 8). All numbers are averaged over 10 runs with different random initializations.

| Epoch | MNIST | | | | CIFAR10 | | | |
|---|---|---|---|---|---|---|---|---|
| | Loss $\mathcal{L}$ | A | P | R | Loss $\mathcal{L}$ | A | P | R |
| 0 | .526 | .998 | .005 | .050 | 1.857 | .907 | .053 | .232 |
| 5 | .067 | .997 | .044 | .137 | 1.352 | .900 | .044 | .195 |
| 10 | .021 | .997 | .116 | .154 | 1.088 | .913 | .041 | .196 |
| 15 | .006 | .997 | .131 | .165 | .768 | .923 | .043 | .206 |
| 20 | .002 | .997 | .136 | .167 | .472 | .931 | .050 | .232 |
| 25 | .001 | .997 | .140 | **.169** | .282 | .935 | .058 | .241 |
| 30 | .001 | .997 | .142 | .168 | .200 | .936 | .062 | **.267** |

Table 4: Success of passive data extraction based on the training stage of the corresponding models. The results depict the percentage of *active neurons* (A), *extraction-precision* (P), and *extraction-recall* (R) for extraction with a mini-batch size of 100 data points from the first layer of the fully-connected network from Table 7. All numbers are averaged over 10 runs with different random initializations.

| | Passive Attack | | | Our Active Attack | | |
|---|---|---|---|---|---|---|
| **B** | A | P | R | A | P | R |
| 20 | .842 | .072 | **.900** | .519 | .610 | **1.000** |
| 50 | .885 | .050 | **.552** | .776 | .376 | **.962** |
| 100 | .909 | .036 | **.254** | .910 | .192 | **.654** |
| 200 | .927 | .030 | **.128** | .978 | .070 | **.255** |

Table 3: Extraction success for the IMDB dataset depends on the size **B** of the mini-batches for passive attack and active attack with adversarial initialization. The results depict the percentage of *active neurons* (A), *extraction-precision* (P), and *extraction-recall* (R). All numbers are averaged over 10 runs with different random and adversarial initialization of the model from Table 8, respectively.

Table 2 reports the *extraction-precision* and *extraction-recall* of training data point extraction from the gradients of randomly initialized models. These gradients are computed over a mini-batch of 100 data points for 1000 neurons (*i.e.* 1000 weight rows' gradients for extraction) in the first fully-connected layer. We later study the impact of these two parameters on the success of reconstruction attacks. Even if this attack is passive, and the central party has not modified any of the weights adversarially, training data extraction is often successful: for the MNIST dataset, around 6% of individual training data points can be directly extracted from the model gradients, whereas for CIFAR10 and ImageNet, roughly 26% and 22% of the training data points can be perfectly extracted. The passive attack for extracting embeddings from the IMDB dataset yields roughly 25% *extraction-recall* for 1000 neurons and mini-batches of 100 data points, see Table 3.

These results also suggest that setting higher standard deviations in random weight distributions alone can already significantly increase the *extraction-recall* of individual data points from the model gradients, see Table 2. This is, most likely, due to the larger span within the weight values.

Additionally, we also set out to investigate how as training progresses, and the model's weights converge, the extraction' success evolves. We initialized the FC-NN from Table 7 with a Xavier Uniform distribution and trained the model on MNIST and CIFAR10 for 30 epochs. Table 4 depicts the results. We observe that the *extraction-recall* increases slightly over the training epochs. Analyzing the distribution of the model

| s | MNIST | | | CIFAR10 | | | ImageNet | | |
|---|---|---|---|---|---|---|---|---|---|
| | A | P | R | A | P | R | A | P | R |
| .400 | .022 | .803 | .114 | 0. | 0. | 0. | .0. | 0. | 0. |
| .500 | .149 | .636 | .354 | 0. | 0. | 0. | 0. | 0. | 0. |
| .600 | .462 | .408 | .526 | 0. | 0. | 0. | 0. | 0. | 0. |
| .700 | .796 | .203 | **.540** | 0. | 0. | 0. | 0. | 0. | 0. |
| .800 | .959 | .062 | .334 | 0. | 0. | 0. | 0. | 0. | 0. |
| .900 | .996 | .010 | .089 | .034 | .946 | .077 | 0. | 0. | 0. |
| .950 | .999 | .003 | .029 | .729 | .412 | **.540** | 0. | 0. | 0. |
| .960 | .999 | .003 | .027 | .925 | .175 | .522 | 0. | 0. | 0. |
| .970 | 1. | .002 | .020 | .993 | .025 | .198 | .002 | .900 | .013 |
| .980 | 1. | .002 | .021 | 1. | .001 | .008 | .043 | .986 | .049 |
| .990 | 1. | .002 | .020 | 1. | 0. | 0. | .655 | .514 | **.457** |
| .995 | 1. | .002 | .018 | 1. | 0. | 0. | .999 | .007 | .055 |
| .999 | 1. | .002 | .017 | 1. | 0. | 0. | 1. | 0. | 0. |

Table 5: Success of our adversarial weight initialization dependent on the hyperparameter $s$, which downscales the positive weights. The results depict the percentage of *active neurons* (A), *extraction-precision* (P), and *extraction-recall* (R) with a mini-batch size of 100 data points from the first fully-connected layer of the respective architectures from Table 7. All numbers are averaged over 10 runs with different adversarial initializations.

weights in Figure 8 in Appendix D shows that over training, the uniformly initialized weight values resemble more a normal distribution and obtain a wider spread, which might be the reason for the increased extraction success.

### 8.3 Evaluate Active Adversarial Initialization

We now turn to our active attack, which adversarially modifies the weight initialization to amplify the vulnerability exploited by passive attacks. This amplification is controlled by the scaling factor $s$ in the trap weights. We first evaluate its impact on the reconstruction quality of individual training data points over a mini-batch of 100 data points and 1000 neurons.

In addition to the success metrics described in the previous section, we also report the percentage of *active neurons*, *i.e.* neurons that are activated by at least one training data point from the mini-batch. Recall that our attack seeks to find an adversarial initialization that balances setting enough neurons' outputs to zero (such that a gradient is more likely to isolate individual points from large mini-batches) with, at the same time, having enough neuron outputs' that are non-zero (oth-

| | MNIST | | | CIFAR10 | | | ImageNet | | |
|---|---|---|---|---|---|---|---|---|---|
| (B, N) | A | P | R | A | P | R | A | P | R |
| (200, 20) | .522 | .436 | **.720** | .454 | .670 | **.695** | .090 | .948 | **.355** |
| (200,50) | .690 | .302 | .428 | .662 | .494 | .452 | .381 | .763 | .304 |
| (200,100) | .782 | .196 | .218 | .846 | .280 | .269 | .653 | .500 | .240 |
| (200,200) | .859 | .121 | .086 | .954 | .124 | .096 | .886 | .233 | .113 |
| (500,20) | .535 | .451 | **.915** | .452 | .689 | **.870** | .096 | .939 | **.490** |
| (500,50) | .697 | .301 | .624 | .653 | .505 | .614 | .387 | .767 | .426 |
| (500,100) | .792 | .205 | .397 | .845 | .290 | .422 | .646 | .508 | .358 |
| (500,200) | .871 | .129 | .185 | .950 | .119 | .177 | .892 | .240 | .199 |
| (1000,20) | .539 | .444 | **.950** | .441 | .703 | **.915** | .102 | .942 | **.595** |
| (1000,50) | .705 | .300 | .760 | .648 | .504 | .724 | .388 | .770 | .516 |
| (1000,100) | .796 | .203 | .540 | .844 | .297 | .556 | .655 | .514 | .457 |
| (1000,200) | .871 | .124 | .293 | .951 | .120 | .256 | .892 | .238 | .288 |
| (3000,20) | .541 | .442 | **1.** | .441 | .696 | **.945** | .101 | .934 | **.640** |
| (3000,50) | .704 | .299 | .888 | .646 | .503 | .812 | .386 | .764 | .586 |
| (3000,100) | .797 | .203 | .746 | .840 | .286 | .711 | .649 | .518 | .579 |
| (3000,200) | .873 | .129 | .504 | .951 | .122 | .414 | .889 | .243 | .404 |

Table 6: Success of our adversarial weight initialization dependent on the mini-batch size **B** and the number of neurons **N** that corresponds to the number of weights rows. The results depict the percentage of *active neurons* (A), *extraction-precision* (P), and *extraction-recall* (R). All numbers are averaged over 10 runs with different adversarial initializations.

erwise, in the limit, no points would be extracted). Thus, *active neurons* provide additional context for the *extraction-precision*: with few *active neurons*, even a high *extraction-precision* might not be able to extract many individual training data points, simply because there are very few gradients to perform data extraction from. However, with many *active neurons*, the *extraction-recall* might become small, due to each neuron being most likely activated by several input data points, preventing individual extraction.

Table 5 depicts the results, averaged over ten different random adversarial initializations. We can see that the best scaling factor for MNIST, when it comes to the *extraction-recall*, is $s = 0.7$. With this scaling factor, we are able to extract on average 54.0% of the individual training data points which were involved in the users' gradient computations. This is an improvement by around factor nine to the passive attack. For CIFAR10 and Imagenet, the best scaling factors concerning *extraction-recall* are $s = 0.95$, and $s = 0.99$, which allow for a perfect reconstruction of 54.0%, and 45.7% of the individual training data points, respectively, for 1000 neurons and a mini-batch size of 100 data points. Thereby, the active attack is more than twice as successful as the passive attack for extracting individual training data points in these datasets. Figure 9, Figure 10a, and Figure 11 in the Appendix D show the visual reconstruction results of the best runs for all three image datasets and the respective hyperparameters. Similar improvements of performance could be achieved for the IMDB dataset. The best extraction was achieved also with $s = 0.99$, for which, with 1000 neurons and mini-batches of 100 data points, we obtained an *extraction-recall* of 65.4%, which is around 2.5 time as high as the passive attack, see Table 3.

As hypothesized above, we furthermore confirm that the *extraction-recall* of our attack is related to the percentage of *active neurons*: When very few neurons are activated, it is not possible to extract large numbers of individual data points due to the lack of gradients to extract them from. However, when the percentage of *active neurons* is high, the *extraction-recall* also becomes very small, which is due to the fact that each neuron gets activated by several input data points, and thereby, individual extraction is impossible.

**Impact of Data Labels.** Additionally, we investigated whether this high reconstruction success could also be achieved when all data points in the mini-batch belong to the same class. This is a particularly challenging setting for prior work on optimization-based attacks that end up reconstructing average points rather than individual points exactly. Instead, Figure 10b in the Appendix D shows how, on CIFAR10, our method remains able to perfectly extract individual data points from the gradients even when all points stem from the same class.

**Impact of Mini-Batch Sizes.** We also set out to investigate the impact of the mini-batch size $B$ and the number of weight rows that we can use for extraction. Table 6 depicts the resulting metrics. The metrics show that the smaller the mini-batch sizes are, and the more weight rows there are for extraction, the more individual training data points can be individually reconstructed. For 3000 weight rows, even up to 50% of the individual training data points for mini-batch sizes as large as 200 in the MNIST dataset can be perfectly extracted. Small mini-batches of 20 training data points are entirely extractable without any loss in this setting. Also for the IMDB dataset, smaller batch-sizes for the same number of neurons yield much higher *extraction-recall*, and embeddings of data from small mini-batches of 20 training data points are perfectly extractable, see Table 3. This suggests that in practice, the success of the extraction attack can be significantly increased by the central party demanding smaller mini-batch sizes from the users or initializing larger models.

**Impact of Batch-Averaging.** Additionally, we looked into the effect of averaging over the gradients of multiple mini-batches, *e.g.* the average of gradients received from multiple users. The results in Table 9 in Appendix D show that through averaging, the attack success is significantly reduced. Already when averaging over 20 mini-batches of size $B = 100$ in the MNIST dataset, the average *extraction-recall* drops from 54.0% to 2.6% because multiple data points overlay in the gradients. This highlights that the central party needs to perform the extraction before the averaging operation. The following section shows that this simple change to the protocol is easily implemented by an actively dishonest central party, even for standard FL libraries.

### 8.4 TensorFlow Federated

We experimented with TensorFlow Federated [2]—a standard open source library for FL deployments. We adapted the implementation of FedAvg [1] provided by the developers to pass each individual gradient update through our reconstruction function. Note that the whole change took only minutes of work and required minimal code changes, such that they could easily be implemented by a dishonest central party. We pre-generated our adversarial initialized shared models with scaling factor $s = 0.99$ and 1000 neurons at the first fully-

connected layer, and passed them to the users. The aggregator then collects the gradients and performs reconstruction.

We find that our attack works consistently well against commonly used FL benchmarks. Over 50 users, for EMNIST [9], our trap weights yield $0.32 \pm 0.07$ *extraction-recall* and $0.05 \pm 0.02$ *extraction-precision*, versus $0.10 \pm 0.05$, and $0.02 \pm 0.01$ in the non-adversarial baseline. For CIFAR100 [30] we get $0.44 \pm 0.05$ *extraction-recall* and $0.22 \pm 0.06$ *extraction-precision*, versus $0.20 \pm 0.04$, and $0.04 \pm 0.01$ in the baseline. These results are comparable to the ones reported in the previous experiments, and confirm that our attack is practical.

## 8.5 Comparison To Previous Work

In order to compare the success of our attack with previous work, we experimented with one of the state-of-the art reconstruction method by Geiping *et al.* [17] and applied it to the MNIST and CIFAR10 datasets. For the sake of correctness, we build on their code base and adopt it to also run with neural architectures we used in other experiments. We use the parameters that Geiping *et al.* found to perform best. Here, we are mainly interested in the quality of the other attack's reconstruction in comparison to our method, and in the number of passes over the model, *i.e.* the computing time required to obtain the reconstructions.

Figure 5 in Appendix D shows an example of the gradient-inversion fidelity obtained with [17] for MNIST with $B = 1$. Within $10^4$ iterations, the pixel-wise $\ell_2$ errors observed go as low as $10^{-4}$ for FC-NNs with architecture as depicted in Table 7 in the Appendix, and $10^{-3}$ for LeNet-5. Similar results for CIFAR10 are available in Figure 6 in Appendix D. In contrast to these results, our attack allows us to extract data points perfectly, *i.e.* with $\ell_2$ error of zero, and without *any* back-propagation iterations. These results make clear that gradient inversion in practice suffers from local minima and requires a very large number of iterations to converge to a comparable reconstruction to our method. On a practical note, reconstruction of a single CIFAR10 image with 32 restarts from a seven-layer FC-NN takes on average 1 hour and 3 minutes on a high-end GPU, in comparison to milliseconds needed for extraction with the help of our trap weights. Most importantly, it is clear that reconstruction is not a lossless process and full data recovery is almost never possible, even in the simple cases where gradients of only a single data point are considered. To better understand limitations of prior literature we refer the reader to [46].

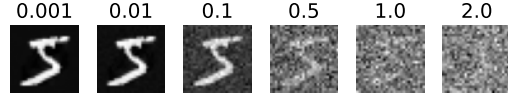## 9 Differential Privacy and Other Defenses

We now evaluate whether applications of DP to FL are an adequate defense against our attack. We then consider other defense mechanisms, notably those that rely on cryptography.

### 9.1 Current Application of DP in FL

In general, FL frameworks implement DP through the introduction of a DP optimizer, such as DPSGD [4] (see Section 2).



(a) Effect of applying different clip norms $c$ at the users' end without adding noise.



(b) Effect of applying different noise scales at the users' end when clipping by $c = 1.0$.

Figure 4: Qualitative effect of gradient clipping and noise addition at the users' end on the reconstructed data. Results depicted for one single data point from the MNIST dataset, recovered at the first fully-connected model layer.

Recall that DPSGD requires that one computes gradients for each data point individually and clips these gradients before noising their average, which is only then applied to update the model. Users, therein, perform the clipping operation, and the central party applies the DP optimizer for training, such as shown in [35]. TensorFlow Federated [2] is an example of a state-of-the-art-framework that integrates DP in this fashion [3].

It is important to note that DP guarantees can only be achieved when both parts, the clipping and the noise addition, are correctly performed. For instance, Figure 4a illustrates that clipping alone does not yield privacy since it is still possible to perfectly extract the original training input from the user's gradients. Hence, in scenarios with a dishonest central party who can simply omit the noise addition, DP in its current application to FL does not represent an effective defense mechanism against attacks described in this paper.

### 9.2 Local DP: Noise Addition at the Users' End

Figure 4b, depicts that noise addition to the users' gradients is effective in order to prevent perfect extractability of individual training data points from the gradients. Thus, a solution to the privacy issue in FL would be for the users to apply calibrated noise to their gradients before sending them to the central party. According to Figure 4b, an adequate amount of noise to mitigate our attack would be a $\sigma > 1$. However, adding noise with such a large scale comes at the costs of reducing model utility. Recent work [45] suggests that it is not possible to train a DP-model for the CIFAR10 dataset above an accuracy of 70%. For more complex datasets, such as ImageNet, it is currently impossible to train a model with DPSGD and a reasonable accuracy, even in non-distributed applications. Yet, as we demonstrate in this paper, local DP remains the only viable solution for achieving privacy with some reported use in the industry [7, 39]. However, how widely it is adopted is unclear.

## 9.3 Other Existing Defenses

**Gradient hiding.** Apart from DP, other defenses to add privacy to the FL protocol have been put forward. Zhu *et al*. [52] proposed an approach for *gradient compression and sparsification* in which gradients with small magnitudes are pruned to zero. In practical evaluation on image data, [44] showed that even under this defense, the extracted data is still recognizable when 80% of the gradients are pruned. Thereby, the defense manages to prevent a perfect extactabiliy, however, it does not prevent the general privacy leakage. Sun *et al*. [44] presented an improvement over this simple pruning strategy where the architecture gets augmented with additional defensive layers to promote bounded gradient invertibility. These defensive layers are trained alongside the main task. Ultimately, these defenses are a special case of noised gradients; thus they share the limitations of DPSGD described previously—with the additional disadvantage that they do not provide any rigorous guarantees of DP.

**Local Training and Averaging.** As described in Section 4, we consider an FL protocol where users share per-batch gradients. Different implementations of FL also give users the possibility of averaging gradients over several data mini-batches or performing local training over several epochs on multiple mini-batches [35] before sharing the resulting gradients with the central party. Our results (see Table 9) suggest that such averaging reduces the success of our attack. However, in practical applications, including the local averaging is often hard since (1) it introduces a form of asynchronism between the users and (2) might stand in conflict with convergence of the shared model. Finally, we find that using larger batches does not provide any privacy *guarantees* against accidental or malicious leakage.

**Privacy in depth.** Kairouz *et al.* surveyed the state of FL, including security and privacy mechanisms that augment FL. The authors advocate for privacy in depth and described means of improving privacy guarantees of FL using protection mechanisms on multiple layers of the technology stack [28]. To this end, they list Trusted Execution Environments (TEE), Zero-Knowledge Proofs (ZKP), Multi-Party Computation (MPC) and Homomorphic cryptography (HC).

At the time of writing, these methods however still largely fail to provide privacy to FL and put trust into the correct and honest operation of either the aggregator, the central party, or another additional component such as a key server. Furthermore, TEE, ZKP, MPC and HC all carry significant overheads in practice, both computationally and in terms of communication rounds; what is more, they still require trust for their operation to provide the privacy guarantees needed.

## 10  Discussion and Future Directions

In this section, we first discuss the detectability of our adversarial initialization and the impact of our attack on training an accurate shared model. We then argue about the potential and capabilities of adversarial weight initialization for future privacy attacks. Finally, we elaborate on the core problem with the expectations to privacy in FL, and propose an outlook on a redesign of the protocol that we argue is necessary.

### 10.1  Detectability of Adversarial Weight Initialization

The main question that we are concerned with when it comes to the topic of detectability is: *Can a user detect our adversarial weight initialization?* The simple answer is no. Indeed, recall that we showed in Section 5 that partial data extraction is possible even with *unmodified* benign model weights when the adversary is simply a passive honest-but-curious central party. This means detection is impossible, at least in the limit. However, for the sake of completeness, we nevertheless consider two detection strategies: (1) analyzing the model weights of the shared model in *one* iteration over the FL protocol, or (2) analyzing the model weights of the shared model over *multiple* iterations of the FL protocol.

**One Iteration.** When receiving the shared model, the users could run a detection mechanism on the weights in order to identify manipulations. This detection mechanism could include deliberately looking for elements of our adversarial initialization, such as the convolutional layers which consist of all *zero* elements and a single *one*, or the presence of slightly larger negative than positive weights at fully-connected layers. The detection could be implemented, for example, in form of binary ML classifiers, or statistical distribution-distance tests such as the Kolmogorov-Smirnov test, to distinguish between benign weights and our adversarially initialized ones. However, such detection mechanisms would have a high false positive rate, since convolutional layers consisting mainly of zeros could also result from compression of ML models through sparsity (*e.g*. [24]) or feature suppression methods (*e.g*. [16]) whilst a surplus of negative weights could just be the result of previous training.

**Multiple Iterations.** Having access to the shared model over multiple iterations of the FL protocol enables users to compare the received shared model's parameters to the parameters from previous FL iterations. Therefore, the success of detection boils down to the following question: *Can a local user tell that a given set of model parameters came from legitimate updates of other local users?*

Unfortunately, the answer to the question above is no, even for non-FL setups where the entire training procedure is transparent. The stochastic nature of training algorithms, combined with the non-determinism of modern hardware, makes it difficult to reproduce training runs [27]. Because of this reproducibility error, an attacker can assemble a mini-batch of natural data points that produce *any* desired gradient update [43]. In other words, given a pair of parameter updates, the user cannot tell if the gradient descent step between these parameters was a result of a legitimate optimization step. This is exacerbated in FL because the data of any given user is invisible to other users, further complicating the verification of gradient descent integrity. As a consequence, a local user always has to put trust into the global shared model sent by the central party and assume it comes with meaningful parameters—in particu-

lar, parameters that do not lead to trivial extraction of training data.

In addition to analyzing the received model weights for detection of our attack, a user could also rely on observing the *functionality* of the shared model. We observe that, for the vast majority of classification tasks, the model's loss across training data points significantly reduces after just a few iterations. Thus, after the initial training iterations, FL clients would expect to encounter low loss values for their own examples. However, research has shown that, in particular for users whose data stems from the tails of the data distribution, FL does not necessarily lead to an improvement of model accuracy on their data [50]. Therefore, detection mechanisms that rely on analyzing the convergence of the model's accuracy over multiple FL iterations would also have a high false positive rate. Additionally, the central party could rely on methods that we discuss in the next Section 10.2, in order to increase model accuracy over time or target different users over different iterations of the FL protocol. Thereby, an individual user would have little to no chance to rely on model performance as an indicator for the presence of our attack.

As a side note, it is interesting to observe that this discussion raises the question of how a central party should take into account noisy updates: a user could decide to send random gradients instead of the actual gradients computed on the dataset. This would affect the utility [6] of the FL solution and shatter its already fragile economic foundations [50].

## 10.2   Training Success and Model Performance

Independent of the attack detectability, increasing the utility of the model over the course of training is important because usually the central party is expected to provide well-performing model after several training iterations.

In order to successfully extract a large proportion of the users' data, the central party needs to perform an adversarial initialization each time before sending the model out to the users. Thereby, the continuous training process is interrupted, and each time, a new model is sent out, it practically results in little to no accuracy gain over multiple FL iterations. The central party has several options in order to still obtain a useful model. (1) Instead of aiming at reconstructing the user data over all communications, the central party could send the adversarially (re-)initialized model out only at a few communication rounds, preferably at the beginning of the training process where lower accuracy is expected. In all other communications, it could send out the actual updated model without adversarially re-initializing it. (2) Instead of targeting all users, the central party could also only target a subset of users, and send out an adversarially initialized model to them, and a continuously trained more accurate model to all other users. Potentially, both strategies could also be combined by sending out an adversarially initialized model only to a subset of users in a few iterations.

## 10.3   The Power of Weight Initialization

In general, even outside of the FL context, our attack shows that controlling and manipulating the weights of neural networks opens a new attack surface against ML. We argue that weight manipulation could be used to design further privacy attacks outside of the FL context. Our adversarial initialization of convolutional and fully-connected model layers is able to transmit input data points to any subsequent layer in the model, practically modulating data perfectly over them. Additionally by setting our trap weights, we can increase the leakage of individual training data points from model gradients. Therefore, the trap weights basically create a simple if-else logic based on $>$ and $<$ relations between weighted inputs to model neurons. Future work could investigate whether the weights could also be set in order to implement more complex logical structures and if-else cases depending on the input. Based on these, it might be possible to craft hybrid attacks that first initialize the model weights and then use that to later extract information, for example, on membership of individual data points, or these data points' sensitive attributes. Note that adversarially setting weights also does not need to be limited to initializing the model weights. Instead, given an already initialized (and trained model), it might be possible to craft additional training data that leads to the weights taking the adversarial values that an attacker wants.

## 10.4   Privacy and Confidentiality—Updating the FL protocol

FL was originally designed as an alternative to centralized ML in which no large datasets would have to be moved from users to a central party in order to train an ML model on the joint data. The approach does not only reduce communication costs but also spares the central party from having to build up the infrastructure by outsourcing training and data storage costs to the users. Indeed, FL is more communication cost effective, and guarantees a form of confidentiality, since the data itself is not shared directly.

Our work, however, shows that the gradients shared by the users leak information that can be used to extract their individual training data points. In particular when dealing with a dishonest central party, the leakage of individual data points can be drastically increased, such that large parts of the users' data can be perfectly extracted. Even worse, we highlight that most current defenses are not designed to deal with such active attackers, and hence yield little to no protection.

We saw that the extraction success of individual private user data points in FL is tightly bounded to the gradient aggregation at the central party. In order to prevent the dishonest and active central party from bypassing the averaging mechanism and extracting data directly from individual user gradients, in a redesigned version of the FL protocol, the central party would have to prove to the users that the averaging was (correctly) executed. This could be implemented by using some cryptography-based protocol elements, inspired by its application for privacy in the setting of centralized ML, *e.g.* [8]. However, under the threat of a dishonest and active central

party, this proof would not be sufficient because the central party would still hold control over the data that goes into the averaging mechanism. Therefore, the central party could, for example, select a target user for the extraction attack, and perform an average over this user's gradients and a number of 'fake' zero gradients and prove having executed the averaging, but, the target user's data would remain perfectly extractable. To overcome this threat, each user who participates in an iteration of the FL protocol would have to prove to all other users that they contributed their gradients to the learning, and the central party would then have to prove that the averaging was executed over exactly these gradients. Such proofs would require many interactions between the users, and thereby cause a significant increase of the communication costs of the protocol. In most applications, this would not be practically implementable. Additionally, given that the users would have to heavily interact with each other, it is unclear what benefit the central party would have in the protocol, and why the users could not just perform a secure aggregation as, for example in secure multi party-computation.

Ultimately, in its current form, FL provides no privacy against both accidental and malicious leakage. The central party will always have an upper hand over how much data a local model will leak, unless an additional control point to validate training provenance is introduced, or the master-puppet paradigm underlying FL communication is replaced. In practice, there are many cases where such power dynamics have led to privacy compromises *e.g.* in communication networks [42]. Therefore, FL needs a re-design: it can still be advantageous to perform training in a distributed fashion, but novel mechanisms are needed to guarantee privacy. Put another way, FL is not a privacy-enhancing technology.

## 11 Conclusion

In this work, we presented a new privacy attack against FL that is based on an active attacker who holds the ability to maliciously manipulate the shared model and its weights. Recent work [46] argues that privacy for FL might not be broken because the attacks have too specific assumptions on small minibatch sizes and on the underlying data distribution. Yet, our results demonstrate the opposite: our attack allows for perfect reconstruction of a significant portion of the users' private training data. Even for very high-dimensional complex datasets, such as ImageNet, we are able to perfectly extract more than 50% of the individual data points from mini-batches of sizes as large as 100. The extraction is computationally highly efficient and even allows to perfectly extract individual training data points from data mini-batches containing all data points from one single class.

Additionally, we highlighted that existing privacy mechanisms like DPSGD are not able to protect the private user data against our data reconstruction attack. This is because the current state-of-the-art implementations of FL operate under the assumption of an honest-but-curious central party. We therefore argue that the FL protocol requires adaptations to counter the threat that arises from active attackers, and to provide meaningful privacy for the users private training data.

## References

[1] Implementation of FedAvg in TensorFlow Federated. `https://github.com/tensorflow/federated/tree/v0.19.0/tensorflow_federated/python/examples/simple_fedavg`. Accessed: 10/2021.

[2] TensorFlow Federated. `https://www.tensorflow.org/federated`. Accessed: 10/2021.

[3] TensorFlow Federated With Differential Privacy. `https://www.tensorflow.org/federated/tutorials/federated_learning_with_differential_privacy`. Accessed: 10/2021.

[4] Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In Edgar Weippl, editor, *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2016.

[5] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, and *et al*. *TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems*. 2016.

[6] Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How to backdoor federated learning, 2019.

[7] Abhishek Bhowmick, John Duchi, Julien Freudiger, Gaurav Kapoor, and Ryan Rogers. *Protection Against Reconstruction and Its Applications in Private Federated Learning*. 2019.

[8] Andrea Bittau, Úlfar Erlingsson, Petros Maniatis, Ilya Mironov, Ananth Raghunathan, David Lie, Mitch Rudominer, Ushasree Kode, Julien Tinnes, and Bernhard Seefeld. Prochlo: Strong privacy for analytics in the crowd. In *Proceedings of the 26th Symposium on Operating Systems Principles*, pages 441–459, 2017.

[9] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. *2017 International Joint Conference on Neural Networks (IJCNN)*, 2017.

[10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[11] Di Xie, Jiang Xiong, and Shiliang Pu. All you need is beyond a good init: Exploring better solution for training

extremely deep convolutional neural networks with orthonormality and modulation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017.

[12] Cynthia Dwork. Differential privacy. 2006.

[13] Liam Fowl, Jonas Geiping, Wojtek Czaja, Micah Goldblum, and Tom Goldstein. Robbing the fed: Directly obtaining private data in federated learning with modified models. *arXiv preprint arXiv:2110.13057*, 2021.

[14] Stephen Gallant. Perceptron-based learning algorithms. *IEEE Transactions on neural networks*, 1(2):179–191, 1990.

[15] Karan Ganju, Qi Wang, Wei Yang, Carl A. Gunter, and Nikita Borisov. Property inference attacks on fully connected neural networks using permutation invariant representations. In David Lie, Mohammad Mannan, Michael Backes, and XiaoFeng Wang, editors, *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, pages 619–633. ACM, 2018.

[16] Xitong Gao, Yiren Zhao, Łukasz Dudziak, Robert Mullins, and Cheng-zhong Xu. Dynamic channel pruning: Feature boosting and suppression. *arXiv preprint arXiv:1810.05331*, 2018.

[17] Jonas Geiping, Hartmut Bauermeister, Hannah Dröge, and Michael Moeller. *Inverting Gradients – How easy is it to break privacy in federated learning?* 2020. 23 pages, 20 figures. The first three authors contributed equally.

[18] Nastaran Gholizadeh and Petr Musilek. Distributed learning applications in power systems: A review of methods, gaps, and challenges. 14(12):3654, 2021. PII: en14123654.

[19] Raja Giryes, Guillermo Sapiro, and Alex M. Bronstein. Deep neural networks with random gaussian weights: A universal classification strategy? 64(13):3444–3457, 2016.

[20] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. pages 315–323, 2011.

[21] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, Montreal, Canada, December 2014.

[22] Kathrin Grosse, Thomas A. Trost, Marius Mosbach, and Michael Backes. Adversarial initialization - when your network performs the way i want -. 2019.

[23] Saqib Hakak, Suprio Ray, Wazir Zada Khan, and Erik Scheme. A framework for edge-assisted healthcare data analytics using federated learning. In *2020 IEEE International Conference on Big Data (Big Data)*. IEEE, 2020.

[24] Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections for efficient neural networks, 2015.

[25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *2015 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2015.

[26] Briland Hitaj, Giuseppe Ateniese, and Fernando Perez-Cruz. Deep models under the gan: information leakage from collaborative deep learning. 2017.

[27] Hengrui Jia, Mohammad Yaghini, Christopher A. Choquette-Choo, Natalie Dullerud, Anvith Thudi, Varun Chandrasekaran, and Nicolas Papernot. Proof-of-learning: Definitions and practice, 2021.

[28] Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, and Rachel Cummings *et al*. Advances and open problems in federated learning, 2021.

[29] Latif U. Khan, Walid Saad, Zhu Han, Ekram Hossain, and Choong Seon Hong. Federated learning for internet of things: Recent advances, taxonomy, and open challenges. 23(3):1759–1799, 2021.

[30] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.

[31] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

[32] Yann LeCun, Corinna Cortes, and C. J. Burges. *MNIST handwritten digit database*. 2010.

[33] Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics.

[34] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. pages 1273–1282, 2017.

[35] H Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially private recurrent language models. In *International Conference on Learning Representations*, 2018.

[36] Luca Melis, Congzheng Song, Emiliano de Cristofaro, and Vitaly Shmatikov. Exploiting unintended feature leakage in collaborative learning. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 691–706. IEEE, 19/05/2019 - 23/05/2019.

[37] Dinh C. Nguyen, Ming Ding, Pubudu N. Pathirana, Aruna Seneviratne, Jun Li, and H. Vincent Poor. *Federated Learning for Internet of Things: A Comprehensive Survey*, volume 23. 2021.

[38] Nicolas Papernot, Martín Abadi, Úlfar Erlingsson, Ian Goodfellow, and Kunal Talwar. Semi-supervised knowledge transfer for deep learning from private training data, 18/10/2016. Accepted to ICLR 17 as an oral.

[39] Matthias Paulik, Matt Seigel, Henry Mason, Dominic Telaar, Joris Kluivers, Rogier van Dalen, Chi Wai Lau, Luke Carlson, Filip Granqvist, Chris Vandevelde, Sudeep Agarwal, Julien Freudiger, Andrew Byde, Abhishek Bhowmick, Gaurav Kapoor, Si Beaumont, Áine Cahill, Dominic Hughes, Omid Javidbakht, Fei Dong, Rehan Rishi, and Stanley Hung. Federated evaluation

and tuning for on-device personalization: System design & applications, 2021.

[40] Le Trieu Phong, Yoshinori Aono, Takuya Hayashi, Lihua Wang, and Shiho Moriai. Privacy-preserving deep learning: Revisited and enhanced. pages 100–110. Springer, Singapore, 2017.

[41] Rudiger Pryss, Manfred Reichert, Jochen Herrmann, Berthold Langguth, and Winfried Schlee. Mobile crowd sensing in clinical and psychological trials – a case study. In *2015 IEEE 28th International Symposium on Computer-Based Medical Systems*. IEEE, 2015.

[42] Altaf Shaik, Ravishankar Borgaonkar, N. Asokan, Valtteri Niemi, and Jean-Pierre Seifert. Practical attacks against privacy and availability in 4g/lte mobile communication systems, 2017.

[43] Ilia Shumailov, Zakhar Shumaylov, Dmitry Kazhdan, Yiren Zhao, Nicolas Papernot, Murat A. Erdogdu, and Ross Anderson. *Manipulating SGD with Data Ordering Attacks*. 2021.

[44] Jingwei Sun, Ang Li, Binghui Wang, Huanrui Yang, Hai Li, and Yiran Chen. Provable defense against privacy leakage in federated learning from representation perspective. *arXiv preprint arXiv:2012.06043*, 2020.

[45] Florian Tramèr and Dan Boneh. Differentially private learning needs better features (or much more data). *arXiv preprint arXiv:2011.11660*, 2020.

[46] Aidmar Wainakh, Ephraim Zimmer, Sandeep Subedi, Jens Keim, Tim Grube, Shankar Karuppayah, Alejandro Sanchez Guinea, and Max Mühlhäuser. Federated learning attacks revisited: A critical discussion of gaps, assumptions, and evaluation setups, 2021.

[47] Yi Wang, Imane Lahmam Bennani, Xiufeng Liu, Mingyang Sun, and Yao Zhou. Electricity consumer characteristics identification: A federated learning approach. 12(4):3637–3647, 2021.

[48] Zhibo Wang, Mengkai Song, Zhifei Zhang, Yang Song, Qian Wang, and Hairong Qi. Beyond inferring class representatives: User-level privacy leakage from federated learning. In *IEEE INFOCOM 2019 - IEEE Conference on Computer Communications*. IEEE, 2019.

[49] Hongxu Yin, Arun Mallya, Arash Vahdat, Jose M. Alvarez, Jan Kautz, and and Pavlo Molchanov. *See through Gradients: Image Batch Recovery via GradInversion*. 2021.

[50] Tao Yu, Eugene Bagdasaryan, and Vitaly Shmatikov. Salvaging federated learning by local adaptation, 2021.

[51] Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. *iDLG: Improved Deep Leakage from Gradients*. 2020.

[52] Ligeng Zhu and Song Han. Deep leakage from gradients. In *Federated Learning*, pages 17–31. Springer, Cham, 2020.

## A  Differential Privacy and DPSGD

DP [12] formalizes the idea that no data point should significantly influence the results of an analysis conducted on a whole dataset. Therefore, the analysis should yield roughly the same results whether or not any particular dataset is included in the dataset or not. Formally, this is expressed by the following definition.

**Definition 1** (($\varepsilon, \delta$)-Differential Privacy)**.** *Let* $A\colon \mathcal{D}^* \to \mathcal{R}$ *a randomized algorithm.* $A$ *satisfies* ($\varepsilon, \delta$)*-DP with* $\varepsilon \in \mathbb{R}_+$ *and* $\delta \in [0,1]$ *if for all neighboring datasets* $D \sim D'$*, i.e. datasets that differ on only one element, and for all possible subsets* $R \subseteq \mathcal{R}$

$$\mathbb{P}\left[A(D) \in R\right] \leq e^\varepsilon \cdot \mathbb{P}\left[M(D') \in R\right] + \delta. \qquad (11)$$

In ML, DP formalizes the idea that any possible training data point in a training dataset cannot significantly impact the resulting ML model [4, 38]. One notable approach to achieve this is Differentially Private Stochastic Gradient Descent (DPSGD) [4]. DPSGD alters the training process to introduce DP guarantees for weight update operations, and thereby protects underlying individual data points. Particularly, the gradient computed for each data point or a mini-batch of data points is first clipped in their $\ell_2$-norm to bound influence. Clipping of the gradient $g(\{x_i\}_b)$ for mini-batch $b$ of data $\{x_i\}_b$ is performed according to a clipping parameter $c$, by replacing $g(\{x_i\}_b)$ with $g(\{x_i\}_b)/\max(1, \frac{\|g(\{x_i\}_b)\|_2}{c})$. This ensures that if the $\ell_2$-norm of the gradients is $\leq c$, the gradients stays unaltered, and if the norm is $> c$, the gradient get scaled down to be in norm of $c$. After the clipping, Gaussian noise with scale $\sigma$ is applied to the gradients of each mini-batch before performing the model updates. DP is commonly used to make FL private – as described in Section 9.1.

---

**Algorithm 2:** Optimization-Based Data Rec. [52].

---

**Input:** Gradients, $G_i^{[t]}$, received from victim user $u_i$ at iteration $t$, Shared model $f_\mathcal{W}^{[t]}(\cdot)$ at iteration $t$.
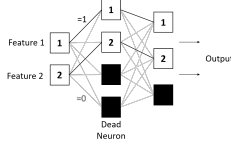**Output:** Reconstructed training data, $(\mathbf{x}_i^*, y_i^*)$

1:  $(\hat{\mathbf{x}}^{[1]}, \hat{y}^{[1]}) \leftarrow (\mathcal{N}(0,1), \mathcal{N}(0,1))$      ▷ Initialize
2:  **for** $\hat{t} \in [1, \hat{T}]$ **do**
3:     $\hat{G}^{[\hat{t}]} = \nabla_\mathcal{W} \mathcal{L}(f_\mathcal{W}^{[t]}(\hat{\mathbf{x}}^{\hat{t}}), \hat{y}^{\hat{t}})$    ▷ Dummy gradients
4:     $D^{[\hat{t}]} = \|G_i^{[t]} - \hat{G}^{[\hat{t}]}\|^2$    ▷ Dummy vs user gradients
5:     $\hat{\mathbf{x}}^{[\hat{t}+1]} \leftarrow \hat{\mathbf{x}}^{[\hat{t}]} - \alpha \nabla_{\hat{\mathbf{x}}^{[\hat{t}]}} D^{[\hat{t}]},$
6:     $\hat{y}^{[\hat{t}+1]} \leftarrow \hat{y}^{[\hat{t}]} - \alpha \nabla_{\hat{y}^{[\hat{t}]}} D^{[\hat{t}]}$
7:  **end for**
8:  $(\mathbf{x}_i^*, y_i^*) \leftarrow (\hat{\mathbf{x}}^{[\hat{T}+1]}, \hat{y}^{[\hat{T}+1]})$

---

## B  Transmitting Data Points in FC-NNs

In this section, we show how to initialize a fully-connected model layer adversarially, such that it transmits its input data to the next layer without altering the input data.

Assuming that the fully-connected model layer has at least as many neurons as there are input data features, a very simple initialization would be the following: Set all components in the weight matrix at that layer to *zero*, apart from one component per input data feature, which is set to *one* in order to

transmit the corresponding feature to the next layer. Thereby, model input could potentially be transmitted all the way from the model's input layer to the output layer:



In the resulting model, all neurons that do not carry an input feature would have a zero in- and output, and, thereby act as dead neurons.

A layer that is initialized like this would be detectable by a users aware of the method. It is, however, possible to hide the adversarial initialization and make it more indistinguishable from random model weights. For the sake of brevity, we won't go into details and only sketch the necessary ideas: (1) The first feature does not need to always be transmitted to the first neuron in the subsequent layer, instead, each feature can be transmitted to a random separate neuron. (2) The weights that transmit features do not need to be set to *one*, but can be set to any random positive value. Thereby, at the layer where the input should be transmitted to, this input arrives scaled feature-wise by a factor depending on all the positive values that the feature was multiplied with on the way. Hence, after extracting the input data from the gradients at that layer, it would only be necessary to rescale it feature-wise according to the inverse of the respective factors. (3) The weights transmitting input features to neurons that do not carry any input feature in the subsequent layer can be set to any random value. In particular, if they contain mainly negative random values, the weighted input to the respective neurons will be negative, and the $\mathrm{ReLU}$ activation will turn these neurons into dead neurons. (4) The weights related to transmitting the output of any dead neuron can be set to any random positive or negative value, since they will be multiplied with the dead neuron's *zero*-output, and thereby, yield *zero* input in the subsequent layer anyways. Only the weight components transmitting input features to neurons carrying different features in the subsequent layer should stay *zero*, such that the different input features are not overlaid and stay separate.

## C  Additional Material

The following table describes the model architectures both for the FC-NNs and CNNs use throughout the paper. Note that the our method could also be applied to much larger CNNs with more layers: in fact, as long as each layer contains as many parameters as the data holds input features, our approach is applicable.

## D  Additional Experimental Results

| FC-NN Architecture | VGG-inspired CNN Architecture |
|---|---|
| Dense(n=1000, act=relu) | Conv(f=128, k=(3,3), s=1, p=same, act=relu) |
| Dense(n=3000, act=relu) | Conv(f=256, k=(3,3), s=1, p=same, act=relu) |
| Dense(n=3000, act=relu) | Conv(f=512, k=(3,3), s=1, p=same, act=relu) |
| Dense(n=2000, act=relu) | Flatten |
| Dense(n=1000, act=relu) | Dense(n=1000, act=relu) |
| Dense(n=#classes, act=None) | Dense(n=#classes, act=None) |

Table 7: Architectures of models used in the experiments on image data. f: number of filters, k: kernel size, s: stride, p: padding act: activation function, n: number of neurons.

| IMDB-Model Architecture |
|---|
| Embedding(feat=10000, dim=250) |
| Dense(n=1000, act=relu) ) |
| Dense(n=1, act=None) |

Table 8: Architecture of models used in the experiments on the IMDB dataset. feat: vocabulary size, dim: embedding size, act: activation function, n: number of neurons.

| B, Num | A | P | R |
|---|---|---|---|
| (20,1) | .544 | .459 | .995 |
| (20,5) | .471 | .540 | .190 |
| (20,10) | .461 | .508 | .100 |
| (20,20) | .484 | .468 | .047 |
| (50,1) | .706 | .317 | .750 |
| (50,5) | .630 | .327 | .169 |
| (50,10) | .716 | .311 | .082 |
| (50,20) | .655 | .332 | .040 |
| (100,1) | .798 | .214 | .542 |
| (100,5) | .822 | .197 | .107 |
| (100,10) | .783 | .230 | .060 |
| (100,20) | .759 | .248 | .026 |

Table 9: Success of our adversarial weight initialization on MNIST under averaging over multiple mini-batches on the same model parameters. The number of mini-batches is denoted by **Num** and their respective size by **B**. The results depict the percentage of *active neurons* (A), *extraction-precision* (P), and *extraction-recall* (R) for extracting from 1000 neurons at the first layer of the FC-NN depicted in Table 7. All numbers are averaged over 10 runs with different adversarial initializations.

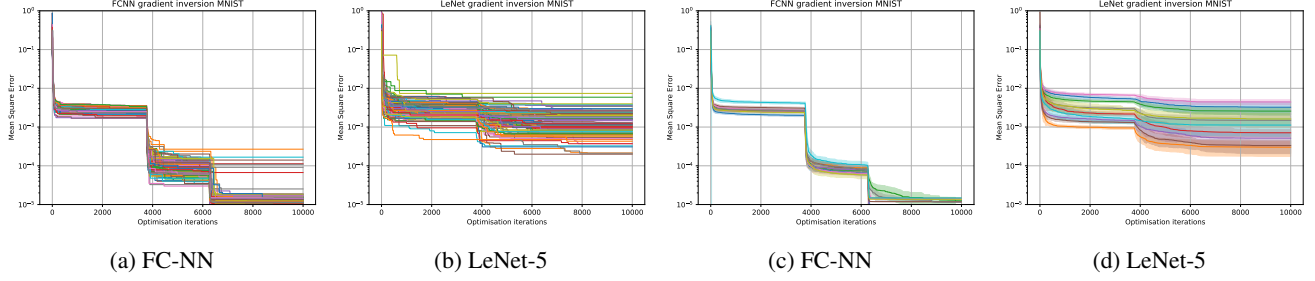(a) FC-NN     (b) LeNet-5     (c) FC-NN     (d) LeNet-5

Figure 5: **Baseline—Prior Work.** Single sample gradient inversion with untrained network using the inversion method proposed by [17] for the first 100 MNIST datapoints. (a) and (b) shows fidelity of individual datapoint reconstruction with no restarts, while (c) and (d) show 32 different optimisation start points. Error bars are a single standard deviation of individual restarts.
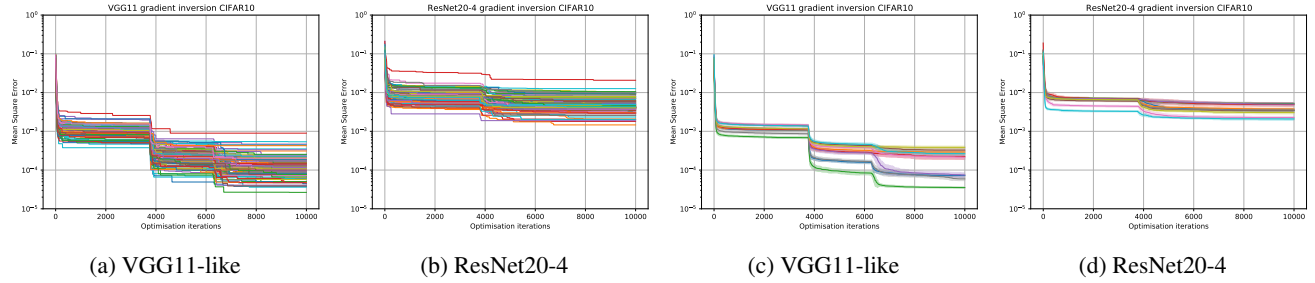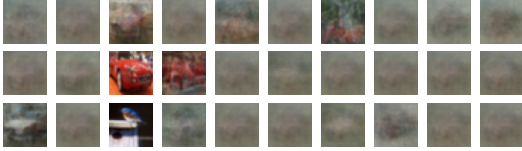


(a) VGG11-like     (b) ResNet20-4     (c) VGG11-like     (d) ResNet20-4

Figure 6: **Baseline—Prior Work.** Single sample gradient inversion with untrained network using the inversion method proposed by [17] for the first 100 CIFAR10 datapoints. (a) and (b) shows fidelity of individual datapoint reconstruction with no restarts, while (c) and (d) show 32 different optimisation start points. Error bars are a single standard deviation of individual restarts.



Figure 7: **Baseline—Passive Attack.** Data from the CIFAR10 dataset, extracted from the gradients of the first 30 weight rows at the first fully-connected layer of a randomly initialized FC-NN with architecture from Table 7.
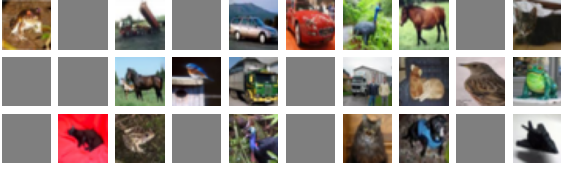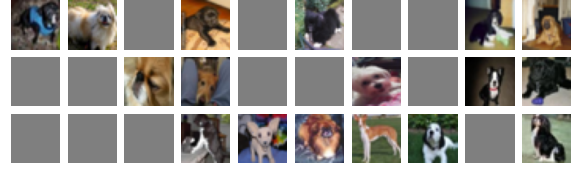


(a) Reconstructed data points.
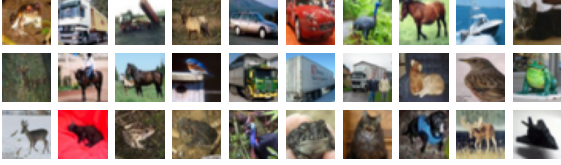


(b) Original data points.

Figure 9: MNIST. Reconstruction success of our adversarial initialization: first 30 images from a mini-batch of 100 data points, extracted at the first fully-connected layer of the FC-NN from Table 7. Gray images indicate that the corresponding original data point could not be extracted individually from the model gradients.



Figure 8: Distribution of the first layer's weights of the FC-NN from Table 7 over training on the MNIST dataset. Weights at epoch zero were initialized with a random uniform distribution.
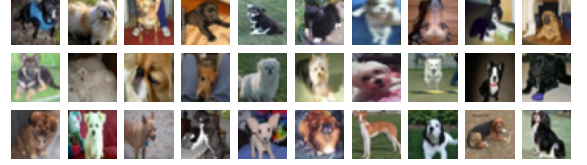
(a) Reconstructed data points from different classes.



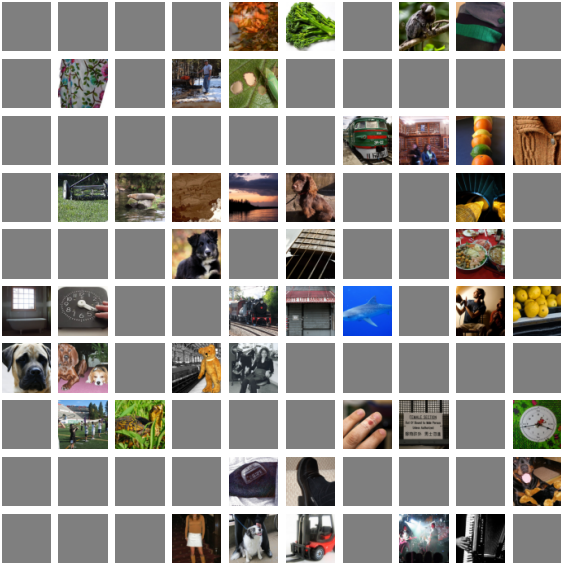(b) Reconstructed data points from class "dog".
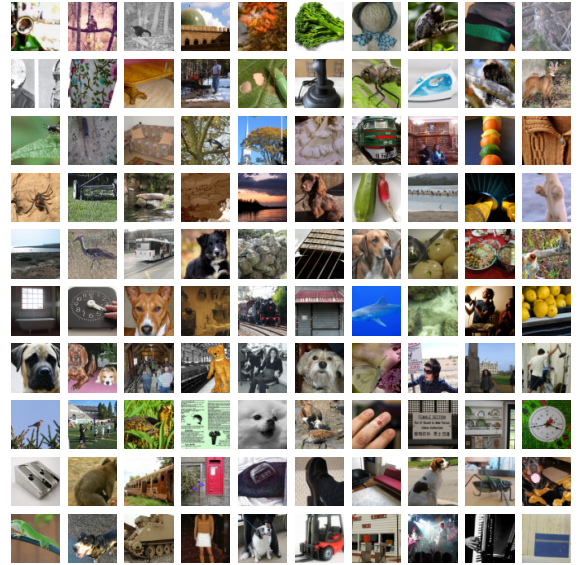


(c) Original data points from different classes.



(d) Original data points from class "dog".

Figure 10: CIFAR10. Reconstruction success of our adversarial initialization: first 30 images from a mini-batch of 100 data points, extracted at the first fully-connected layer of the CNN from Table 7. Gray images indicate that the corresponding original data point could not be extracted individually from the model gradients.



(a) Reconstructed data points.



(b) Original data points.

Figure 11: ImageNet. Reconstruction success of our adversarial initialization: all reconstructed data points from a mini-batch of 100 data points, extracted at the first fully-connected layer of the CNN from Table 7. Gray images indicate that the corresponding original data point could not be extracted individually from the model gradients..

19