# Deep Leak

MIT HAN LAB

Normal Client
Differentiable Model $f(x, W)$ → $Pred$ → $Loss$ ← [1.0, 0, 0]
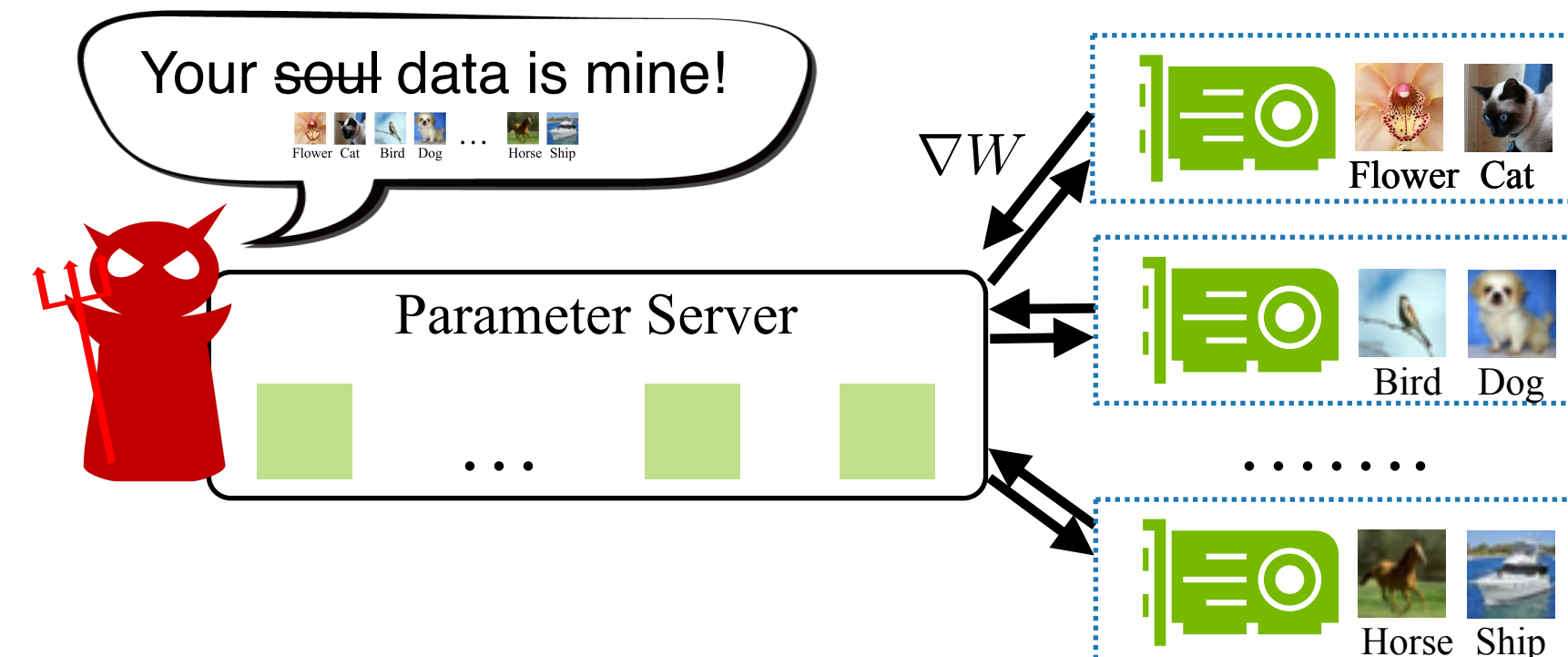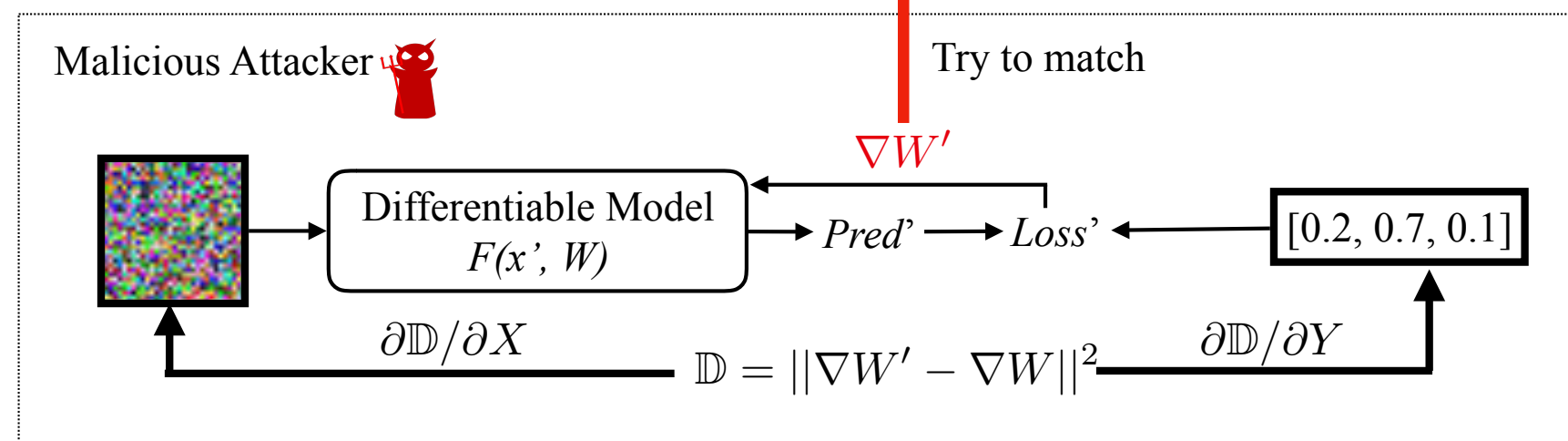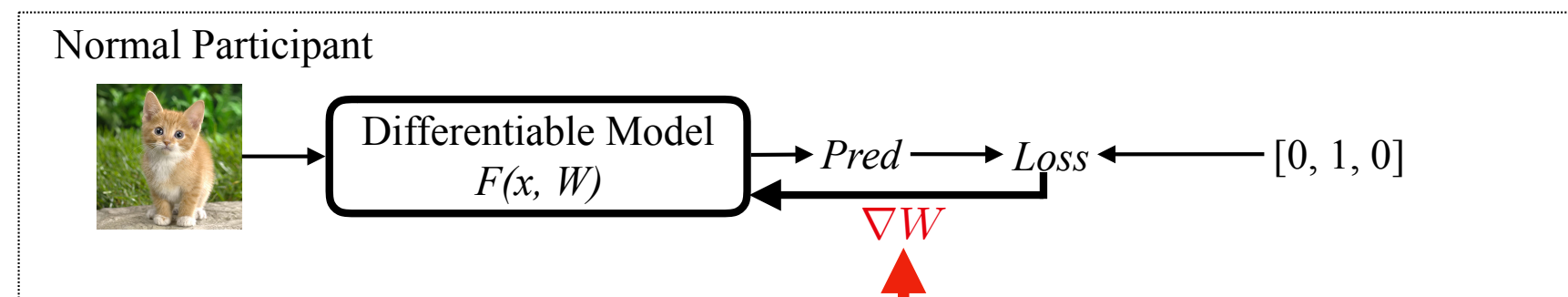$\nabla_W Loss$

Ligeng
Massach

## Rethink the Safety of Gradients Exchange

People believe it is safe to share gradients since numerical gradients seem have nothing related with semantic training set.
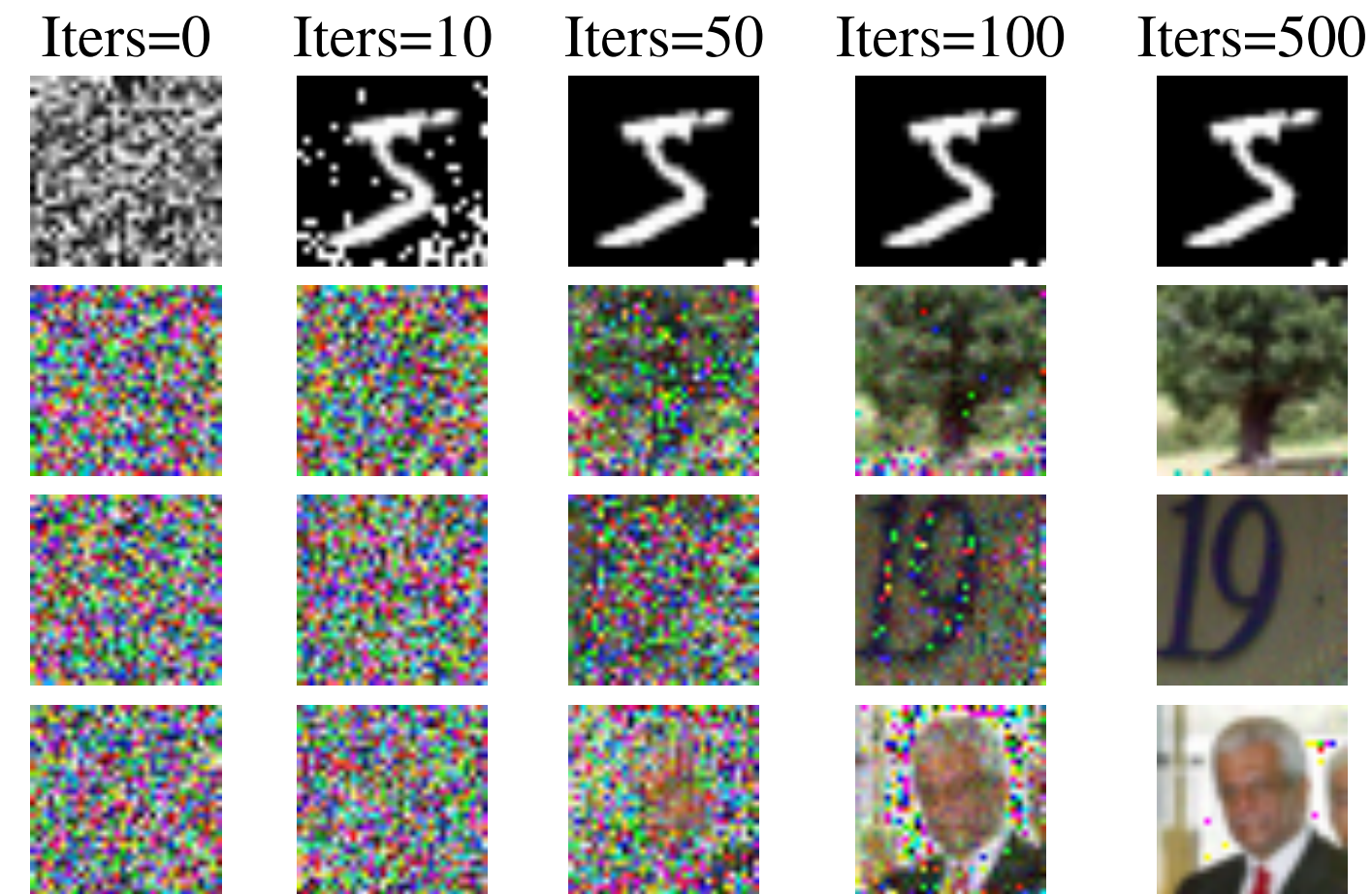
Based on the believe, collaborative & federated learning are developed, where *data never leaves local device* and only *public gradients are communicating between*.

But does the scheme really protect training data? Conventionally wisdom suggests yes. But we show that actually **it is NOT safe**.
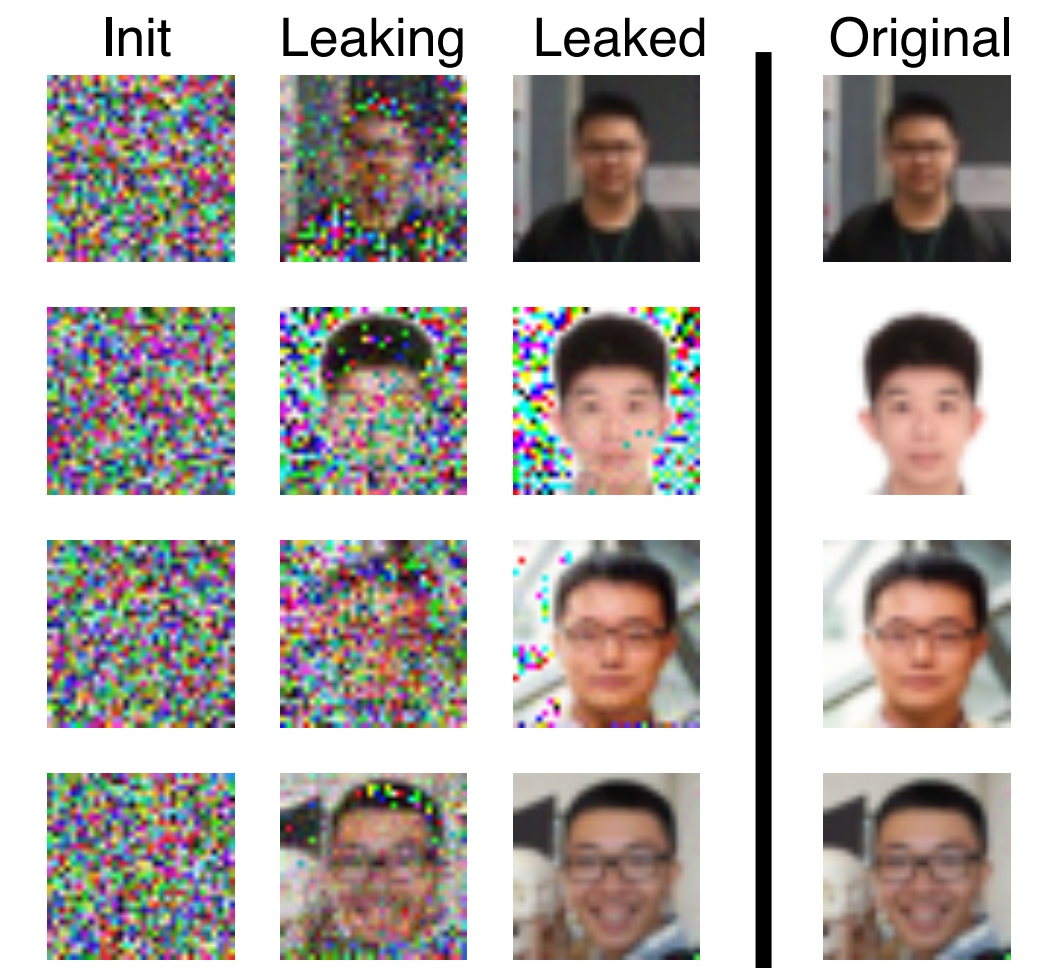
Malicious Attacker
$\Delta_W \mathbb{D}$   $\nabla_Y \mathbb{D}$
$\nabla_W Loss'$
Differentiable Model $f(x, W)$ → $Pred$ → $Loss$ ← [0.17, 0.2, 0.73]

Iters=0

## *D*eep *L*eakage from *G*radients (DLG)

Normal Participant
Differentiable Model $F(x, W)$ → $Pred$ → $Loss$ ← [0, 1, 0]
$\nabla W$

Malicious Attacker
Differentiable Model $F(x', W)$ → $Pred'$ → $Loss'$ ← [0.2, 0.7, 0.1]
Try to match
$\nabla W'$
$\partial \mathbb{D}/\partial X$   $\mathbb{D} = ||\nabla W' - \nabla W||^2$   $\partial \mathbb{D}/\partial Y$

Your ~~soul~~ data is mine!

Flower Cat Bird Dog ... Horse Ship

$\nabla W$

Parameter Server

Flower  Cat
Bird  Dog
Horse  Ship

Table 1

## DLG for batched data



Init   L2 Leaking   Leaked   Original

Note: the order may not be same as original's.

loss at log scale axis: $10^0$, $10^{-2}$, $10^{-4}$, $10^{-6}$, $10^{-8}$, $10^{-10}$, $10^{-12}$
iterations: 0, 50, 100, 150, 200

## DLG for NLP task

**[iter=0]**
tilting fill given **less word **itude fine **nton over- heard living vegas **vac-/vation *f forte **dis ce- rambycidae ellison **don yards marne **kali

**[iter=10]**
tilting fill given **less full solicitor other ligue shrill living vegas rider treatment carry played sculptures life- long ellison net yards marne **kali

**[Iter=20]**
registration , volunteer applications , at student travel application open the ; week of played ; child care will be glare

**[Iter=30]**
registration , volunteer applications , and student travel application open the first week of september . child care will be available .

**[Ground Truth]**
Registration, volunteer applications, and student travel application open the first week of September. Child care will be available.

Table 1

| | MNIST | CIFAR | SVHN | LFW |
|---|---|---|---|---|
| Ours | 0.0038 | 0.0069 | 0.0051 | 0.0055 |
| Melis | 0.2275 | 0.2578 | 0.2771 | 0.2951 |

■ Melis  ■ Ours

Mean Square Errors: 0.3, 0.2, 0.1
MNIST  CIFAR  SVHN  LFW

## defend

No leak
Leak with artifacts
Deep Leakage

original
laplacian-$10^{-4}$
laplacian-$10^{-3}$
laplacian-$10^{-2}$
laplacian-$10^{-1}$
Iterations: 0, 200, 400, 600, 800, 1000, 1200

No leak
Leak with artifacts
Deep Leakage

original
prune-ratio-1%
prune-ratio-10%
prune-ratio-20%
prune-ratio-30%
prune-ratio-50%
prune-ratio-70%
Iterations: 0, 200, 400, 600, 800, 1000, 1200

Table 1

| | MNIST | CIFAR | SVHN | LFW |
|---|---|---|---|---|
| Ours | 0.0038 | 0.0069 | 0.0051 | 0.0055 |
| Melis | 0.2275 | 0.2578 | 0.2771 | 0.2951 |