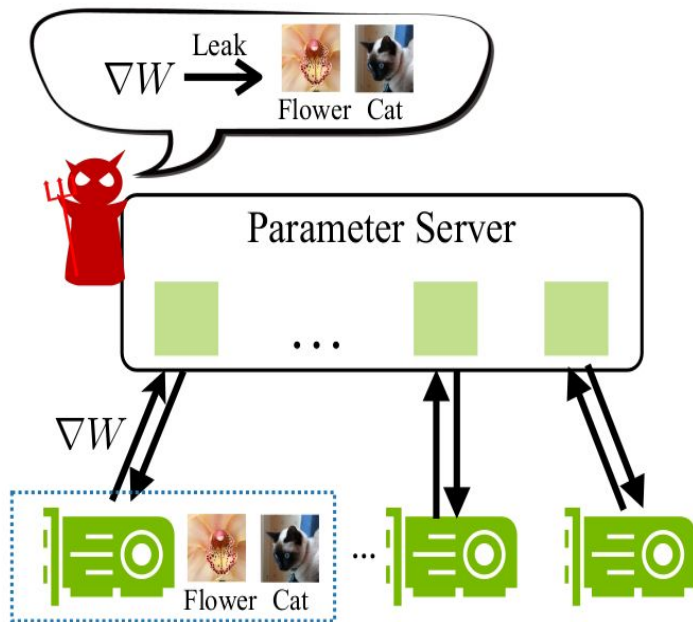

Deep Leakage from Gradients

— By: Wafaa Mohammed —

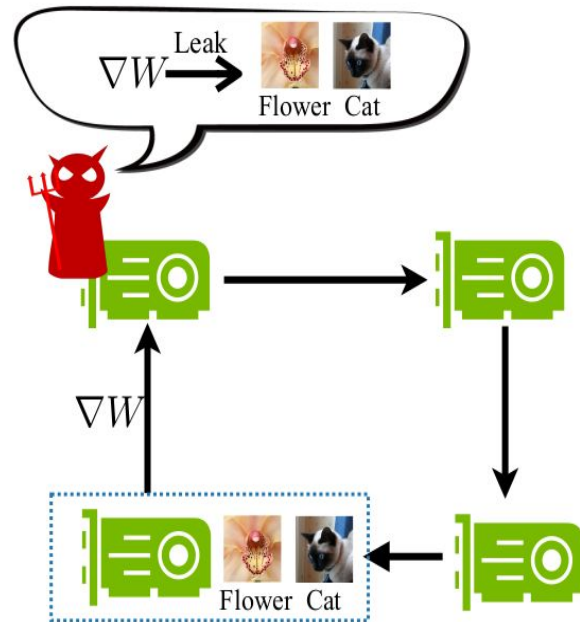
The Idea

Main Ideas

- For a long time, people believed that gradients are safe to share (the training data will not be leaked).
- The paper shows that it is possible to obtain the private training data from the publicly shared gradients.
- Without changes on training setting, the most effective defense method is gradient pruning.



(a) Distributed training with a centralized server



(b) Distributed training without a centralized server

The DLG Algorithm

Normal Participant



Differentiable Model
 $F(x, W)$

$Pred$

$Loss$

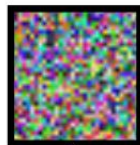
$[0, 1, 0]$

∇W

Malicious Attacker



Try to match



Differentiable Model
 $F(x', W)$

$Pred'$

$Loss'$

$[0.2, 0.7, 0.1]$

$\nabla W'$

$\partial \mathbb{D} / \partial X$

$\mathbb{D} = \|\nabla W' - \nabla W\|^2$

$\partial \mathbb{D} / \partial Y$

Algorithm 1 Deep Leakage from Gradients.

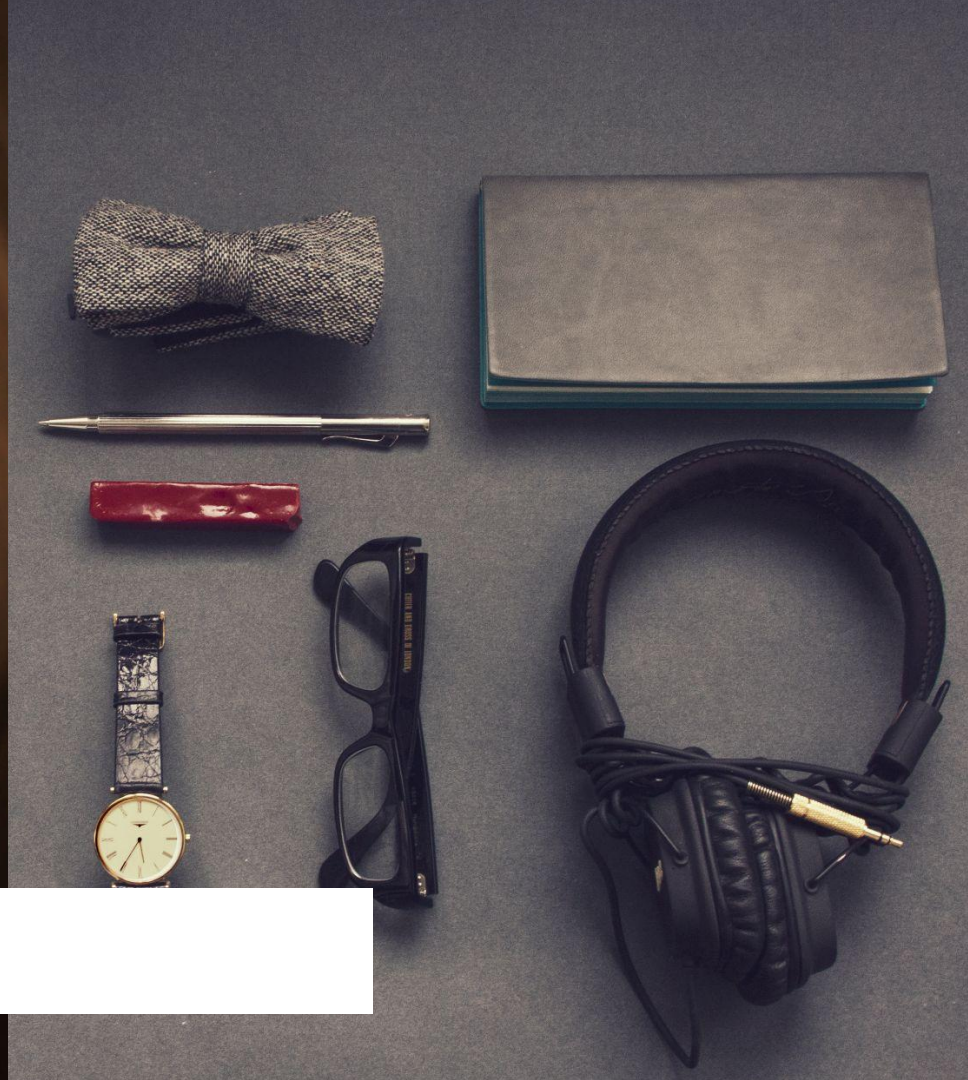
Input: $F(\mathbf{x}; W)$: Differentiable machine learning model; W : parameter weights; ∇W : gradients calculated by training data

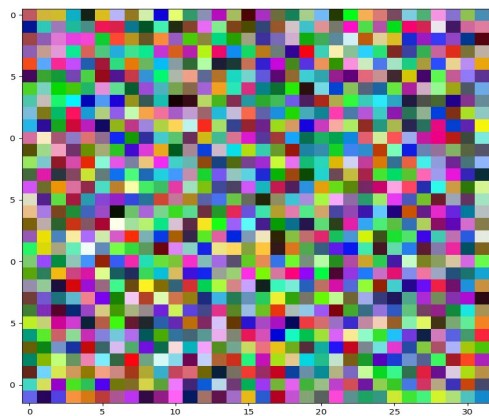
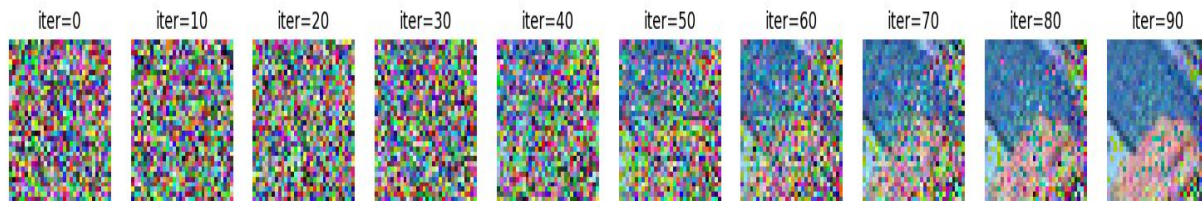
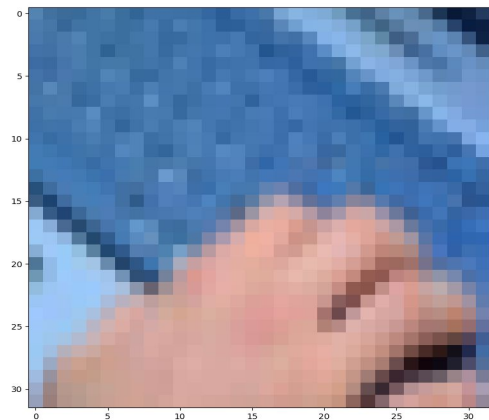
Output: private training data \mathbf{x}, \mathbf{y}

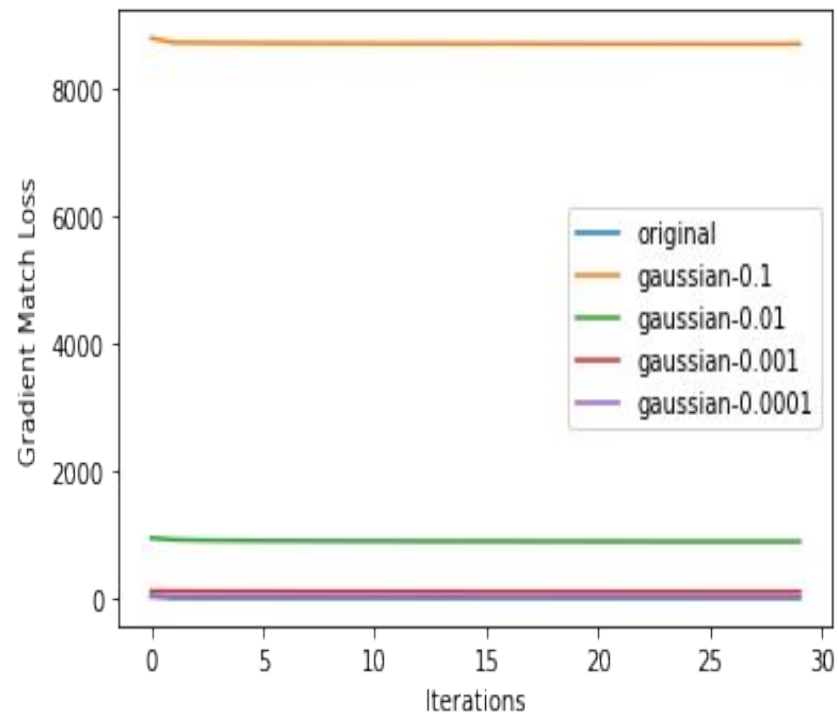
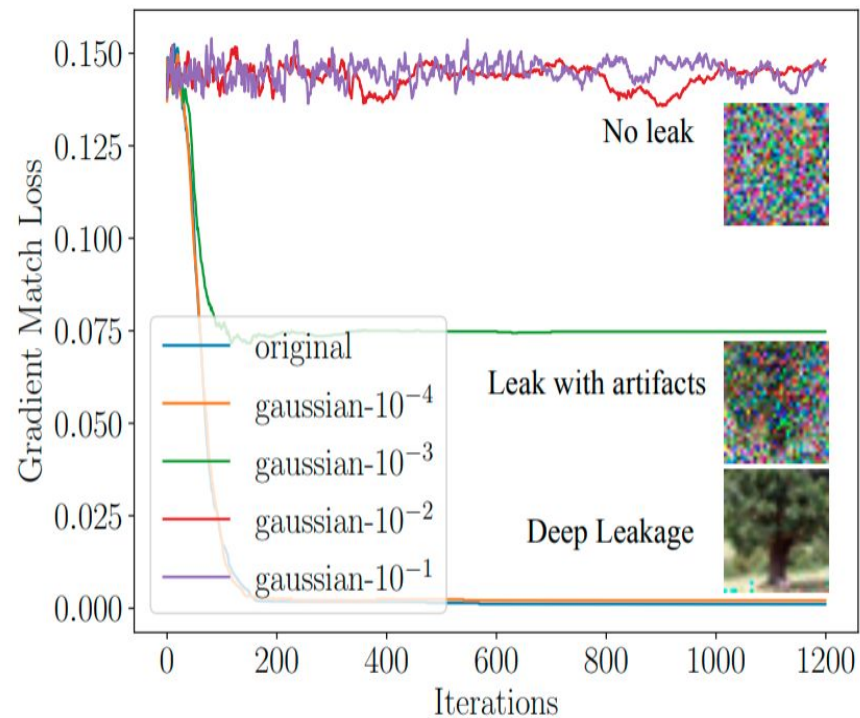
```
1: procedure DLG( $F, W, \nabla W$ )  
2:    $\mathbf{x}'_1 \leftarrow \mathcal{N}(0, 1), \mathbf{y}'_1 \leftarrow \mathcal{N}(0, 1)$  ▷ Initialize dummy inputs and labels.  
3:   for  $i \leftarrow 1$  to  $n$  do  
4:      $\nabla W'_i \leftarrow \partial \ell(F(\mathbf{x}'_i, W_t), \mathbf{y}'_i) / \partial W_t$  ▷ Compute dummy gradients.  
5:      $\mathbb{D}_i \leftarrow \|\nabla W'_i - \nabla W\|^2$   
6:      $\mathbf{x}'_{i+1} \leftarrow \mathbf{x}'_i - \eta \nabla_{\mathbf{x}'_i} \mathbb{D}_i, \mathbf{y}'_{i+1} \leftarrow \mathbf{y}'_i - \eta \nabla_{\mathbf{y}'_i} \mathbb{D}_i$  ▷ Update data to match gradients.  
7:   end for  
8:   return  $\mathbf{x}'_{n+1}, \mathbf{y}'_{n+1}$   
9: end procedure
```



Experiments and results







Conclusions

- Deep Leakage from Gradients (DLG): an algorithm that can obtain the local training data from public shared gradients.
- DLG does not rely on any generative model or extra prior about the data.
- Defense strategies include noisy gradients, gradient perturbation (half precision), and gradient compression (pruning).