
Privacy-Preserving Split Learning via Patch Shuffling over Transformers

Dixi Yao, Liyao Xiang, Hengyuan Xu, Hangyu Ye, Yingqi Chen

John Hopcroft Center, Shanghai Jiao Tong University, China

Sept., 2022

饮水思源 · 爱国荣校

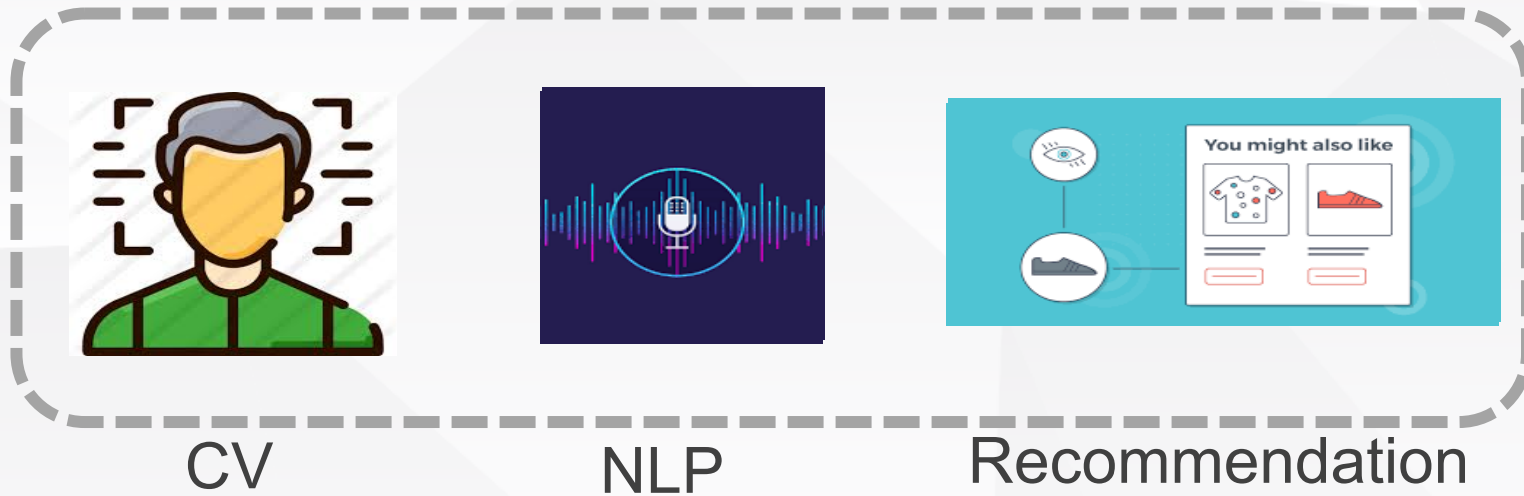


01

Background



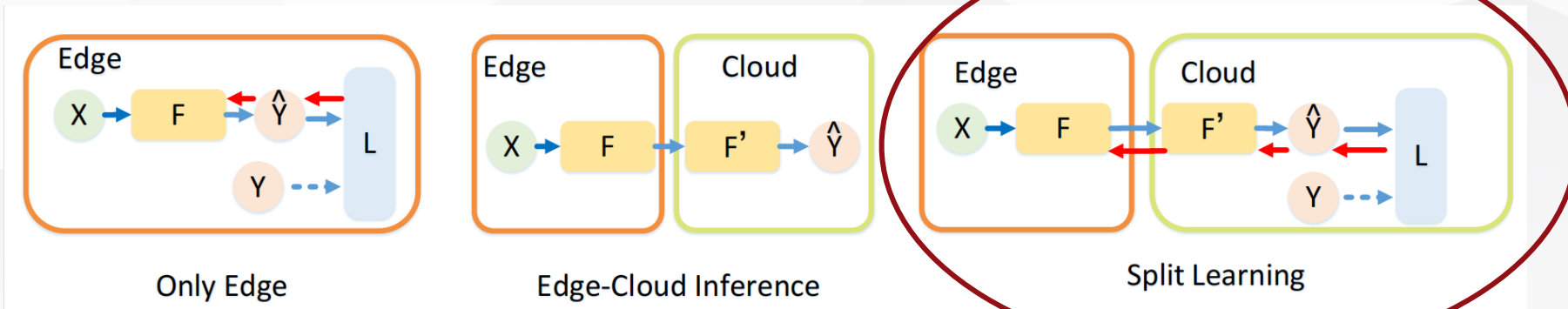
New Computational Paradigm



Compute on edge: resource constrained



Upload to cloud: privacy leakage





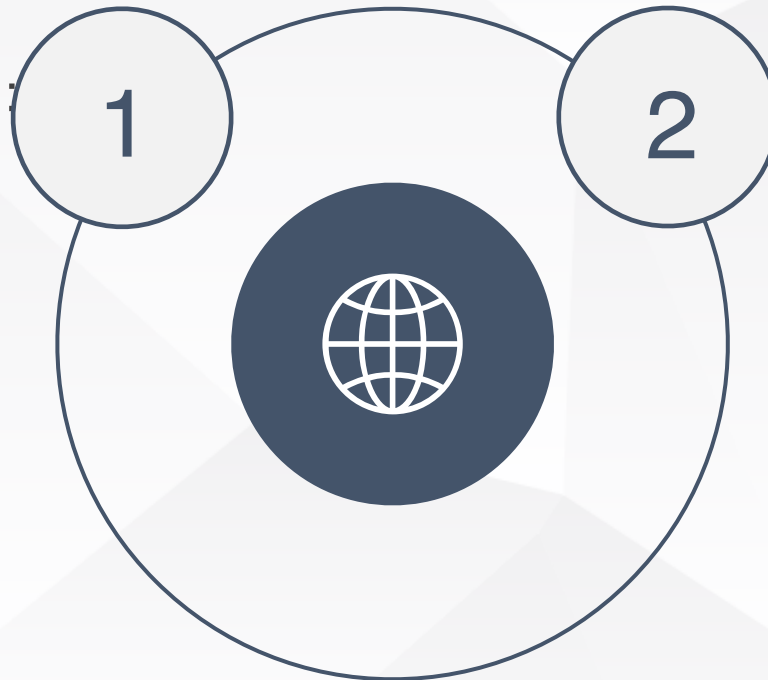
Is Split learning perfect?



Challenge 1

Unprotected intermediate results :

leak privacy of input !



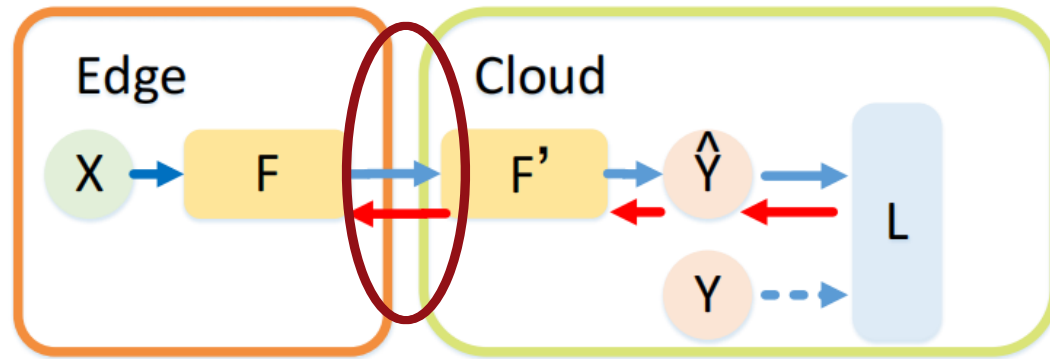
Challenge 2

Protect label privacy :

Labels should not leave cloud
if labels are proprietary



An Example



Split Learning



Facial images:
private on edges

Forward loop:
intermediate features

Backward loop:
error gradients

Identity: **Bob**
belongs to a proprietary
enterprise database



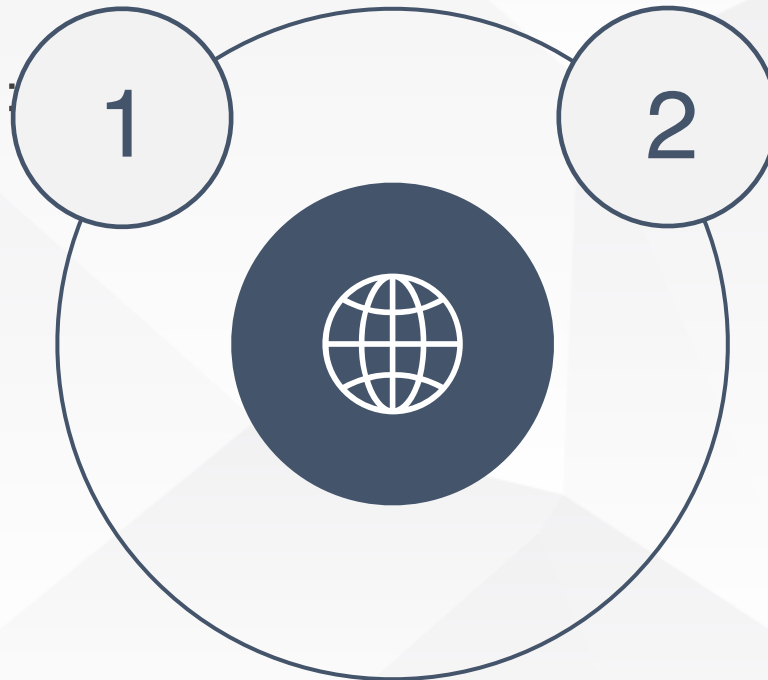
Is Split learning perfect?



Challenge 1

Unprotected intermediate results :

Leak privacy of input !



Challenge 2

Protect label privacy :

Labels should not leave cloud if labels are proprietary

Challenge 3

Privacy in training

Leakage would occur in each iteration





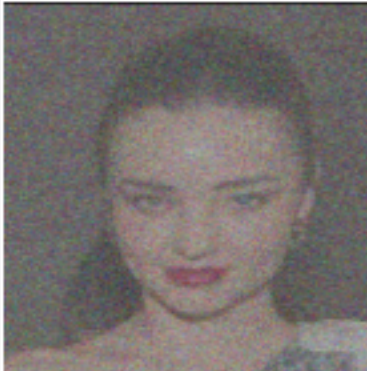
Protecting training data privacy is hard



Inference: one-time transmission

Training: multiple forward & backward rounds

Privacy should be guaranteed throughout training!



Add Noise

Adding Gaussian noise barely works

Adversarial learning based methods:

Protection is effective only at convergence 😄





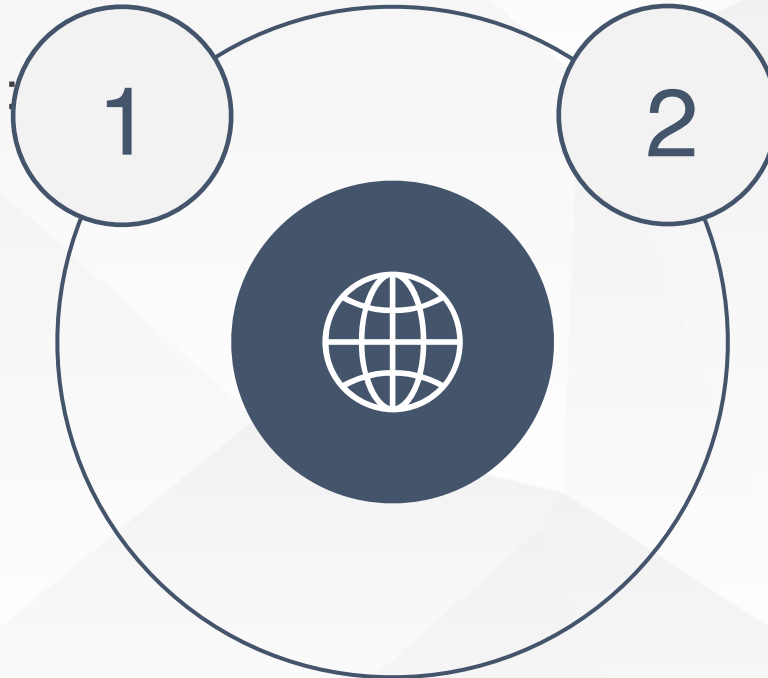
Is Split learning perfect?



Challenge 1

Unprotected intermediate results :

Leak privacy of input !



Challenge 2

Protect label privacy :

Labels should not leave cloud if labels are proprietary

Challenge 3

Privacy in training

Leakage would occur in each iteration

Challenge 4

Practicality in deployment





DNN on thin edge devices:

Low in efficiency --- cryptographic tools including homomorphic encryption, multi-party computation

High training performance:

Sacrifice of accuracy --- differential privacy

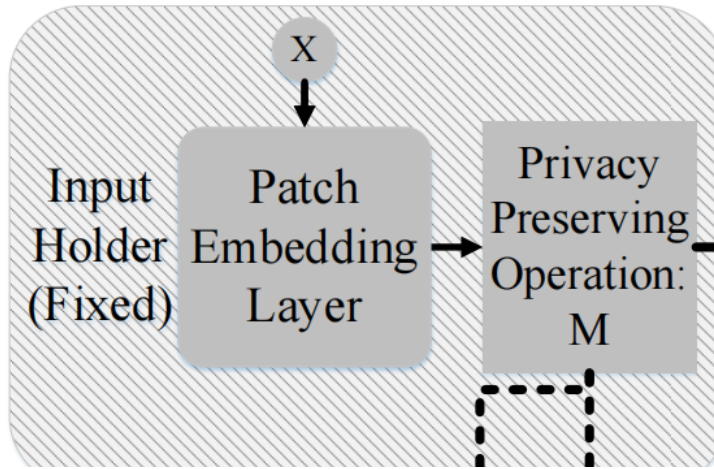
A photograph of a modern building with a white, faceted facade and large glass windows, set against a blue sky with light clouds. The building is the central focus of the upper half of the slide.

02

Threat Model & Methodology



Objective: minimize task loss and maximize attacker reconstruction loss





White-box attack

Attacker's prior:

- ✓ Intermediate features
- ✓ Model weights

Black-box attack

Attacker's prior:

- ✓ Intermediate features
- ✓ Auxiliary datasets
- × Model weights

Adaptive attack

Similar to Black-box

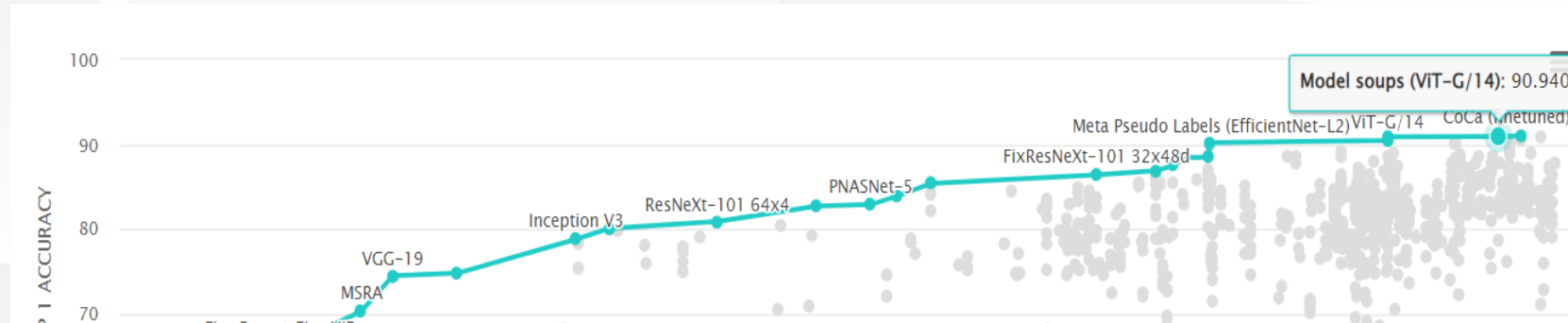
Use features from multiple rounds

Attacker's prior:

- ✓ multiple features
- ✓ Auxiliary datasets
- × Model weights



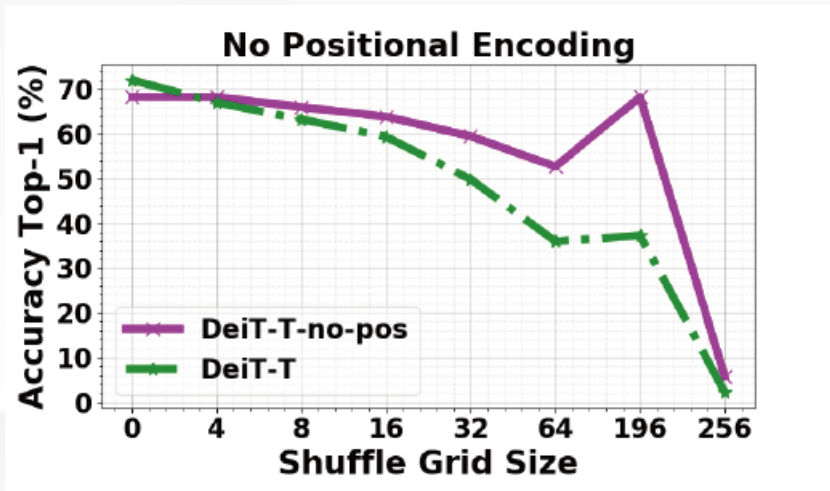
Property of Transformer



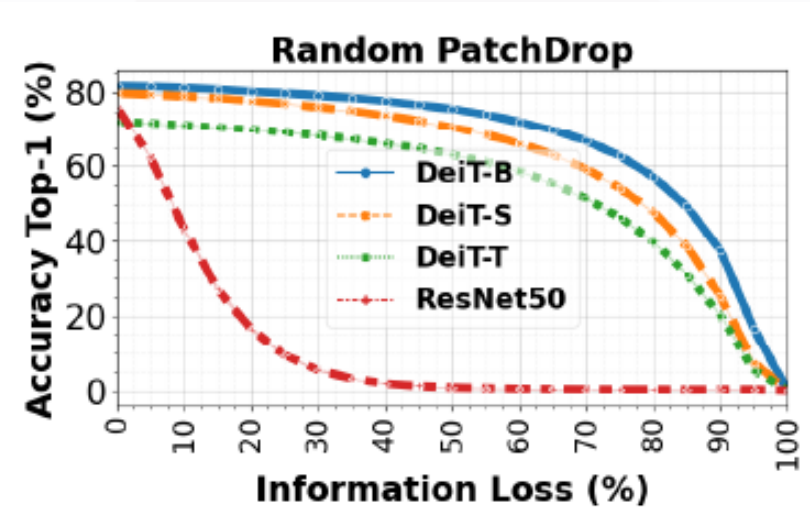
Transformer has shown a superior accuracy

ImageNet-1k (from paperswithcode.com)

Shuffling Invariance



Robustness against Patch Dropping





Privacy Definition



a permutation (1,4,8,9,7,2,3,5,6)

1	2	3
4	5	6
7	8	9

1	2	3
4	5	6
7	8	9

1	2	3
4	5	6
7	8	9

1	6	7
2	8	9
5	3	4

4	9	8
7	3	2
6	1	5

8	4	9
3	7	2
6	5	1

1	7	7
5	6	2
3	5	2

5	9	8
7	4	9
1	3	6

9	4	1
4	8	3
6	2	8

(a) Original Input

(b) Output of Patch Shuffling

(c) Output of Batch Shuffling

Definition 1. (Neighbouring Permutations) We divide a single instance into N patches, and the permutations of these N patches constitute S . Any two permutation $\sigma, \sigma' \in S$ are defined to be neighboring.

Definition 2. (σ -privacy) Given private dataset X and a set of permutations S , a randomized mechanism $\mathcal{A} : f(X) \mapsto \mathcal{V}$ is σ -private if for all $x \in X$, neighbouring permutations σ and σ' and any $z \in \mathcal{V}$, we have

$$\Pr[\mathcal{A}(\sigma(f(x))) = z] = \Pr[\mathcal{A}(\sigma'(f(x))) = z]. \quad (6)$$

Each permutation has the same likelihood to generate z .





Patch Shuffling



1	2	3
4	5	6
7	8	9

1	2	3
4	5	6
7	8	9

1	2	3
4	5	6
7	8	9

1	6	7
2	8	9
5	3	4

4	9	8
7	3	2
6	1	5

8	4	9
3	7	2
6	5	1

(a) Original Input

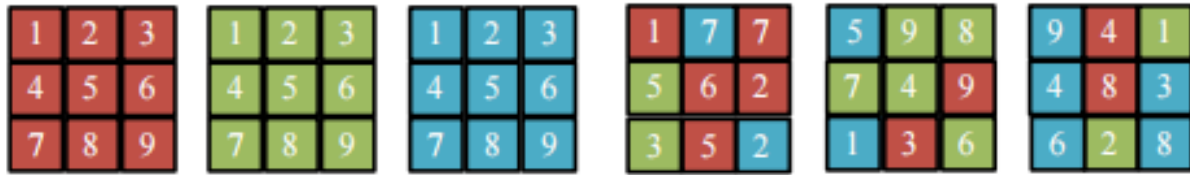
(b) Output of Patch Shuffling



Apply a permutation to shuffle patches within an image



Each permutation has $P_r = 1/N!$ (e.g., $N=196$) to produce z



(a) Original Input

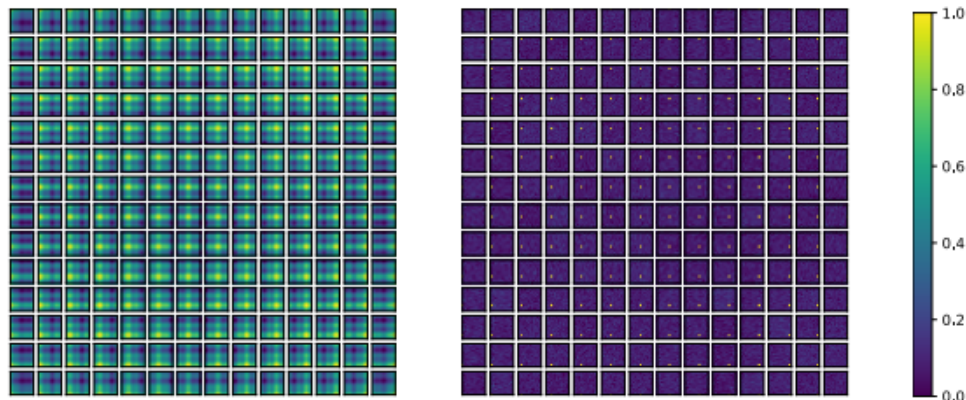
(c) Output of Batch Shuffling

Batch Shuffling:

Parameters:

- Proportion of patches shuffled across diff. images within a batch
- Proportion of patches shuffled across diff. batches

Spec



Time Domain

Frequency Domain

Position Embedding

before patch shuffling

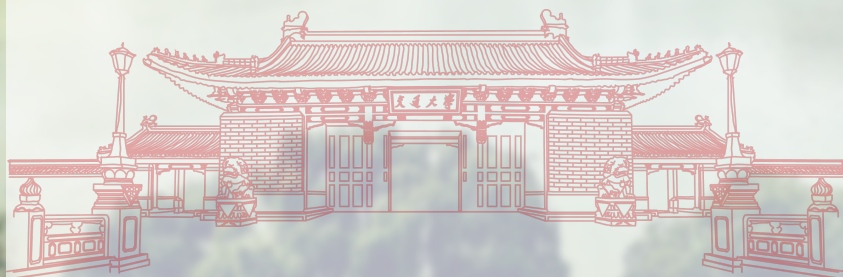
Further eliminate positional correlation between patches

So that each permutation has equal prob. to occur

(c) Output of Spectral Shuffling

03

Evaluation



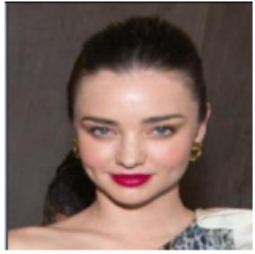


Black-Box Attack (MAE Decoder)

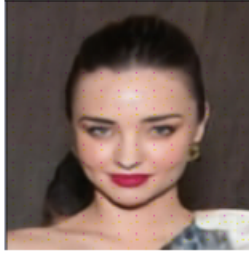


Accuracy VS Privacy: BS --- Batch Shuffling, PS --- Patch Shuffling, PS+ --- Spectral Shuffling

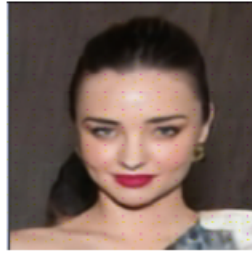
➤ Visualization effect of CelebA reconstruction



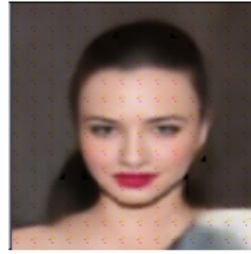
(a) Input



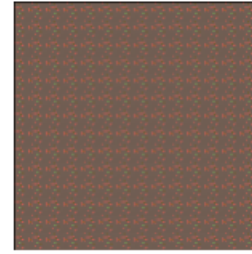
(b) SL



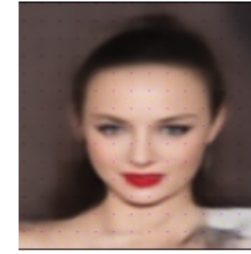
(c) Adv



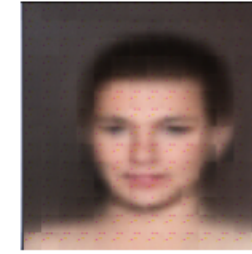
(d) Blur



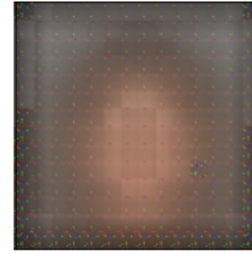
(e) DP



(f) GN



(g) Our BS



(h) Our PS+

Accuracy(%)

91.05

90.36

89.58

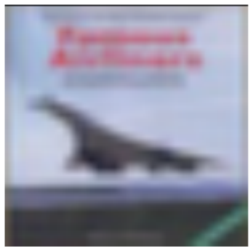
80.67

87.35

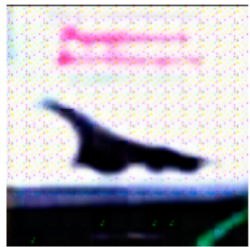
89.18

88.21

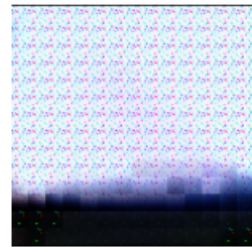
➤ Visualization effect of CIFAR10 reconstruction



(a) Input



(b) SL



(c) Our PS



(d) Our BS 75

Accuracy(%)

98.36

96.99

96.16

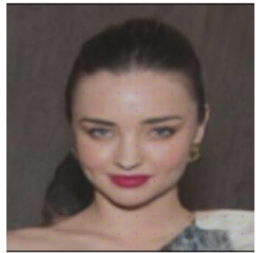
Criteo

Methods	Utility: Acc ↑	Privacy: MSE ↑
SL	77.81	0.0012
Our PS	77.78	0.0015
GN	77.28	0.0012

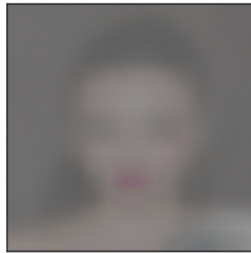




Attacker is aware of the model weights, but not the permutation order



(a) SL/Adv



(b) Blur



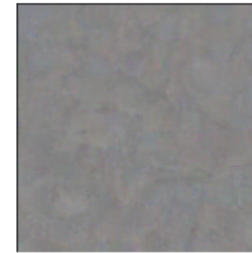
(c) DP



(d) GN



(e) Our PS



(f) Our BS



(g) Our PS+



(h) Jigsaw to
Our BS

A stronger threat: Jigsaw solving

Train a model to guess the permutation order

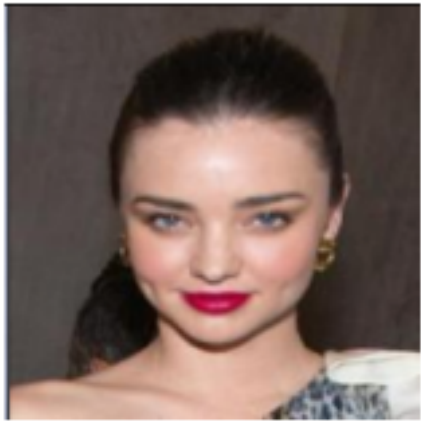


Failed due to random permutation

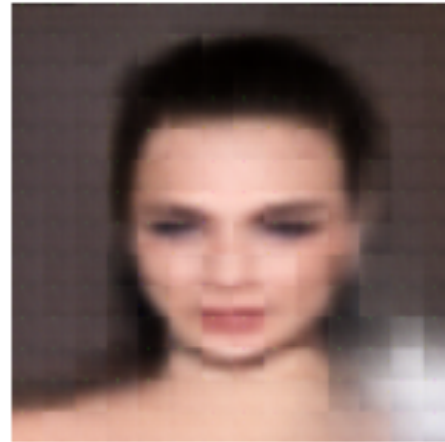


Attackers intercept the intermediate results throughout the whole training process

- We use 30 rounds of intermediate results to attack



(a) Input



(a) Our BS

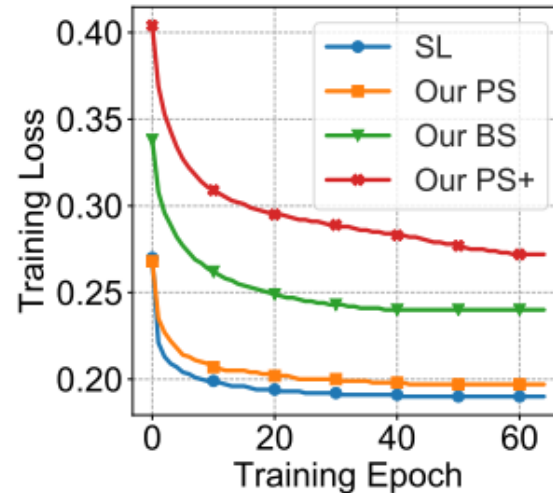
Failed to recover the original images



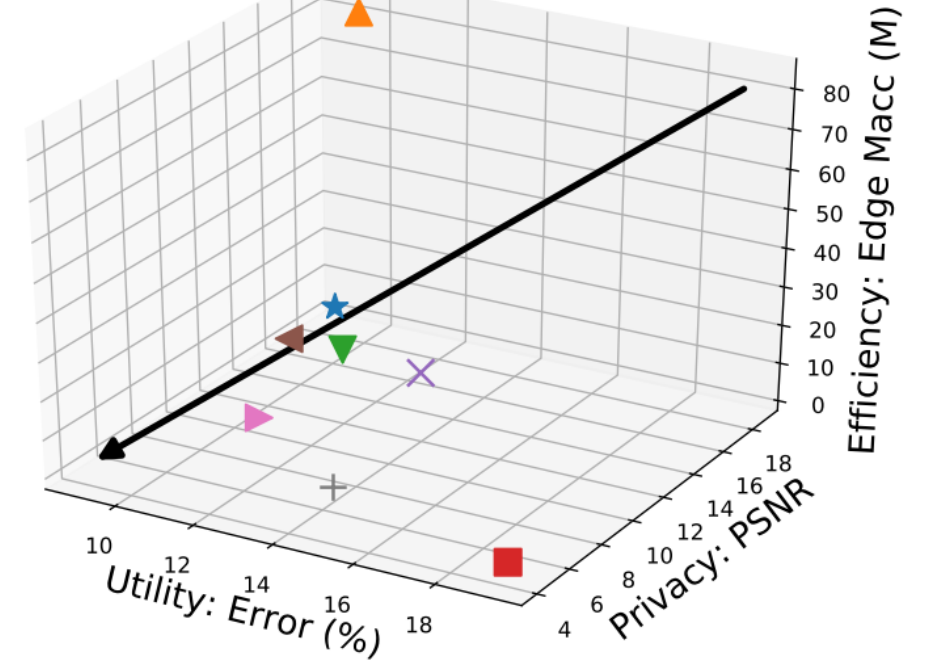
Efficiency, CelebA

Computational and memory costs at the edge, lower is better

Methods	Macc Edge (M)↓	Mem Edge (G)↓
SL / Transform	3.10	0.97
Adv	81.63	2.43
Our PS/BS	3.10	0.97
Our PS+	1.18	1.01



Convergence curves



Privacy, Utility & Efficiency, CelebA:

➤ Our methods achieve ideal tradeoffs

➤ Our methods have negligible impact to standard split learning

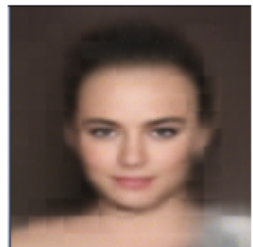




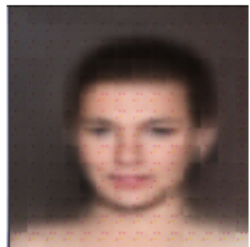
Ablation Studies



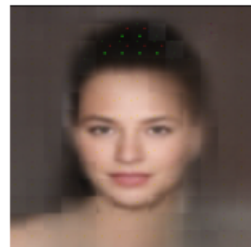
➤ k: Proportion of patches shuffled across diff. images within a batch



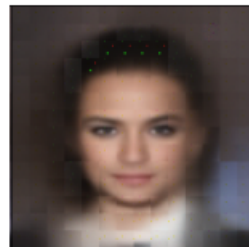
k: 0.5



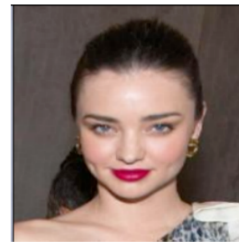
0.6



0.75



0.85



Original Input

Acc.(%): 90.29

89.18

88.54

88.76

- k = 0.6 exhibits the best tradeoff
- a smaller k leads to better reconstruction and higher accuracy

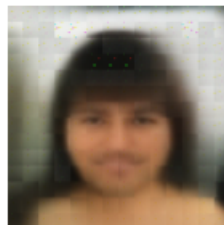
➤ Transferability: against black-box attacks with auxiliary datasets



(a) Input

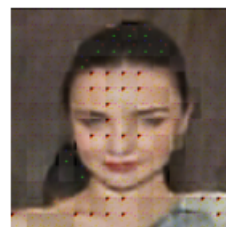


(b) SL

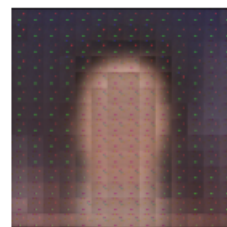


(c) Our BS

Auxiliary set: CelebA
Private set: LFW



(a) SL

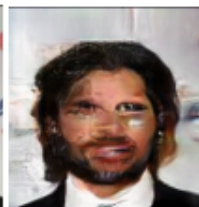


(b) Our BS

Auxiliary set: LFW
Private set: CelebA



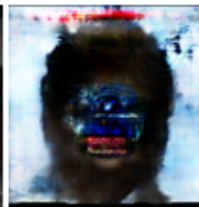
(a) Input



(b) SL



(c) BS



(d) PS+

➤ Adaptability: change attack model to CNN model --- Pix2Pix





An efficient privacy-preserving approach in split learning



A formal privacy guarantee based on patch shuffling



Eliminating positional correlation by spectral shuffling



上海交通大學

SHANGHAI JIAO TONG UNIVERSITY

Thanks!

飲水思源 愛國榮校