## DP Training & Synthe

Felix Zhou, Samson Zhou, Vahab Mirrokni, Alessandro Epaste

## **Differential Privacy**

A (randomized) algorithm  $\mathcal{A}$  is  $(\varepsilon, \delta)$ -differentially private if for all neighboring datasets D and D' and all sets of possible outputs *Y*:

$$e^{-\varepsilon} - \delta \leq \frac{\Pr[\mathcal{A}(D) \in Y]}{\Pr[\mathcal{A}(D') \in Y]} \leq e^{\varepsilon} + \delta$$



Epasto, Vincent Cohen-Addad	Istering Vale
Example: CIFAR-10 (Above+Right) • Left: Synthetic CIFAR-10 data • $\varepsilon = 8, \delta = 1e - 5$	Results: Synthetic Data (Below) Task: Generate DP synthetic images
<ul> <li>Right: CIFAR-10 training data</li> <li>Each row is a single class</li> </ul>	<ul> <li>Generate DP synthetic images</li> <li>Train ResNet50 on images</li> <li>Test directly on original test set</li> <li>Better results than DP synthetic data baselines! <ul> <li>Private evoluation [LGK+24]</li> <li>DP-Diffusion [GBG+23]</li> </ul> </li> </ul>
<ul> <li>Results: Private Training (Below)</li> <li>Task: Train DP image classifier</li> <li>Generate DP synthetic embeddings</li> </ul>	
<ul> <li>Train 2-layer neural network on embeddings</li> <li>Test on embeddings of original test set</li> <li>Bottor results than DB finatuning baseling!</li> </ul>	90 - 85 -
• [DBH+22] DATASET $\varepsilon$ SOTA OURS CIFAR-10 8 96.6 97.0 + 0.01	80 - 75 - 75 - 75 -
CIFAR-1008 $81.8$ $80.5 \pm 0.177$ CIFAR-100891.1 $93.1 \pm 0.067$	$   \begin{bmatrix}             2 & 4 & 6 & 8 & 10   \end{bmatrix}   $

Epsilon









## [ethodology]

**vel Approach**: DP Gaussian Mixture Model embeddings

Encode the images to embedding space

Run DP Clustering on embeddings to obtain centers (means)

Privately estimate covariance per cluster

Sample (unlimited) new "embeddings" from Gaussian mixture model

Decode "embeddings" into images

Work done while Felix Zhou was a student researcher at Google