

SPEAR: Exact Gradient Inversion of Batches in Federated Learning

Code: <https://github.com/eth-sri/SPEAR>



Dimitar I. Dimitrov



Maximilian Baader



Mark Müller



Martin Vechev

INSAIT

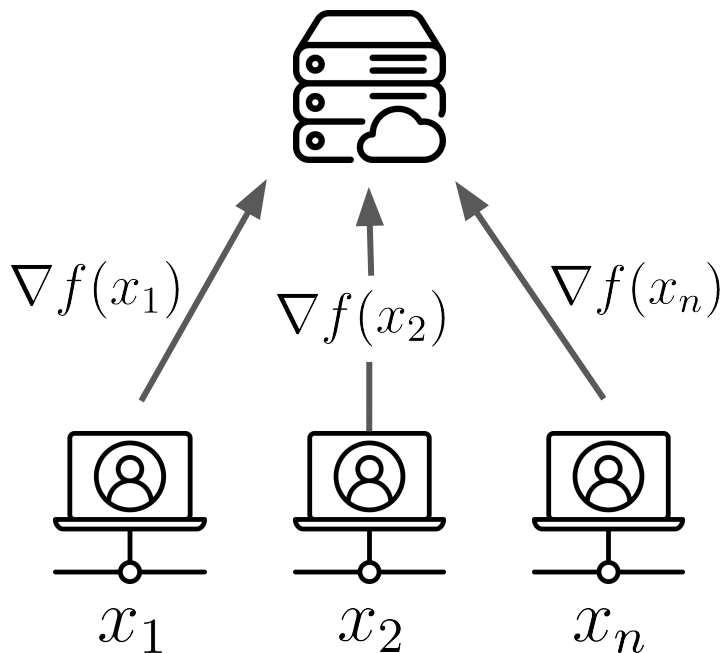
 **SRILAB**

ETH zürich

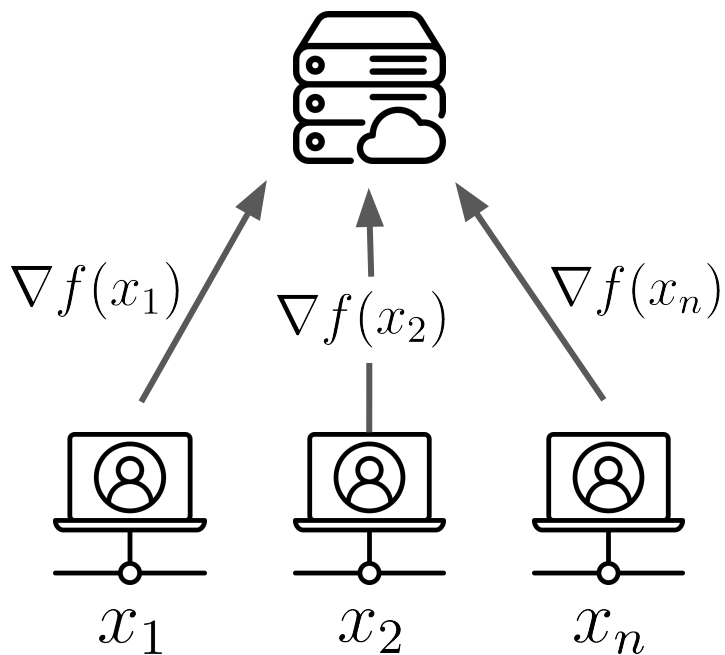


Federated Learning and Gradient Inversion

In **federated learning** a model is trained in a distributed fashion where clients only send **gradient updates** to the server in order to **preserve their data privacy**.



Federated Learning and Gradient Inversion

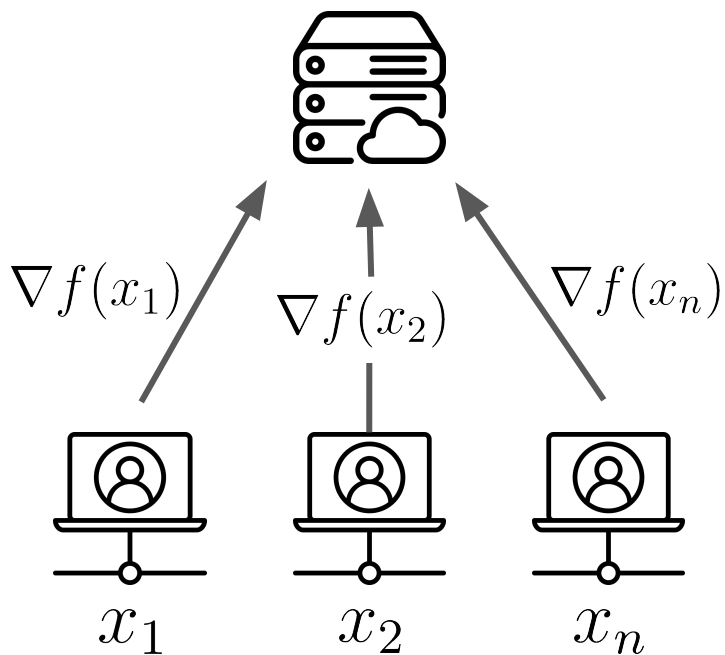


In **federated learning** a model is trained in a distributed fashion where clients only send **gradient updates** to the server in order to **preserve their data privacy**.

Unfortunately, an **adversarial server** can **reconstruct client data x** from received gradient updates:

$$\operatorname{argmin}_z \mathcal{E}(\nabla f(x), \nabla f(z))$$

Federated Learning and Gradient Inversion



In **federated learning** a model is trained in a distributed fashion where clients only send **gradient updates** to the server in order to **preserve their data privacy**.

Unfortunately, an **adversarial server** can **reconstruct client data** x from received gradient updates:

$$\operatorname{argmin}_z \mathcal{E}(\nabla f(x), \nabla f(z))$$

Key Question: When is it possible to exactly recover the client data?

Low Rankness of Gradient Updates

Assume linear layer with ReLU activation:

$$\begin{aligned} Z &= W \cdot X + b & X &\in \mathbb{R}^{N \times B} \\ Y &= \text{ReLU}(Z) & W &\in \mathbb{R}^{M \times N} \end{aligned}$$

Low Rankness of Gradient Updates

Assume linear layer with ReLU activation:

$$\begin{aligned} Z &= W \cdot X + b & X &\in \mathbb{R}^{N \times B} \\ Y &= \text{ReLU}(Z) & W &\in \mathbb{R}^{M \times N} \end{aligned}$$

The gradient of the linear layer can be written in the form:

$$\frac{\partial \mathcal{L}(f(x), y)}{\partial W} = \frac{\partial \mathcal{L}(f(x), y)}{\partial Z} \cdot X^T$$

Low Rankness of Gradient Updates

Assume linear layer with ReLU activation:

$$\begin{aligned} Z &= W \cdot X + b & X &\in \mathbb{R}^{N \times B} \\ Y &= \text{ReLU}(Z) & W &\in \mathbb{R}^{M \times N} \end{aligned}$$

The gradient of the linear layer can be written in the form:

$$\underbrace{\frac{\partial \mathcal{L}(f(x), y)}{\partial W}}_{M \times N} = \underbrace{\frac{\partial \mathcal{L}(f(x), y)}{\partial Z}}_{M \times B} \cdot \underbrace{X^T}_{B \times N}$$

Key Observation: The gradient of W is not full rank for $B < \min(N, M)$.

Sparse Gradients of ReLU Activations

Assume linear layer with ReLU activation:

$$\begin{aligned} Z &= W \cdot X + b & X &\in \mathbb{R}^{N \times B} \\ Y &= \text{ReLU}(Z) & W &\in \mathbb{R}^{M \times N} \end{aligned}$$

Gradients of ReLU activations:

$$\frac{\partial \mathcal{L}(f(x), y)}{\partial Z} = \frac{\partial \mathcal{L}(f(x), y)}{\partial Y} \underbrace{\odot \mathbb{I}_{[Z > 0]}}_{\text{0/1 mask based on the sign of Z}}$$

0/1 mask based on the sign of Z

Sparse Gradients of ReLU Activations

Assume linear layer with ReLU activation:

$$\begin{aligned} Z &= W \cdot X + b & X &\in \mathbb{R}^{N \times B} \\ Y &= \text{ReLU}(Z) & W &\in \mathbb{R}^{M \times N} \end{aligned}$$

Gradients of ReLU activations:

$$\frac{\partial \mathcal{L}(f(x), y)}{\partial Z} = \frac{\partial \mathcal{L}(f(x), y)}{\partial Y} \underbrace{\odot \mathbb{I}_{[Z > 0]}}_{\text{0/1 mask based on the sign of Z}}$$

0/1 mask based on the sign of Z

Key Observation: The ReLU activation makes gradient of Z sparse

Gradient Decomposition

Use **SVD** to create **low-rank decomposition** of the **gradient of W**:

$$\frac{\partial \mathcal{L}(f(x), y)}{\partial W} = U \cdot \Sigma \cdot V^T = \overbrace{U \cdot \sqrt{\Sigma}}^{L} \cdot \overbrace{\sqrt{\Sigma} \cdot V^T}^R$$

$M \times B$ $B \times N$

Gradient Decomposition

Use **SVD** to create **low-rank decomposition** of the **gradient of W**:

$$\frac{\partial \mathcal{L}(f(x), y)}{\partial W} = U \cdot \Sigma \cdot V^T = \underbrace{U \cdot \sqrt{\Sigma}}_{M \times B} \cdot \underbrace{\sqrt{\Sigma} \cdot V^T}_{B \times N}$$

We show that under mild assumptions: there exists a **unique** $Q \in \mathbb{R}^{B \times B}$ s.t.:

$$\frac{\partial \mathcal{L}(f(x), y)}{\partial Z} = L \cdot Q \quad X^T = Q^{-1} \cdot R$$

Key Observation: The Low-Rankness simplifies the problem from $N \times B$ to $B \times B$

Exploiting Gradient Sparsity (Example)

Assume that **first 3 neurons are not activated** for some input x in a batch with $B=3$:

$$\frac{\partial \mathcal{L}(f(x), y)}{\partial z} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0.1 \\ \vdots \\ 0.5 \end{bmatrix} = L \cdot q$$

Exploiting Gradient Sparsity (Example)

Assume that **first 3 neurons are not activated** for some input x in a batch with $B=3$:

$$\frac{\partial \mathcal{L}(f(x), y)}{\partial z} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0.1 \\ \vdots \\ 0.5 \end{bmatrix} = L \cdot q \quad \Rightarrow \quad \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = L_{[1,2,3]} \cdot q$$

Exploiting Gradient Sparsity (Example)

Assume that **first 3 neurons are not activated** for some input x in a batch with $B=3$:

$$\frac{\partial \mathcal{L}(f(x), y)}{\partial z} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0.1 \\ \vdots \\ 0.5 \end{bmatrix} = L \cdot q \quad \Rightarrow \quad \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = L_{[1,2,3]} \cdot q$$

If we also assume that $\text{rank}(L_{[1,2,3]}) = 2$, we have:

$$q = c \cdot \text{basis}(L_{[1,2,3]})$$

Exploiting Gradient Sparsity

Theorem: For **randomly initialized networks** with probability approximately $\frac{1}{2^b}$ a vector \mathbf{q}' obtained that way is a column of \mathbf{Q} (**up to scale**).

Key Idea: Sample random submatrices of \mathbf{L} and obtain the respective \mathbf{q}

Exploiting Gradient Sparsity

Theorem: For **randomly initialized networks** with probability approximately $\frac{1}{2^b}$ a vector \mathbf{q}' obtained that way is a column of **Q (up to scale)**.

To discard bad vectors \mathbf{q}' we compute **the sparsity** of their associated gradient:

$$\left\| \frac{\partial \mathcal{L}(f(x), y)}{\partial z'} \right\|_0 = \left\| \mathbf{L} \cdot \mathbf{q}' \right\|_0 < \tau \cdot M$$

Key Idea: Sample random submatrices of L and obtain the respective \mathbf{q}'

Recovering Scale

Let $q_i = c_i q'_i$, where c_i is **unknown**. Then one **can recover** c_i by:

$$\begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_B \end{bmatrix} = L^{-L} \cdot \overline{Q}^{-1} \cdot \frac{\partial \mathcal{L}(f(x), y)}{\partial b} \quad \text{with} \quad \overline{Q} = \begin{bmatrix} | & | & \dots & | \\ q'_1 & q'_2 & \dots & q'_B \\ | & | & \dots & | \end{bmatrix}$$

Key Observation: The gradient of the bias b can be used to recover scale.

Experimental Results

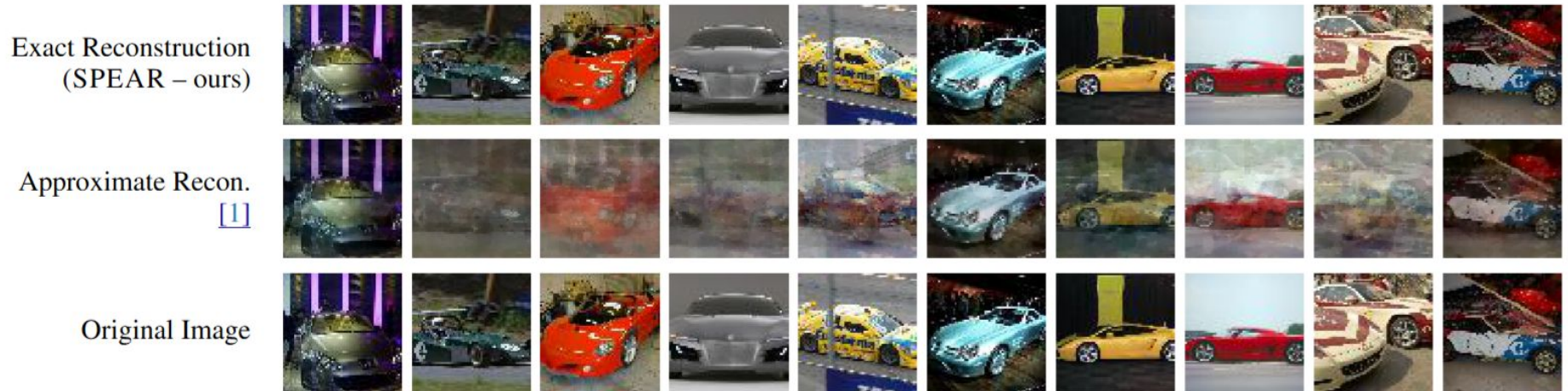
Main Results

We run our experiments on **6 layers FCNN** with **batch size B=20** on **image data**:

Method	PSNR \uparrow	Time/Batch
CI-Net [12] Sigmoid	38.0	1.6 hrs
CI-Net [12] ReLU	15.6	1.6 hrs
Geiping et al. [1]	19.6	18.0 min
SPEAR (Ours)	124.2	2.0 min

Dataset	PSNR \uparrow	LPIPS \downarrow	Acc (%) \uparrow	Time/Batch
MNIST	99.1	NaN	99	2.6 min
CIFAR-10	106.6	1.16×10^{-5}	99	1.7 min
TINYIMAGENET	110.7	1.62×10^{-4}	99	1.4 min
IMAGENET 224×224	125.4	1.05×10^{-5}	99	2.1 min
IMAGENET 720×720	125.6	8.08×10^{-11}	99	2.6 min

Main Results - Visualisation



Tabular Data Results

We also work with **tabular data** where FCNNs are common:

Method	Discr Acc (%) \uparrow	Cont. MAE \downarrow	Time/Batch
Tableak [8]	97	4922.7	2.6 min
SPEAR (Ours)	100	20.4	0.4 min

Convolutional Network Results

We first recover the **input features** to the **linear layers** of the **CNN** and then execute a combination of **feature inversion** and **gradient inversion attack**:

Method	LPIPS ↓	Feature Sim ↑
Geiping et al. [1]	0.562	-
CPA[9] + FI + Geiping et al. [1]	0.388	0.939
SPEAR + FI + Geiping et al. [1]	0.362	0.984

Further details can be found in the paper.

NeurIPS 2024 page:



OpenReview:

