# All Rivers Run to the Sea: Private Learning with Asymmetric Flows

Yue Niu[1], **Ramy E. Ali**[2], Saurav Prakash[3], Salman Avestimehr[1]

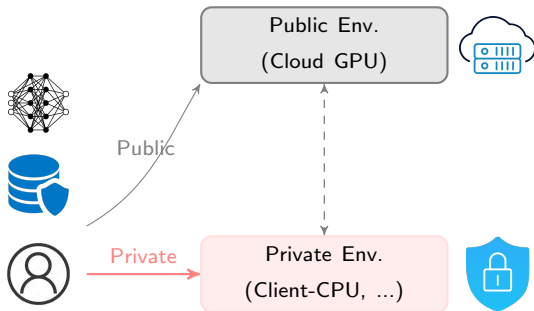[1]University of Southern California (USC)

[2]Samsung (was with USC)
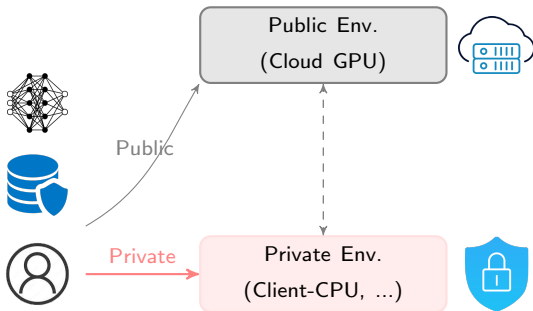
[3]University of Illinois Urbana-Champaign

Feb. 22, 2024
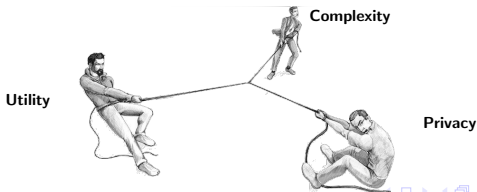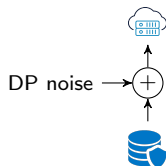
# Outline

# Outline

# How to leverage cloud ML while ensuring privacy?

# How to leverage cloud ML while ensuring privacy?



The Utility-Privacy-Complexity Trilemma

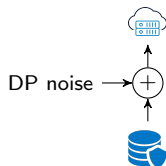(Naive) DP-based ML



DP noise $\longrightarrow$ $(+)$

- Provable guarantee
- Severe accuracy drop

(Naive) DP-based ML

Crypto-based ML

- Provable guarantee
- Severe accuracy drop

- Strong protection
- High complexity

# Privacy-Preserving ML Approaches



(Naive) DP-based ML

Crypto-based ML

Secure Enclaves

- Provable guarantee
- Severe accuracy drop

- Strong protection
- High complexity

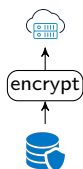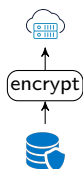- Hardware security
- Long running time

# Privacy-Preserving ML Approaches
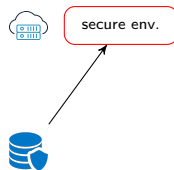


**(Naive) DP-based ML**

- Provable guarantee
- Severe accuracy drop

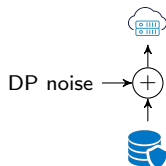**Crypto-based ML**

- Strong protection
- High complexity

**Secure Enclaves**

- Hardware security
- Long running time

Our work: Leverage both DP & Trusted hardware (local CPU, ...)
→ Overcome accuracy drop of naive-DP & long running time of TEEs

# Privacy-Preserving ML Approaches

<u>(Naive) DP-based ML</u>

<u>Crypto-based ML</u>

<u>Secure Enclaves</u>



- Provable guarantee
- Severe accuracy drop

- Strong protection
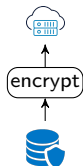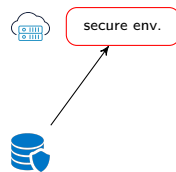- High complexity

- Hardware security
- Long running time

Related works leveraging TEEs

- Slalom'18: Inference only → This work: Inference and Training
- 3LegRace'21: Layerwise TEE-GPU communication → This work: No layer-wise communication

# Outline

## What does Delta do?

Decompose model & data into a low-dimensional part & a residual part

1. Lightweight model (client-side, TEEs, ...)

    - Fed with the low-dimensional information-sensitive part of the data
    - Confidential computing (no DP noise needed)

# Delta: Private Learning with Asymmetric Flows

## What does Delta do?

Decompose model & data into a low-dimensional part & a residual part

1. Lightweight model (client-side, TEEs, ...)

    - Fed with the low-dimensional information-sensitive part of the data
    - Confidential computing (no DP noise needed)

2. Large model (offloaded to cloud)

    - Fed with the quantized residual part of the data
    - The residual data is protected by a DP noise

# Delta: Private Learning with Asymmetric Flows

## What does Delta do?

Decompose model & data into a low-dimensional part & a residual part

1. Lightweight model (client-side, TEEs, ...)

   - Fed with the low-dimensional information-sensitive part of the data
   - Confidential computing (no DP noise needed)

2. Large model (offloaded to cloud)

   - Fed with the quantized residual part of the data
   - The residual data is protected by a DP noise

   $\Rightarrow$ Delta provides better utility-privacy trade-off than naive-DP methods

Forward Propagation: Asymmetric Data Decomposition



$\rightarrow$ To leverage the low-rank structure of the data

Forward Propagation: Perturbation & Binary Quantization



$\rightarrow$ To ensure privacy and reduce communication cost

Forward Propagation: Model Decomposition



$\rightarrow$ To ensure low complexity in the private environment

## Asymmetric Data Decomposition



- SVD $\rightarrow$ asymmetric decomposition along channel dimension
- DCT $\rightarrow$ asymmetric decomposition along spatial dimension

Why asymmetric decomposition?

SVD Approximation Error



Error

0.1

0

12.5%    25%    37.5%    50%

$$\frac{\|x - x_{lr}\|}{\|x\|}$$

Fraction of principal channels in $X_{lr}$

DCT Approximation Error



Error

0.1

0

8%    18%    32%    50%

$$\frac{\|x - x_{lf}\|}{\|x\|}$$

Fraction of low-freq components in $X_{lf}$

# Delta: Detailed Procedure

## Random Binary Quantization



$$\mathsf{IR_{quant}}(\cdot) = \mathrm{BinQuant}(\mathsf{IR_{noisy}}(\cdot)) = \begin{cases} 0 & \mathsf{IR_{noisy}}(\cdot) < 0 \\ 1 & \mathsf{IR_{noisy}}(\cdot) \geq 0 \end{cases}$$

# Delta: Detailed Procedure

## Random Binary Quantization



$$\mathsf{IR}_{\mathsf{quant}}(\cdot) = \mathrm{BinQuant}(\mathsf{IR}_{\mathsf{noisy}}(\cdot)) = \begin{cases} 0 & \mathsf{IR}_{\mathsf{noisy}}(\cdot) < 0 \\ 1 & \mathsf{IR}_{\mathsf{noisy}}(\cdot) \geq 0 \end{cases}$$

**Theorem:** Delta ensures that any operation in the public environment satisfy $(\epsilon, \delta)$-DP given noise $\mathcal{N}(0, p\Delta/\epsilon \cdot \sqrt{2\log(1.25/\delta)})$ and mini-batch size $b$, where $p = b/N$ is the sampling probability.

## Private Backpropagation



$$\mathcal{M}_{\text{main}} : \boldsymbol{o}_{\text{tot}}(i) = \frac{e^{\boldsymbol{z}_{\text{main}}(i) + \boldsymbol{z}_{\text{res}}(i)}}{\sum_{j=1} e^{\boldsymbol{z}_{\text{main}}(j) + \boldsymbol{z}_{\text{res}}(j)}} \quad \text{for} \quad i = 1, \cdots, L$$

$$\mathcal{M}_{\text{res}} : \boldsymbol{o}_{\text{res}}(i) = \frac{e^{\boldsymbol{z}_{\text{res}}(i)}}{\sum_{j=1} e^{\boldsymbol{z}_{\text{res}}(j)}} \quad \text{for} \quad i = 1, \cdots, L,$$

# Delta: Full Picture



- Asymmetric data decomposition
- Efficient model design
- Random binary quantization
- Private backpropagation

# Outline

# Experiment Highlights: Model Utility



$\rightarrow$ Lightweight model achieves good accuracy, but still residuals are useful
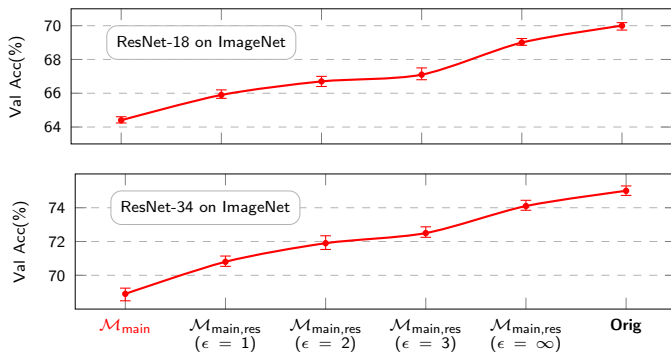
Setting: ResNet-18 with $\epsilon = 1$

|  | Delta: perturb $\text{IR}_{\text{res}}$ | naive-DP: perturb IR |
|---|---|---|
| CIFAR-10 | 92.4% | 69.6% ($\downarrow -22.8$) |
| CIFAR-100 | 71.4% | 48.3% ($\downarrow -23.1$) |
| ImageNet | 65.9% | 34.4% ($\downarrow -31.5$) |

$\rightarrow$ Delta improves accuracy by up to 31.5%

# Experiment Highlights: Model Complexity

MACs of the modules in `Delta`

|  | $\mathcal{M}_{bb}+\mathcal{M}_{main}$ | SVD | DCT | $\mathcal{M}_{res}$ |
|---|---|---|---|---|
| ResNet-18 | 48.3 M | 0.52 M | 0.26 M | 547M |
| ResNet-34 | 437 M | 1.6 M | 0.7 M | 3.5G |

- Small model $\mathcal{M}_{main}$ only costs 10% complexity of $\mathcal{M}_{res}$
- Costs of SVD and DCT are marginal

# Experiment Highlights: Speedup

Running time with one single input

|  | Priv-only | 3LegRace | Slalom | Delta |
|---|---|---|---|---|
| Train (ms/speedup) | 1372 | 237 (6×) | - | 62 (22×) |
| Inference (ms/speedup) | 510 | 95 (5×) | 84 (6×) | 20 (25×) |

3LegRace [Niu, et al, PETs 2022]: layer-wise feature decomposition on linear layers
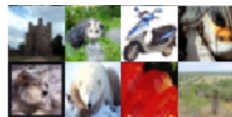Slalom [Tramer, et al, ICLR 2019]: layer-wise computation distribution on linear layers

- Significant speedup compared to solely using private envs
- Faster compared to baselines due to reduced communication
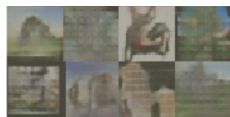
# Experiment Highlights: Protection Against Attacks

<u>Procedure</u>: Train a GAN with the quantized residuals
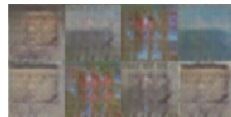<u>Setting</u>: ResNet-18, CIFAR-100

Against model inversion attack [`SecretRevealer`, CVPR'20]
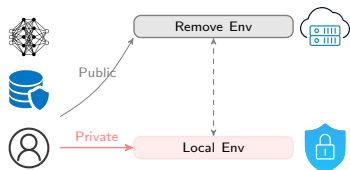


Original samples       Reconstruction (no noise)       Reconstruction ($\epsilon = 1$)

- Attack can succeed on certain samples (e.g., row 1, col 3)
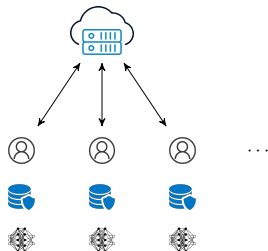- Random quantization provide further protection

# Outline

# Future Works

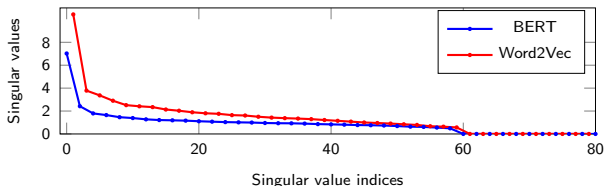## Extend to More General Settings

### User-Server Setting



### Federated Setting

## Extend to LMs

LMs' embedding also exhibits a low-ranks structure



Original text (top) and approximated (bottom) text with 1/5 principal vectors.

Large Language Models are foundational machine learning models that use deep learning algorithms to process and understand natural language. These models are trained on massive amounts of text data to learn patterns and entity relationships in the language.

Large Language Models can perform many types of language tasks, such as translating languages, analyzing sentiments, chatbot conversations, and more. They can understand complex textual data, identify entities and relationships between them, and generate new text that **is** coherent and grammatically accurate.

Large Language Models are foundational machine learning models that use deep learning algorithms to process and understand natural language. These models are trained on massive amounts of text data to learn patterns and entity relationships in the language.

Large Language Models can perform many types of language tasks, such and translating languages, analyzing sentiments, chatbot conversations, and more. They can understand complex textual data, identify entities and relationships between them, and generate new text that **are** coherent and grammatically accurate.

Questions?