# Can Individuals Trust Privacy Mechanisms for Machine Learning?

# A Case Study of Federated Learning

Franziska Boenisch
franziska-boenisch.de

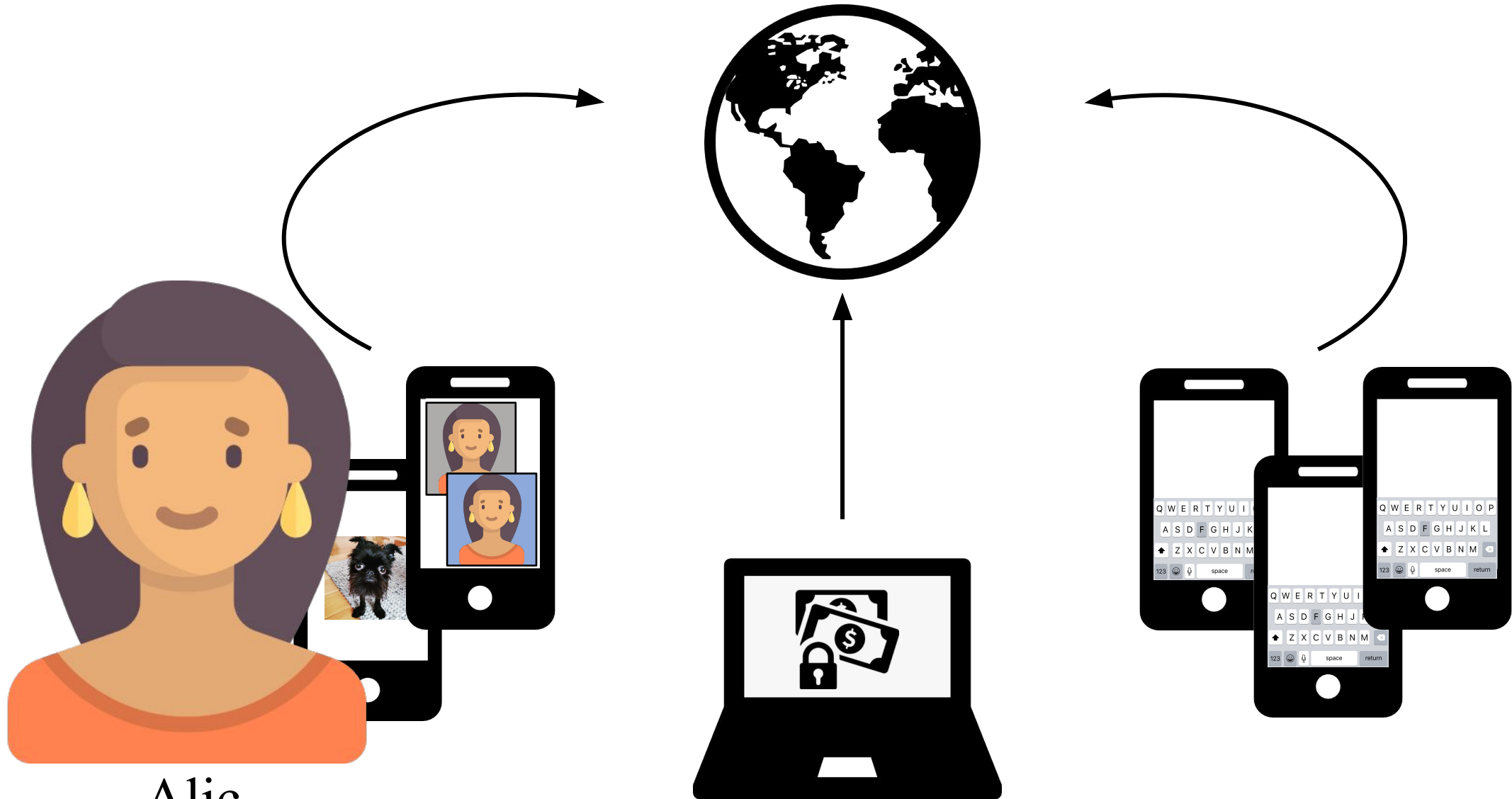ML in PL Conference
26 - 29 OCTOBER 2023

Apply for Open
Positions in my
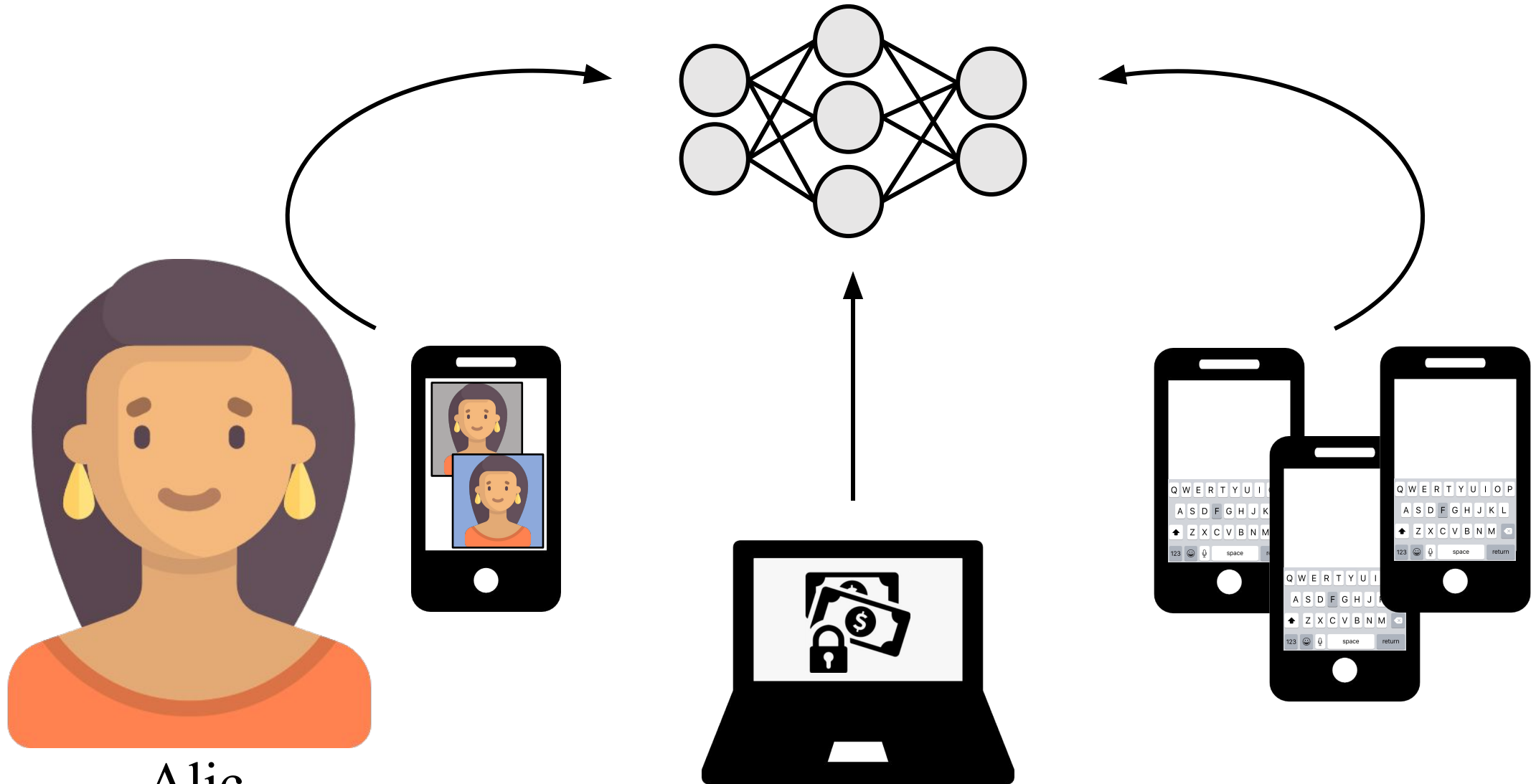
CISPA
HELMHOLTZ CENTER FOR
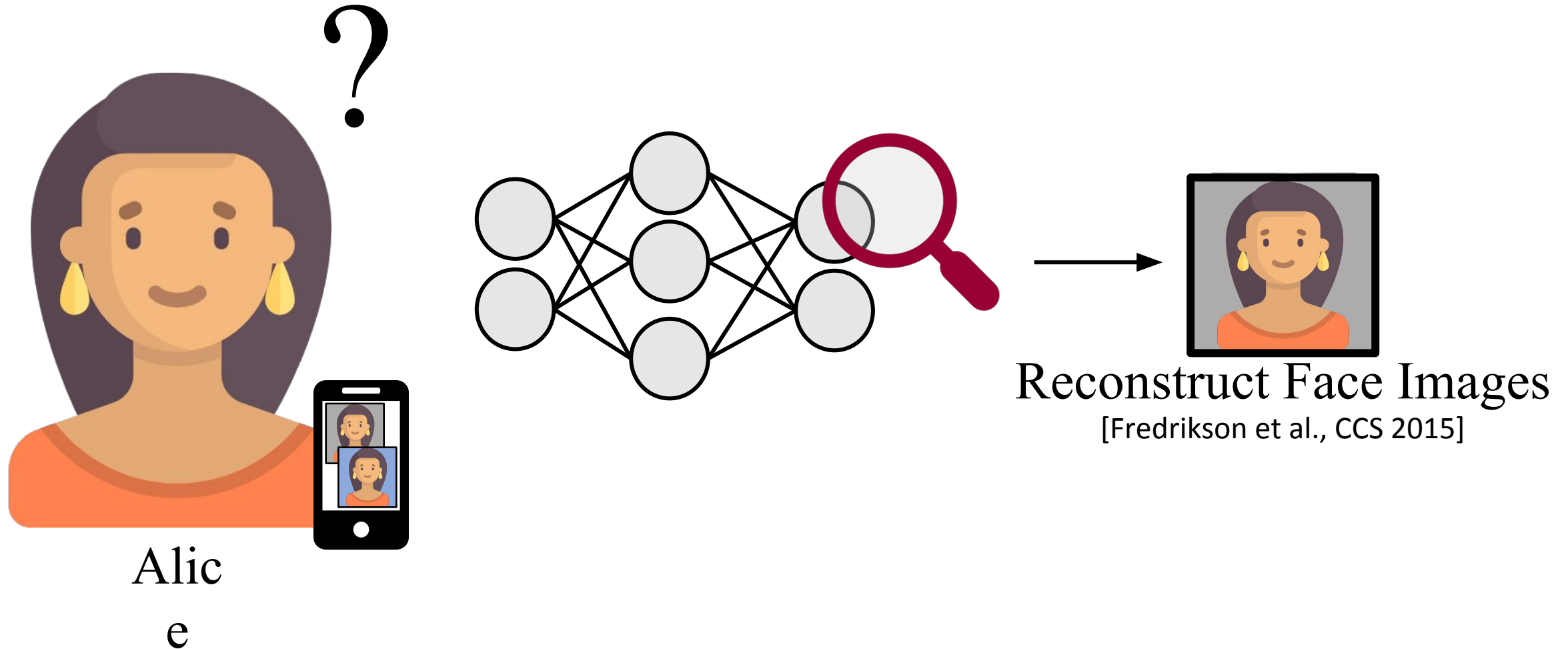INFORMATION SECURITY

# Individuals Generate Sensitive Data



Alice

# Companies apply Machine Learning
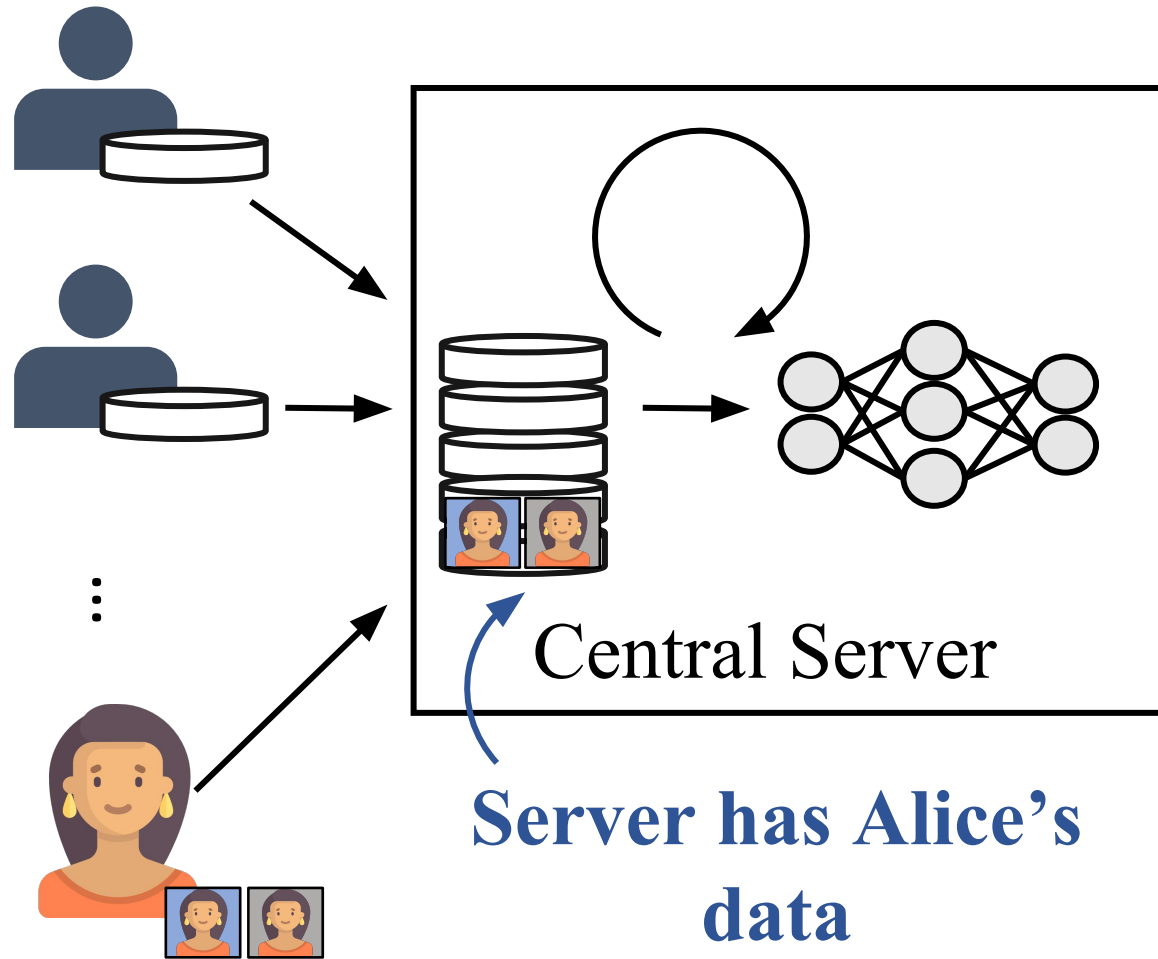


Alice

# ML Models Leak Private Information



Reconstruct Face Images
[Fredrikson et al., CCS 2015]

Alice

# Centralized vs. Federated Learning



Server has Alice's data

Centralized Learning

Gradients
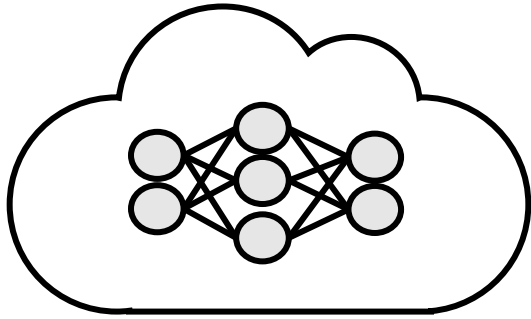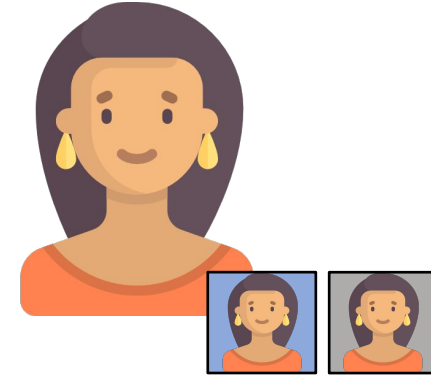
Central Server

Federated Learning

# Key Properties of Federated Learning



Central Server
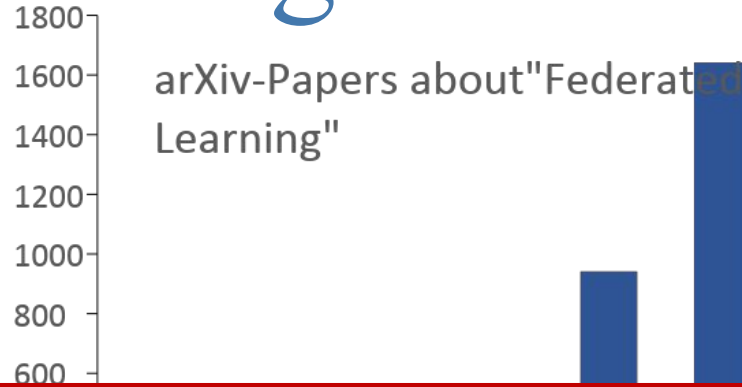
+ Heterogenous data
+ Efficient communication
+ Low costs

Individual User

- Performs compute
- Provides storage
+ Keeps data locally

**Privacy?!**
**?**

# Federated Learning is Extremely Popular

arXiv-Papers about "Federated Learning"

Federated Learning: A Game-Changer for Secure and Accurate AI in He...

Collaboration between Intel, Aster DM H... the launch of India's first-of-its-kind se... based health data platform

Authored by: TN Tech Desk

Features | October 14, 2022

## Can federated learning unlock AI in clinical trials without breaching privacy?

...ves brain tumour

Hi how are you!

Health Access

Health

In A New AI Research, Federated Learning Enables Big Data For Rare Cancer Boundary Detection

By Aneesh Tickoo - December 13, 2022

Reddit · Y · in · ✈

0 SHARES

# What Trust Model is Needed for Privacy?



Federated Learning

# What Trust Model is Needed for Privacy?



Federated Learning

# Federated Learning



Central Server

Shared Model

Calculate Gradients

**Sampled by server** M Users

# Federated Learning



Update Model

Central Server

Aggregation

Gradients

Gradients

Gradients

Gradients

Gradients

Gradients

M Users

# Alice's Privacy Relies purely on the Gradients



Central Server

Shared Model

**Should hide Alice's data**

Gradients

M Users

# Prior Data Reconstructions Attacks are Limited

We can reconstruct data…
… from different classes
… from small mini-batches
… that is of
… at high computational costs low complexity



We can extract data:
… from mini-batches of size = 1



[Zhu et al., NeurIPS 2019]

[Geiping et al., NeurIPS 2020]

# We Extract Large Amounts of Data Perfectly

Original Data



Extracted Data



… from all kinds of class distribution
… from large mini-batches with hundreds of data points
… with high complexity
… at near-zero computational costs

Franziska Boenisch, Adam Dziedzic, Roei Schuster, Ali Shahin Shamsabadi, Ilia Shumailov, Nicolas Papernot. *When the Curious Abandon Honesty: Federated Learning Is Not Private*, 2021. [IEEE Euro S&P '23a]
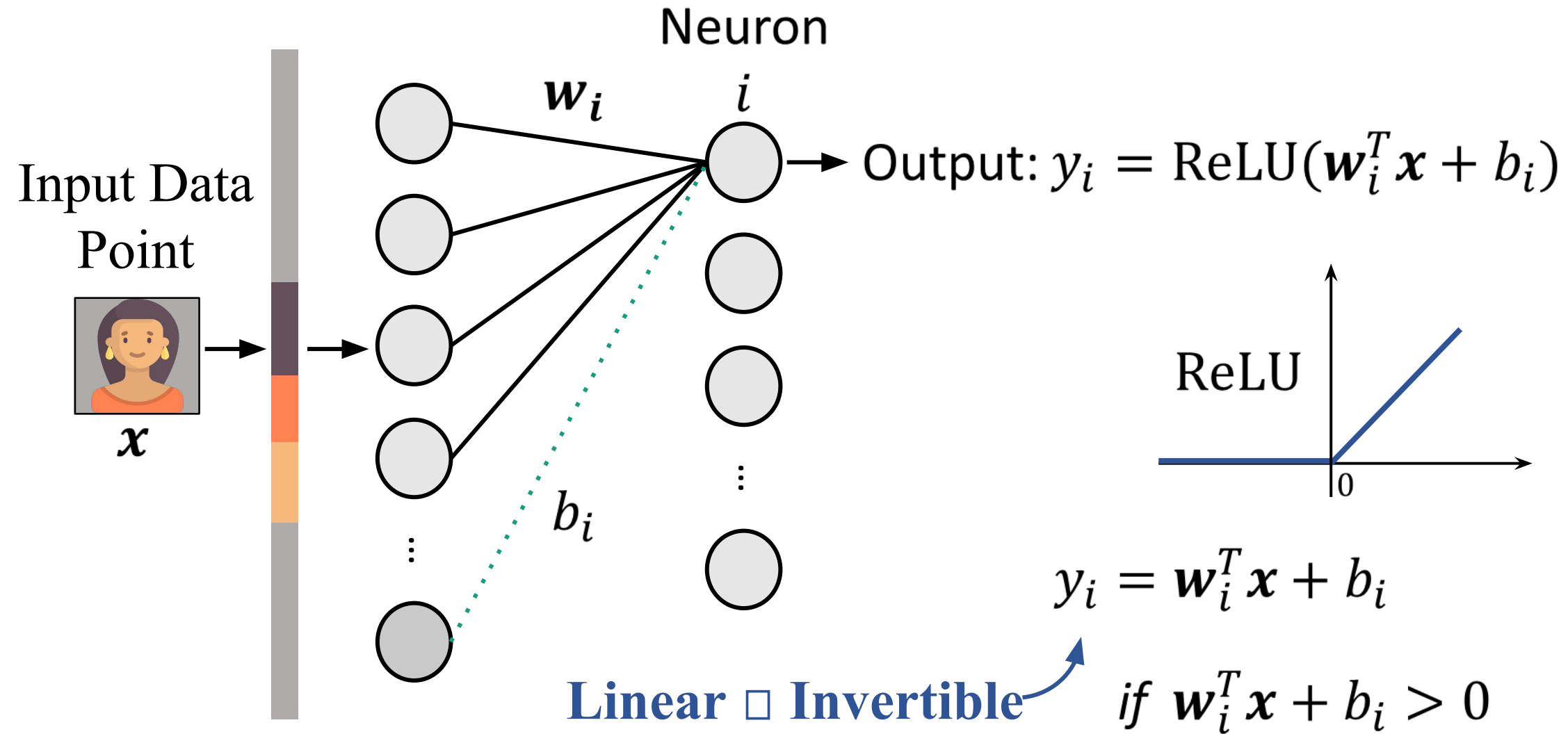
# Forward Pass through Fully-Connected Layer

Neuron

Input Data Point

$\boldsymbol{w_i}$

$i$

Output: $y_i = \mathrm{ReLU}(\boldsymbol{w}_i^T \boldsymbol{x} + b_i)$

$\boldsymbol{x}$

$b_i$

ReLU

$y_i = \boldsymbol{w}_i^T \boldsymbol{x} + b_i$

**Linear □ Invertible**

$if \; \boldsymbol{w}_i^T \boldsymbol{x} + b_i > 0$

# Prior Extraction Works only for Single Data Points

$$\rightarrow \frac{\partial \mathcal{L}}{\partial \boldsymbol{w}_i^T} = \frac{\partial \mathcal{L}}{\partial b_i} \boldsymbol{x}$$

$$y_i = \boldsymbol{w}_i^T \boldsymbol{x} + b_i$$

$$\frac{\partial y_i}{\partial \boldsymbol{w}_i^T} = \boldsymbol{x}$$

$$\frac{\partial y_i}{\partial b_i} = 1$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{w}_i^T}$$

$$\frac{\partial \mathcal{L}}{\partial b_i}$$

$$=$$

**Contains scaled input data point**

**Contains scaling factor**
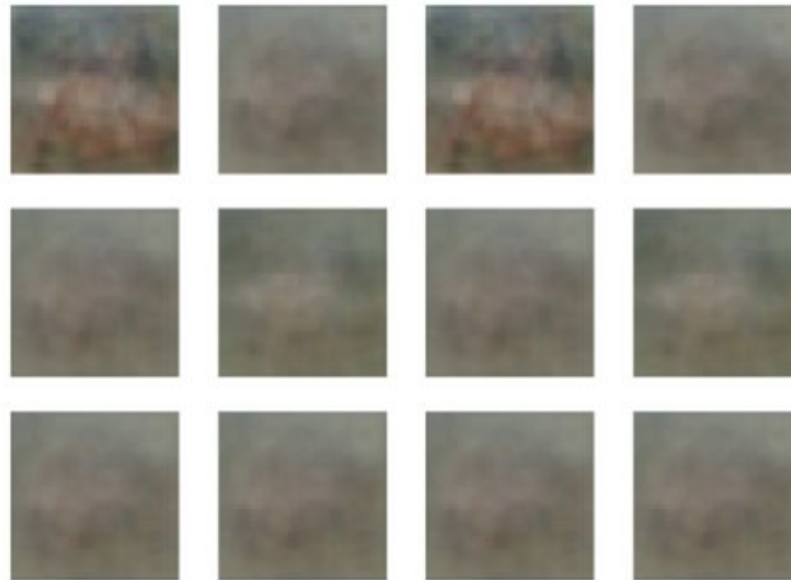
$\boldsymbol{x}$

[Geiping et al., NeurIPS 2020]

# Extraction for Large Mini-Batches Should Fail

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{w}_i^T} = \sum_{j=1}^{B} \frac{\partial \mathcal{L}}{\partial y_{i,j}} \frac{\partial y_{i,j}}{\partial \boldsymbol{w}_i^T}$$

**Mini-batch gradient**



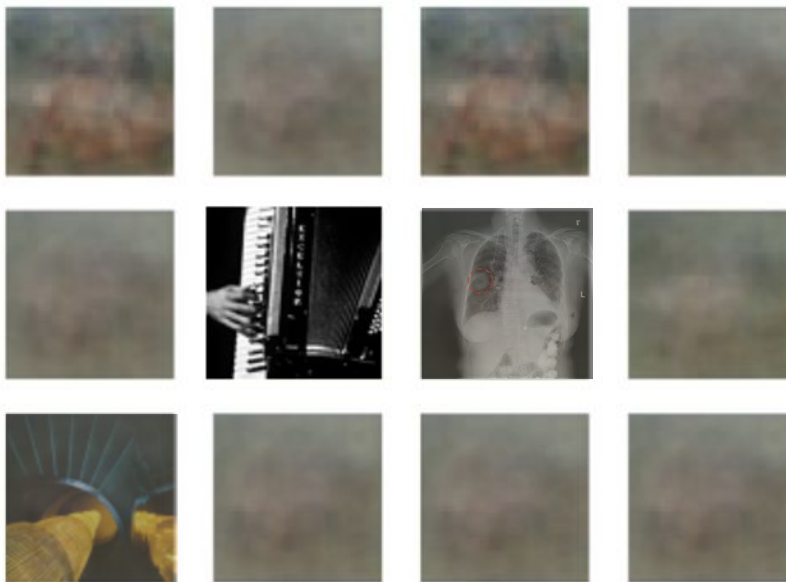**We believe rescaled gradients look like this….**

# Data Leaks Directly from Model Gradients

```
weights_gradient = gradients[0].numpy()
inverse_bias = 1 / gradients[1].numpy()
extracted_data = inverse_bias * weights_gradient
plot(extracted_data, num_rows = 3, num_cols = 6)
```

$$x = \left(\frac{\partial \mathcal{L}}{\partial b_i}\right)^{-1} \frac{\partial \mathcal{L}}{\partial w_i}$$
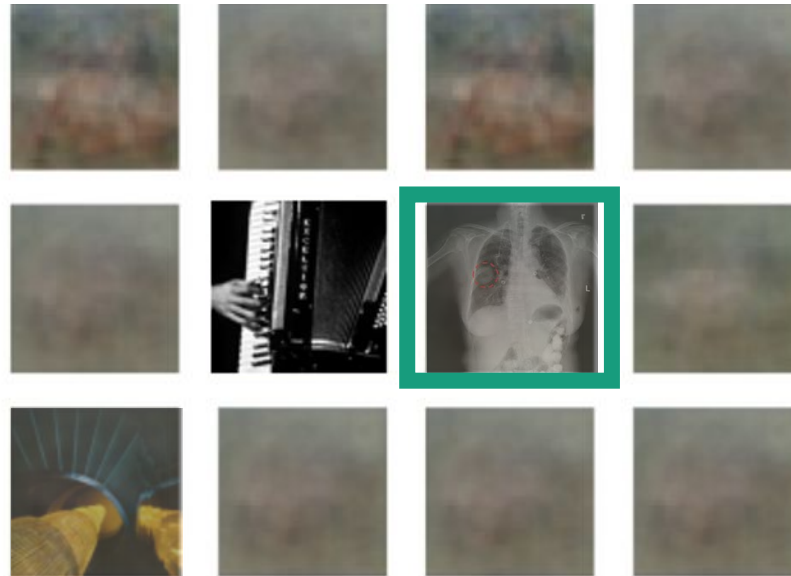
All you need is matplotlib

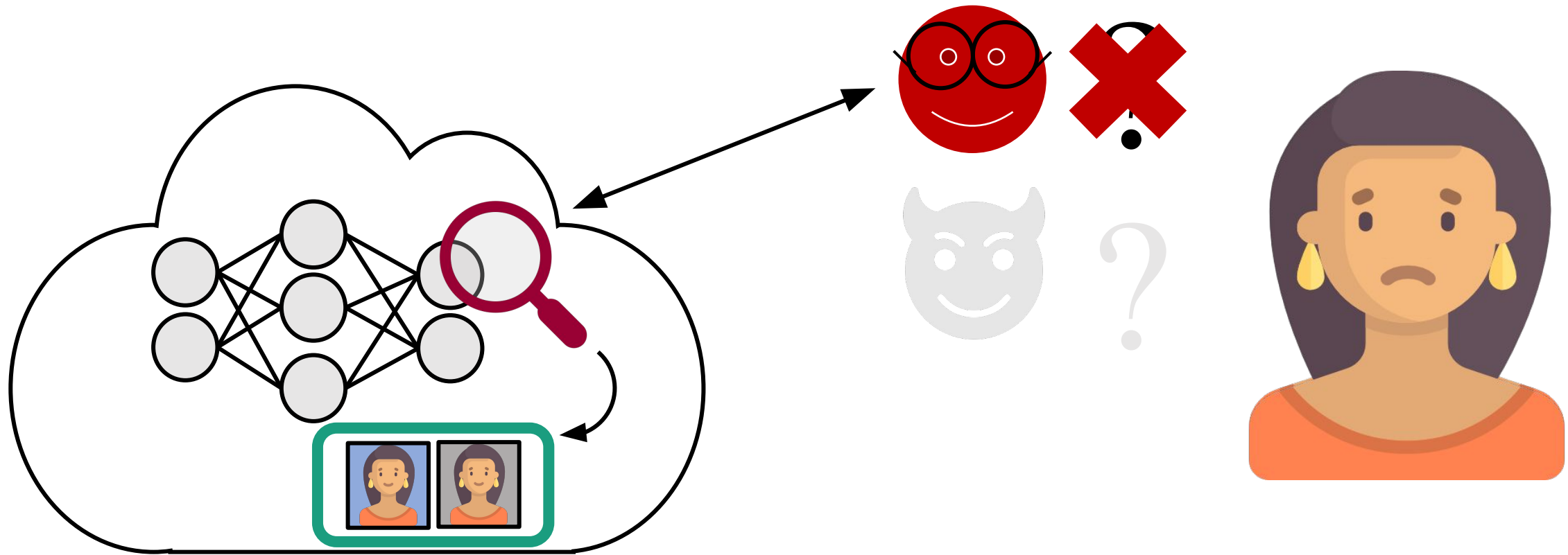**… but they actually look like that!**

mini-batch size=100

# Gradients can Leak Single Data Points

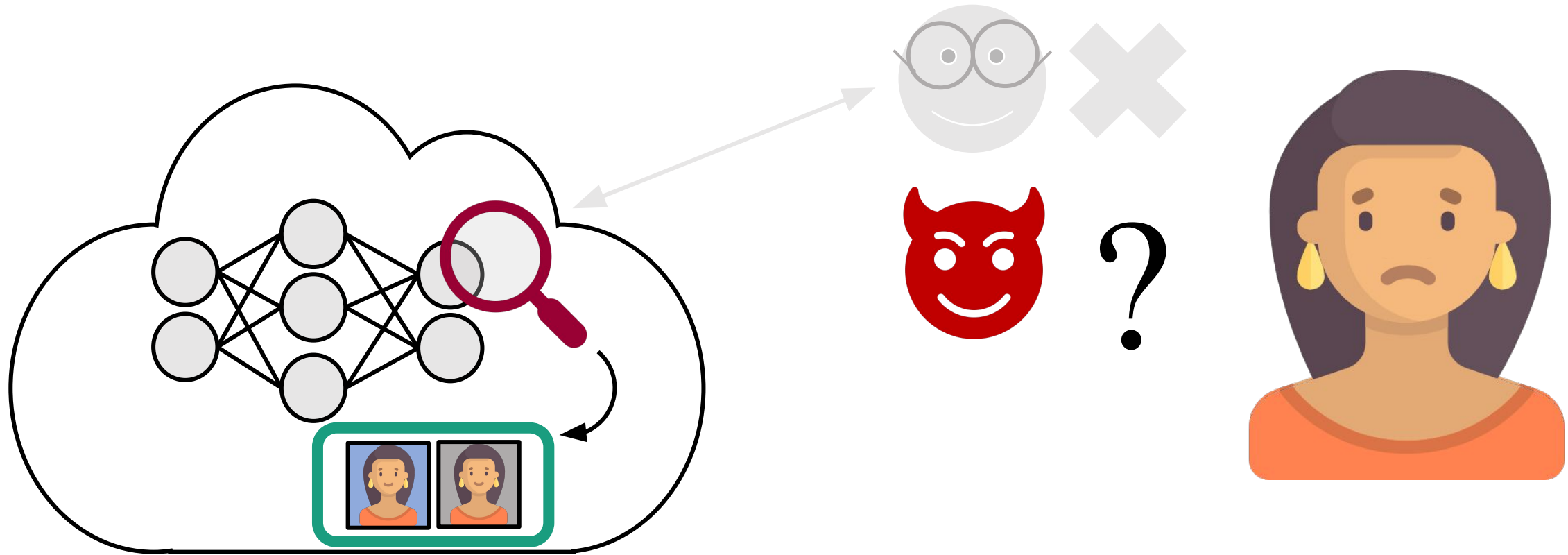Why can we still extract individual data points $\boldsymbol{x}$?



**Gradient of a single data point**

# What Trust Model is Needed for Privacy?

Even a passive, honest-but-curious attacker can extract
a significant amount of sensitive user-data.
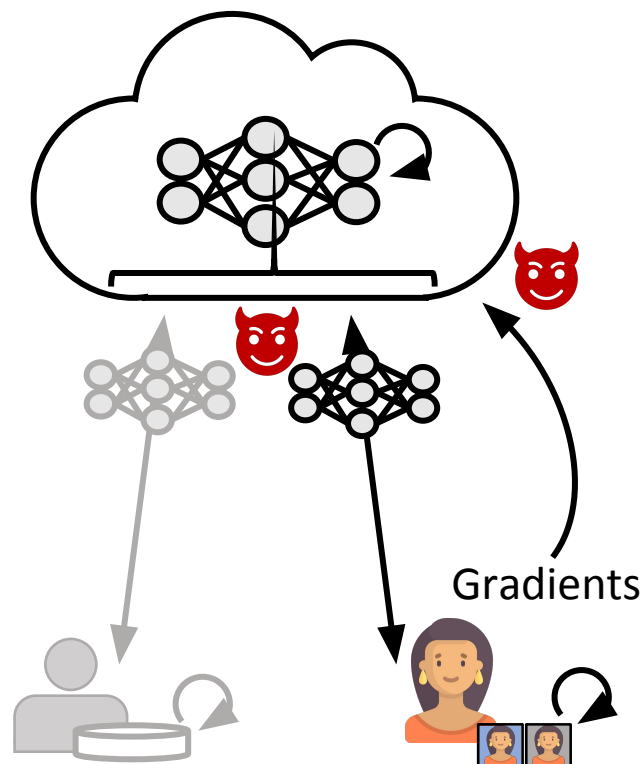
# What Trust Model is Needed for Privacy?



Even a passive, honest-but-curious attacker can extract a significant amount of sensitive user-data.

# Our Trap Weights Increase Natural Leakage

**Trap Weights:** Induce $x^T w_i + b_i \leq 0$ for most input data points $x$

**Makes other points extractable**

1) Initialize model weights at random
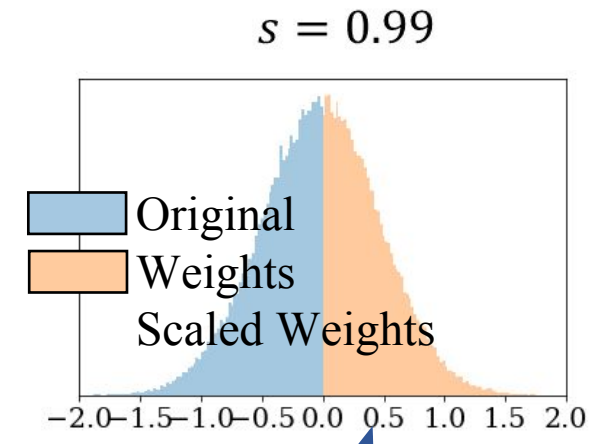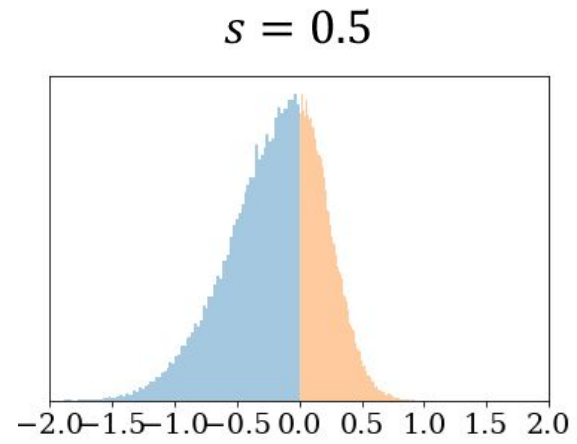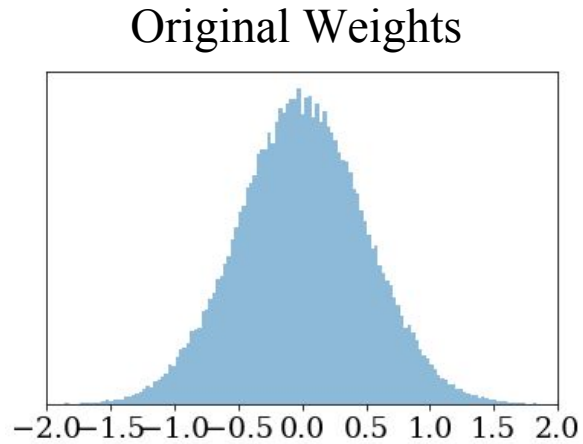2) Scale positive components down by $s < 1$

$\rightarrow (x^T s w_i^+) + (x^T w_i^-) + b_i \leq 0$ more often

Gradients

Assumes input features $x$ in range $[0, 1]$

**Standard pre-processing**

# Influence of Scaling Factor "s"



Original Weights

$s = 0.5$

$s = 0.99$

Original Weights
Scaled Weights

**Inconspicuous**

**Active Extraction**

**Baseline: Passive Extraction**

| Scaling Factor (s) | Activated Neurons (by 1 data point) (%) ❗ | Extracted Data (%) |
|---|---|---|
| 0.4 | 0 | 0 |
| 0.5 | 0 | 0 |
| 0.9 | 0 | 0 |
| 0.99 | 65.5 (51.4) | 45.7 |
| 1.0 | 99.9 (4.4) | 21.8 |

ImageNet Extraction: Mini-Batch Size = 100, 1000 Neurons

23

# Our Trap Weights Improve Extraction

|  | Passive | Active |
|---|---|---|
| MNIST | 5.8 | 54 |
| CIFAR10 | 25.5 | 54 |
| ImageNet | 21.8 | 45.7 |
| IMDB | 25.4 | 65.4 |

Extracted Data (%),
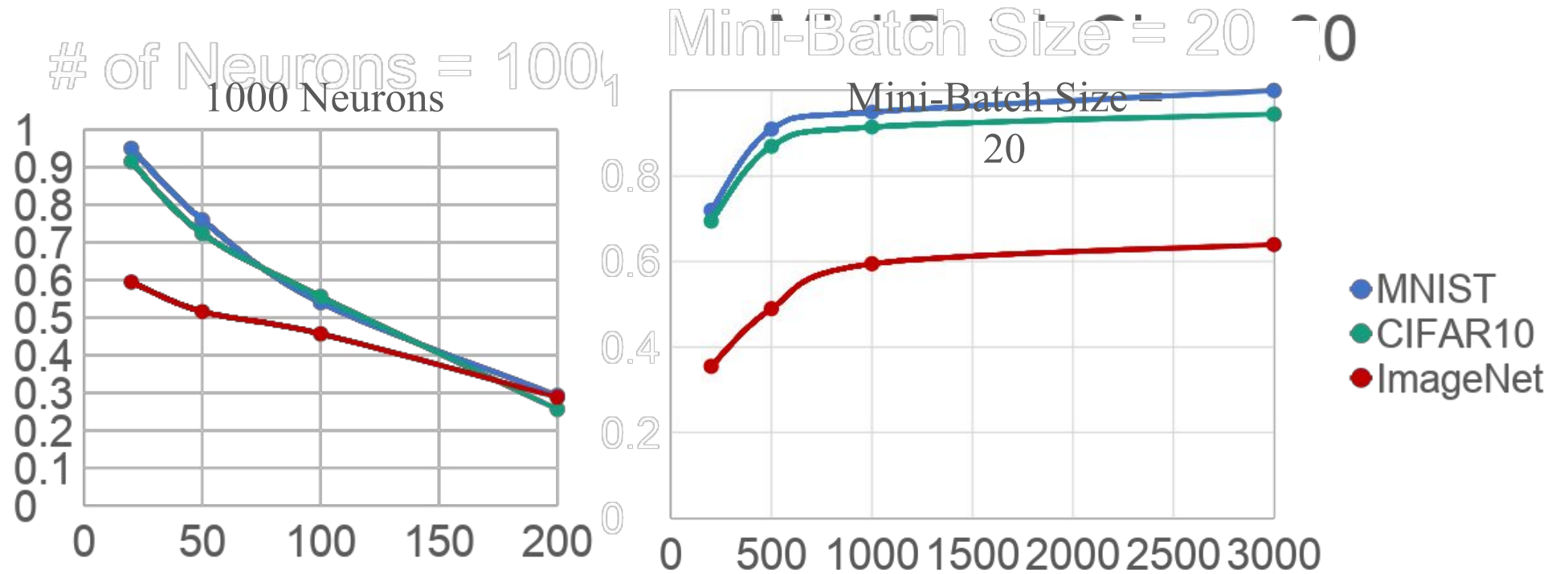Mini-Batch Size = 100,
1000 Neurons



CIFAR10
(Non-IID)
**Extracted from gradients within < 1 second**

# More Neurons and Smaller Mini-Batches Let us Extract More Data



Extraction Recall
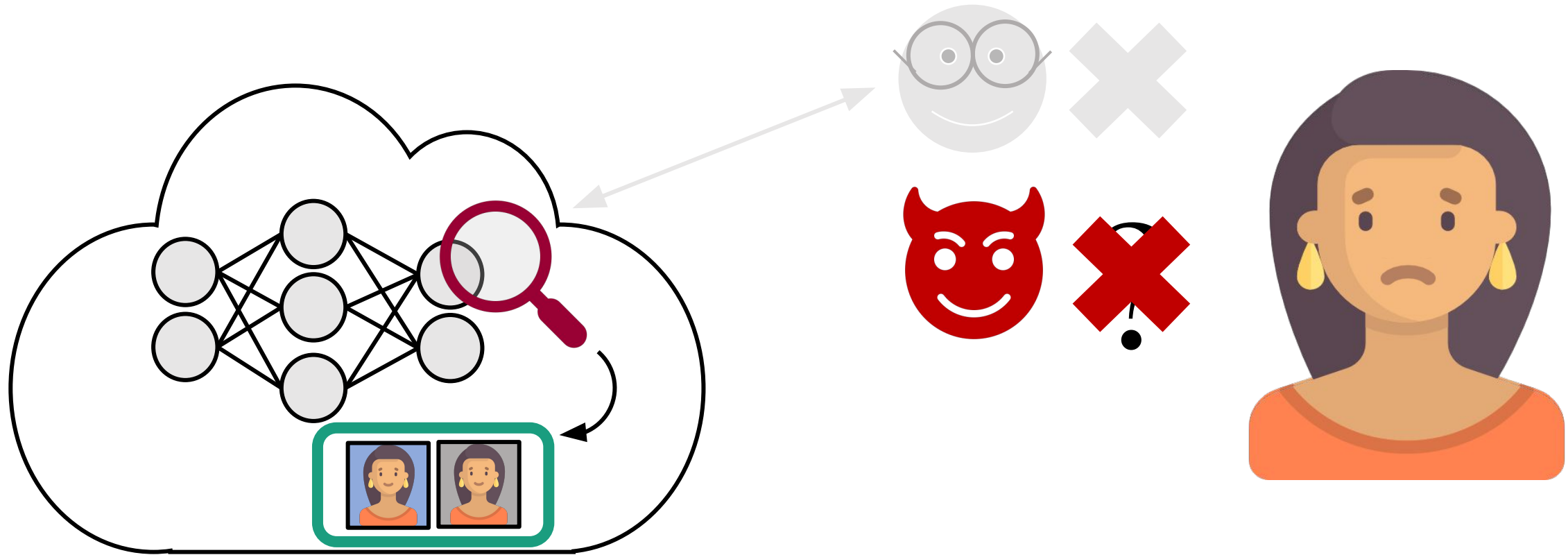
Mini-Batch Size

1000 Neurons

Mini-Batch Size = 20

MNIST
CIFAR10
ImageNet

# of Neurons

**Specified by the server**

# What Trust Model is Needed for Privacy?



An active, malicious attacker can significantly increase privacy risks for users.

# Conclusion for Privacy in FL



Participate only in Protocols
with Trusted Server

Replace Trust by
Verifiable
Mechanisms

Franziska Boenisch, Adam Dziedzic, Roei Schuster, Ali Shahin Shamsabadi, Ilia Shumailov, Nicolas Papernot. *Is Federated Learning a Practical PET Yet?*, 2023. [IEEE Euro S&P '23a]

# Backup Slides

# Defending FL is Complex and Costly



User Sampling

Model Initialization

ZK ✔

Gradient Calculation and Aggregation

Computational Costs

Noise Addition

# Power Imbalance Makes FL Vulnerable

Server wants
Utility

User Provisioning
& Sampling

Model
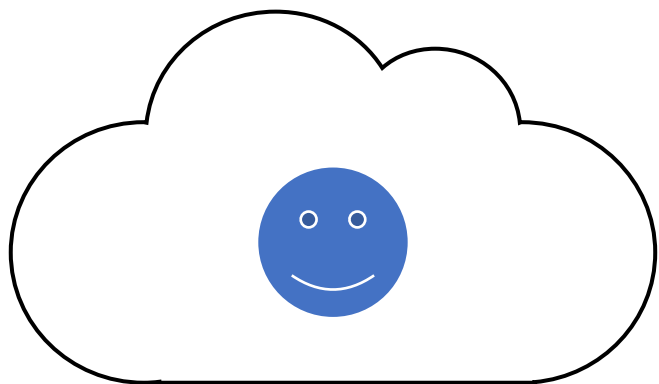Manipulations

Users need
Privacy

Unknown
Collaborators

Unverified shared model
and computations

# What Trust Model is Needed for Privacy?

An active, malicious attacker can significantly increase privacy risks for users.

# Differential Privacy Protects Individual Data



$$\frac{\text{Pr} \left( \begin{array}{ccc} & & \end{array} \right) \square}{\text{Pr} \left( \begin{array}{ccc} & & \end{array} \right) \square} \leq e^{\varepsilon}$$

Pr (Train

Pr (Train

**Adjacent datasets** $\mathcal{C}$

$\sigma$

(1) Clip Gradients

(2) Noise Gradients

# Differential Privacy in Federated Learning



Central DP: Server adds noise

Distributed DP: Users add noise

**After aggregation**

Local DP: Users add noise

$$\mathcal{N}\left(0, \sigma^2 c^2 \frac{\sigma^2}{(M-1)}\right)$$

Noised Clipped Gradients

# Aggregate via Secure Aggregation

Global Noise
$$\mathcal{N}(0, \sigma^2 c^2)$$



Release
Aggregate

Gradient    Gradient

**Alice's data seems protected**

Overhead:
- Computation
- Communication
- Storage
- Availability of PKI

Local Noise: $\mathcal{N}\left(0, \frac{\sigma^2}{(M-1)} c^2\right)$

[Bonawitz et al., CCS 2017]

# Attacking FL protected by DDP+SA



Sybil Users

Legitimate Users

Server

$0 \longrightarrow 0$

$G_{Alice} \longrightarrow \bar{G}_{Alice}$

$\bar{G}_{Alice}$

$0 \longrightarrow 0$

Franziska Boenisch, Adam Dziedzic, Roei Schuster, Ali Shahin Shamsabadi, Ilia Shumailov, Nicolas Papernot. *Is Federated Learning a Practical PET Yet?*, 2023. [IEEE Euro S&P '23a]

# DDP Reduces to LDP with Low Privacy Levels

Test Acc.

User



$\varepsilon = \text{inf}$  $\varepsilon = 2e4$  $\varepsilon = 592$  $\varepsilon = 33.97$  $\varepsilon = 5.41$  $\varepsilon = 2.39$

… believes to $\mathcal{N}(0, \sigma^2 c^2)$

**Not private enough**

**Too little utility**

# What Trust Model is Needed for Privacy?



Even in hardened variants of the protocol, a malicious attacker can breach individual users' privacy.

# My Research

Goal: Develop mechanisms that provide individualized notions of privacy for machine learning



| | | | |
|---|---|---|---|
| | **Federated Learning** | Data Extraction in Federated Learning | EuroS&P'23a |
| | | Reconstruction in Hardened Protocols | EuroS&P'23b |
| | **Individualized Privacy** | Individualized Privacy with PATE | PoPETs'23a |
| | | Individually Private SGD Training | Submission'24 |
| | **Privacy Auditing** | Private Prompt Tuning to Query LLMs | Submission'24 |
| | | Side-Channels in Private Query Systems | CCS'21 |
| | | | SPSC'22 |
| | | Model Inversion in Speaker Recognition | PoPETs'23b |

# My Research

Goal: Develop mechanisms that provide individualized notions of privacy for machine learning

**Individualized Privacy**

Individualized Differential Privacy — PoPETs'23 a

GDPR-Aligned Privacy Assessment — PoPETs'23 b

**Privacy Auditing**

Side-Channels in Private Query Systems — CCS'21

Model Inversion in Speaker Recognition — SPSC'22

Bounding Membership Inference — arXiv'2

**Federated Learning**

Data Extraction in Federated Learning — EuroS&P'23 a

Reconstruction in Hardened Protocols — EuroS&P'23

# My Research

Goal: Develop mechanisms that provide individualized notions of privacy for machine learning



**Individualized Privacy**

Individualized Differential Privacy — PoPETs'23a

GDPR-Aligned Privacy Assessment — PoPETs'23b

**Privacy Auditing**

Side-Channels in Private Query Systems — CCS'21

Model Inversion in Speaker Recognition — SPSC'22

**Federated Learning**

Bounding Membership Inference — arXiv'2

Data Extraction in Federated Learning — EuroS&P'23a

Reconstruction in Hardened Protocols — EuroS&P'23

# Side-Channel Attacks against Query Systems



SQ
L

New
SQL

Anon. Answer
or Error

**Chorus (Uber)**

**Diffix**

**Analys**

**Anonymizing Proxy**

**Private Database**

IF     *Name='Alice'*
AND   *Disease='Cancer'*

THEN *SQRT(age − 1000)*

Franziska Boenisch, Reinhard Munz, Marcel Tiepelt, Simon Hanisch, Christiane Kuhn, and Paul Francis. *Side-channel attacks on query-based data anonymization.* In Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security, 2021. [CCS'21]

# FL Sources

# Mitigation Methods

# Differential Privacy

Goal: produce statistically indistinguishable outputs on any pair of datasets that only differ by any single data point.

**Differential Privacy**: a randomized mechanism $M$ with domain $D$ and range $R$ satisfies $(\varepsilon, \delta)$-differential privacy if for any subset $S \subseteq R$ and any adjacent datasets $d, d' \in D$, i.e., $\|d - d'\|_1 \leq 1$, the following inequality holds:

$$\Pr[M(d) \in S] \leq e^{\varepsilon} \Pr[M(d') \in S] + \delta$$

# Secure Multi Party Computation (MPC)

**Setup:** given participants $p_1, p_2, p_3$ and their private data $x_1, x_2, x_3$.

**Task:** compute value of a private function $F(x_1, x_2, x_3)$.

**Example:** compute the maximum or average salary of the participants, without revealing the individual salaries.

**Machine Learning:** shareholders can compute **any function** of inputs without seeing anything but shares and the final output.**Properties:**
(1) input privacy – no information about private data can be inferred from messages exchanged during MPC, and

(2) honest parties either compute correct output or abort.

# Secure Multi-Party Computation (MPC)



Alice
100K

Bob
200K

Carol
300K

# Secure Multi-Party Computation (MPC)

Alice
100K

Bob
200K

Carol
300K

Generate
Random
Shares

| Alice | Bob | Carol |
|-------|------|-------|
| 50 | -80 | 0 |
| 30 | 100 | 350 |
| 20 | 180 | -50 |

# Secure Multi-Party Computation (MPC)

# Secure Multi-Party Computation (MPC)

| | **A**lice | **B**ob | **C**arol |
|---|---|---|---|
| | 100K | 200K | 300K |

| | Alice | Bob | Carol |
|---|---|---|---|
| **Generate Random Shares** | 50 | -80 | 0 |
| | 30 | 100 | 350 |
| | 20 | 180 | -50 |
| **Secret Sharing** | 50 | 30 | 20 |
| | -80 | 100 | 180 |
| | 0 | 350 | -50 |
| **Add Shares** | -30 | 480 | 150 |

Sum: 600     Mean: 200

# Homomorphic Encryption

1. Addition

$$Enc(x) + Enc(y) = Enc(x + y)$$
$$Enc(x) + y^* = Enc(x + y)$$

2. Multiplication

$$Enc(x) * Enc(y) =$$
$$x^e \bmod n * y^e \bmod n =$$
$$x^e \, y^e \bmod n =$$
$$(xy)^e \bmod n =$$
$$Enc(x * y)$$

# Attacker Models

- Honest-but-curious – adversary follows the protocol but tries to infer information from the protocol transcript.

- Malicious – adversary actively deviates from the protocol

- Occasionally Byzantine – adversary acts honest most of the time and only acts maliciously on occasions

# Secure Aggregation



- Robustness **(Malicious Server)**
  - Can collaborate with up to $n/3-1$ clients
  - Tolerates up to $n/3-1$ dropouts of clients

Bonawitz, Keith, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H. Brendan McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. "Practical secure aggregation for privacy-preserving machine learning." In *proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1175-1191. 2017.

# Secure Aggregation

## Bonawitz et al., 2017

- Computation:
  - User : $O(n^2 + mn)$
  - Server : $O(mn^2)$

- Communication:
  - User : $O(n + m)$
  - Server : $O(n^2 + mn)$

- Storage:
  - User : $O(n + m)$
  - Server : $O(n^2 + m)$

## Bell et al., 2020

- Computation:
  - User : $O(\log^2 n + \log n)$
  - Server : $O(n(\log^2 n + \log n))$

- Communication:
  - User : $O(\log n + m)$
  - Server : $O(n(\log n + m))$

Bonawitz, Keith, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H. Brendan McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. "Practical secure aggregation for privacy-preserving machine learning." In *proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1175-1191.

Bell, James Henry, Kallista A. Bonawitz, Adrià Gascón, Tancrède Lepoint, and Mariana Raykova. "Secure single-server aggregation with (poly) logarithmic overhead." In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1253-1269. 2020.

# Distributed Differential Privacy

**Distributed Discrete Gaussian**

- discretizes the data and adds discrete Gaussian noise before performing secure aggregation
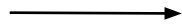
**Skellam Mechanism**

- based on the difference of two independent Poisson random variables

Kairouz, Peter, Ziyu Liu, and Thomas Steinke. "The distributed discrete gaussian mechanism for federated learning with secure aggregation." In *International Conference on Machine Learning*, pp. 5201-5212. PMLR, 2021.

Agarwal, Naman, Peter Kairouz, and Ziyu Liu. "The skellam mechanism for differentially private federated learning." *Advances in Neural Information Processing Systems* 34 (2021): 5052-5064.

# Forwarding over Convolutional Layers

| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

Input

| 0 | 1 | 0 |
| 0 | 0 | 0 |

Filter

| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

Feature Maps

Input

Filters

Feature Maps

Filter 1

Filter 2

Filter 3

Input

| 5 | 6 | 7 | 8 |
|----|----|----|----|
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

Feature Maps

Filter 1

Filter 2

Filter 3

| 0 | 1 | 0 |
| 0 | 0 | 0 |

Filter Input

| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

Feature Maps

| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

Filter 1

Filter 2

Filter 3

Filter 4

Input

Feature Maps

| 1 | 3 |
|---|---|
| 9 | 2 |

4

| 10 | 5 | 7 |
|----|---|---|

| 13 | 6 | 8 |
|----|---|---|

| 14 | 16 |
|----|----|

Filter 1

Filter 2

Filter 3

Filter 4

Input

Feature Maps

# Forwarding over
# Fully Connected Layers

Feat.1 1 =1

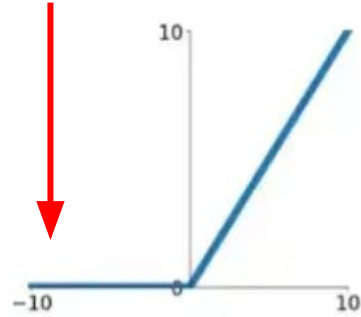Feat.2 2

Noise

=0

1

2

1

2

⋯

Class 1

⋮

Class k

# Other Activation Functions

**ReLU**

$\max(0, x)$
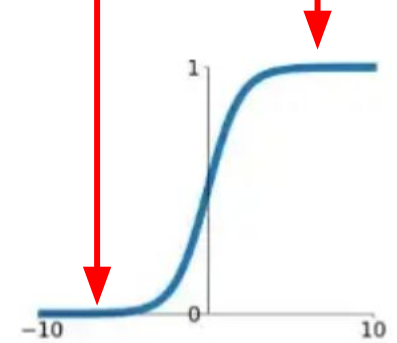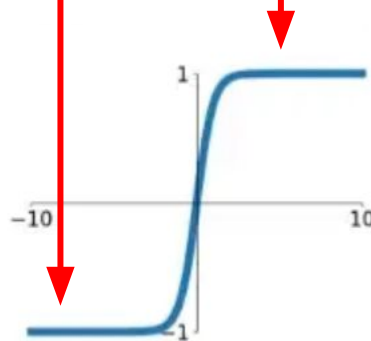
Zero-Gradients

**Sigmoid**
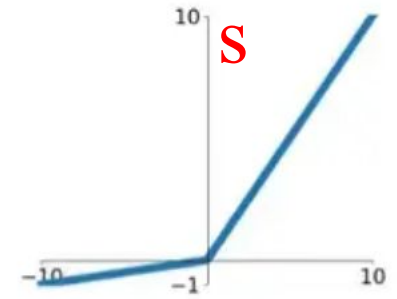
$\sigma(x) = \frac{1}{1+e^{-x}}$

Zero-Gradients

**tanh**

$\tanh(x)$

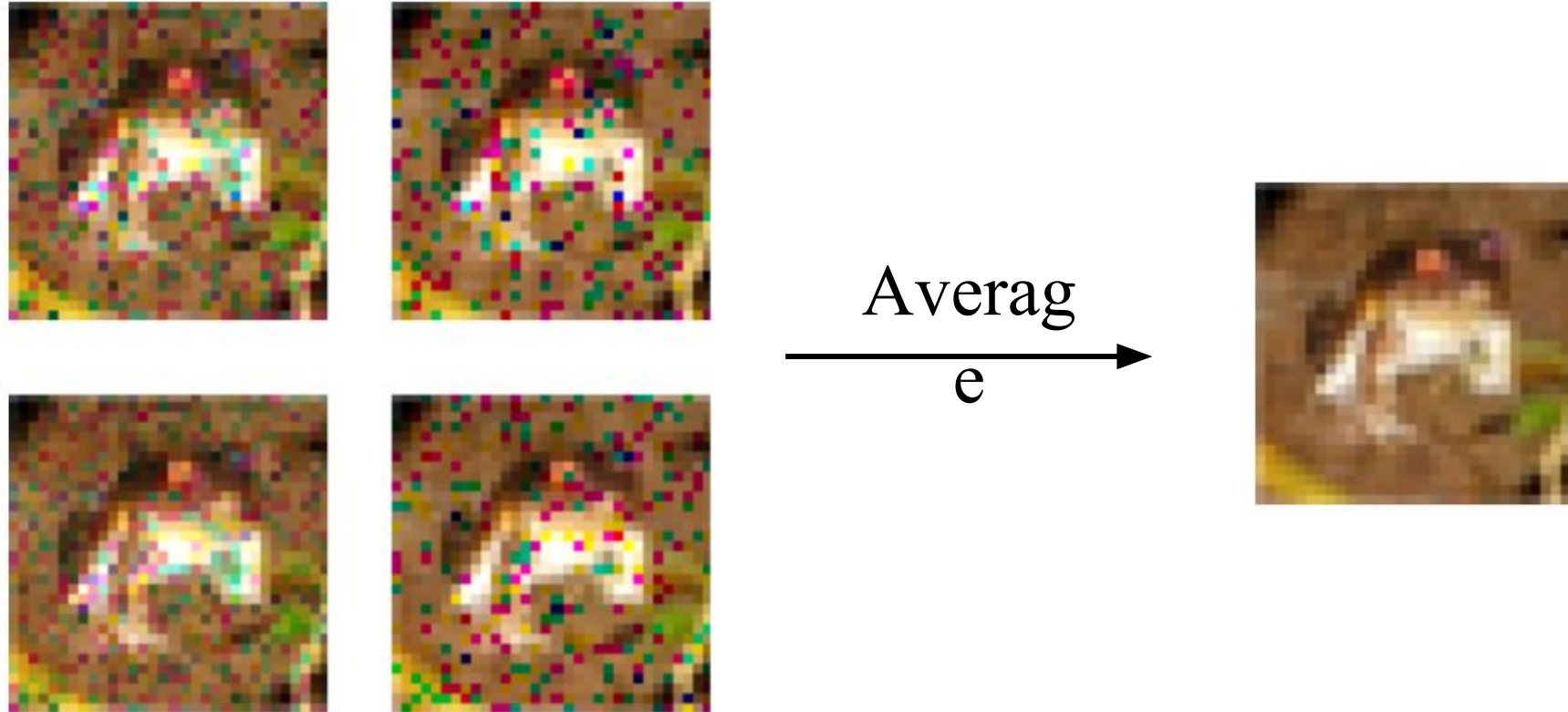Zero-Gradients

**Leaky ReLU**

$\max(0.1x, x)$

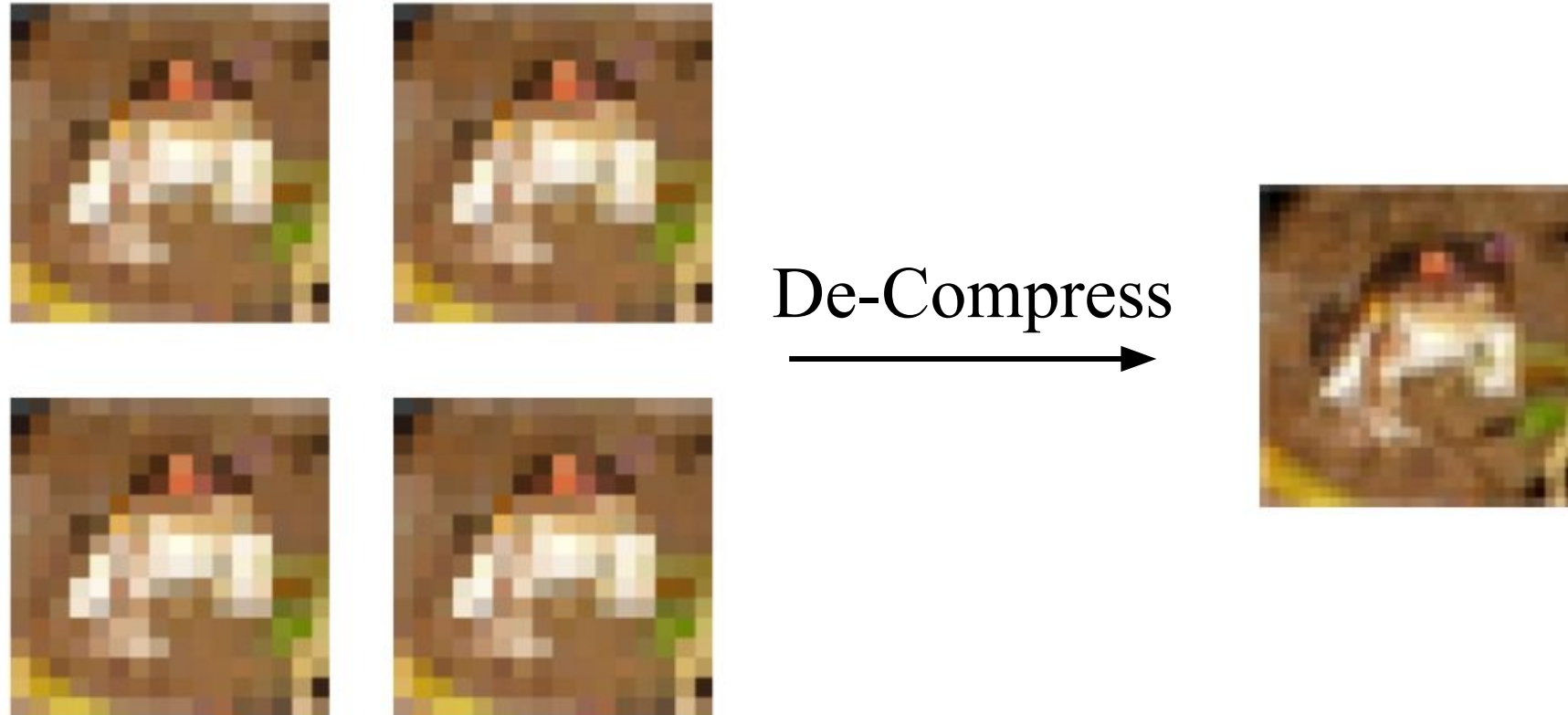Non-Zero-Gradient

s

… but less sparsity.

# Lossy Architecture

Conv(f=32, k=(3,3), s=1, p=same, act=relu)
MaxPool()
Conv(f=64, k=(3,3), s=1, p=same, act=relu)
Dropout()
Flatten
Dense(n=1000, act=relu)
Dropout()
Dense(n=#classes, act=None)

# Effect of Dropout



Average

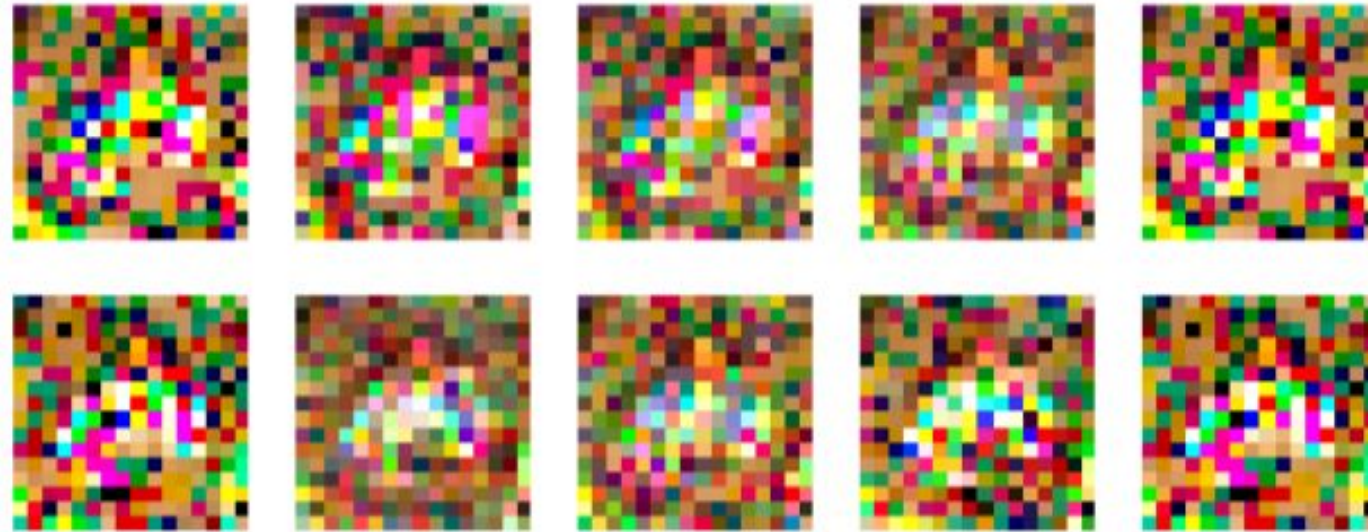Same Data Point Extracted at 4 Different Gradients (Dropout Rate = 0.1)

# Effect of Pooling



De-Compress

Same Data Point Extracted at 4 Different Gradients (Max Pooling with 2x2)

# Heavy Dropout and Pooling



(c) Dropout with $p = 0.3$ and pooling.

# Individual Activation Neurons

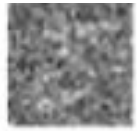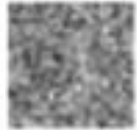| | % Individual Act. | |
|---|---|---|
| | Passive | **Active** |
| MNIST | 0.6% | **20.3%** |
| CIFAR10 | 5.8% | **41.2%** |
| ImageNet | 4.4% | **51.4%** |
| IMDB | 3.6% | **19.2%** |

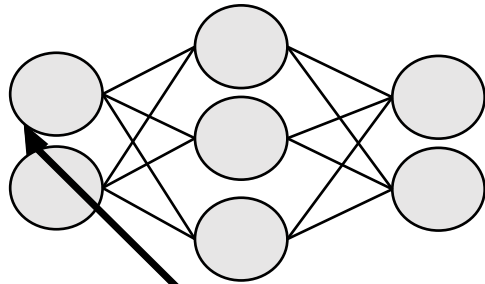# Extractable Datapoints

# Extractable Datapoints

Related Work (FL)

# Optimization-based Data Reconstruction

"Gradient-Matching"



$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{w}_i} \frac{\partial \mathcal{L}}{\partial b_i}$$

$$\frac{\partial \mathcal{L}'}{\partial \boldsymbol{w}_i} \frac{\partial \mathcal{L}'}{\partial b_i}$$

**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
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: $(\mathbf{x}_i^*, y_i^*) \leftarrow (\hat{\mathbf{x}}^{[\hat{T}+1]}, \hat{y}^{[\hat{T}+1]})$

[1] Zhu, Ligeng, Zhijian Liu, and Song Han. "Deep leakage from gradients." *Advances in neural information processing systems* 32 (2019).

# Limitations and Summary of Passive Attackers

- Computationally expensive

- Low fidelity

- Non-complex data
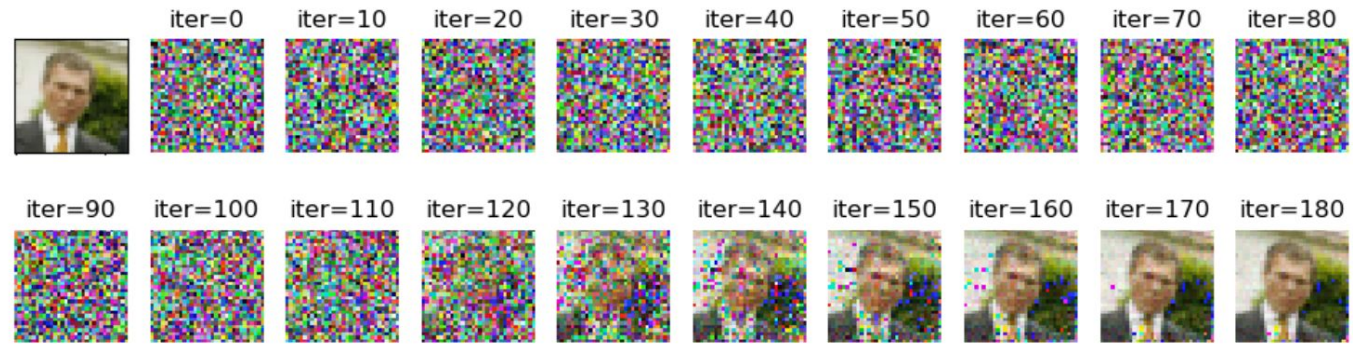
- Small mini-batch sizes, different classes



Figure taken from [2].

Even a *passive attacker* in vanilla FL
can reconstruct private user data.

[2] Zhao, Bo, Konda Reddy Mopuri, and Hakan Bilen. "idlg: Improved deep leakage from gradients."
*arXiv preprint arXiv:2001.02610* (2020).

# Imprinting User Data in Model Gradients

- Observation: bias term controls if a data point activates a neuron

$$y_i = ReLU\left(\boldsymbol{w}_i^T \boldsymbol{x} + b_i\right)$$
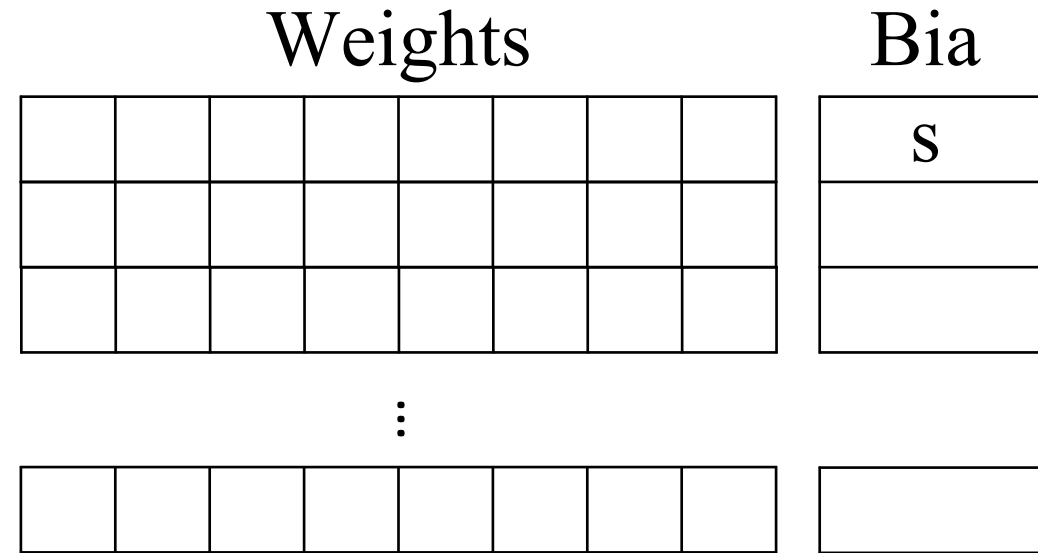$$y_i = 0 \ \text{ if } b_i < 0 \text{ and } |b_i| > |\boldsymbol{w}_i^T \boldsymbol{x}|$$

- Approach:
  - Control which data points activate what neurons
  - Turn model weights into linear function (e.g. average pixel brightness: $\frac{1}{m}$, $m$: number of features)
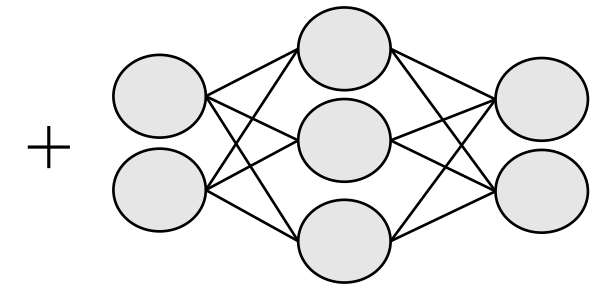  - Iteratively extract data

[Fowl et al., 2021, ICLR]

# Imprinting User Data in Model Gradients

Weights

Bia

$\gamma > 0$

s

$\vdots$

+

Mini-Batch

Imprint Module
(Fully Connected
Layer)

[Fowl et al., 2021, ICLR]

# Imprinting User Data in Model Gradients



Weights

Bia

s

Mini-Batc
h

Imprint Module
(Fully Connected
Layer)

# Data Extraction Success
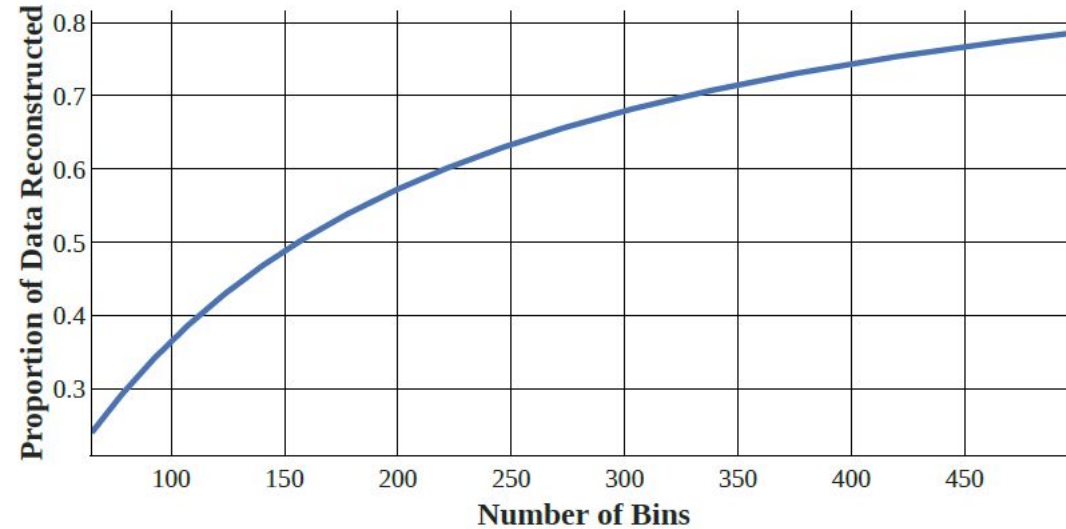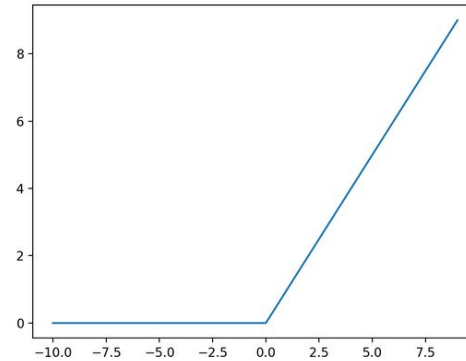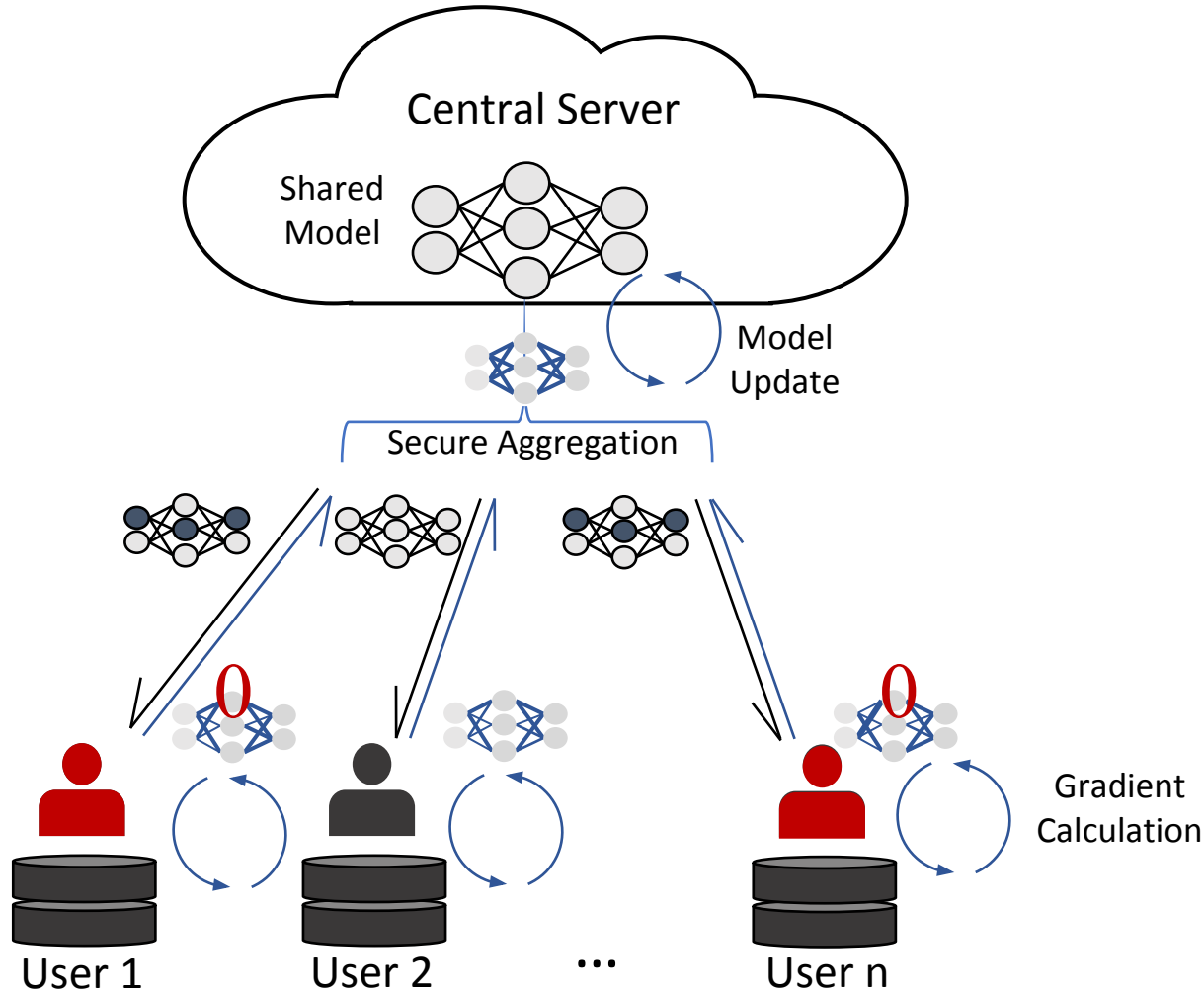


Figure taken from [4]. Results for mini-batch size of 64.

Extraction success increases with increasing the number of bins

[Fowl et al., 2021, ICLR]

# Eluding Secure Aggregation



ReLU

Model Inconsistency

Gradient Suppression

$$ReLU(w_i^T x + b_i) = 0$$

[Pasquini et al., 2022, CCS]

# Summary of my Contributions

1. Even with large mini-batches of high-dimensional data, significant proportions of private user data can be leaked to a passive attacker.

2. Active attackers can amplify this leakage even without performing highly noticeable changes to the model architecture / parameters.

3. Prior work has still largely underestimated the privacy risk of (hardened) FL.