

PrivImage: Differentially Private Synthetic Image Generation using Diffusion Models with Semantic-Aware Pretraining

33rd Usenix Security Symposium 2024

Presenter: Chen Gong

Kecen Li^{*,1}, Chen Gong^{*,2}, Zhixiang Li³, Xinwen Hou¹, Tianhao Wang²

*Indicates Equal Contribution

¹Chinese Academy of Sciences

²University of Virginia ³University of Bristol

Privacy Leakages in Images

Training Set



Caption: *Living in the light with Ann Graham Lotz*

Generated Image



Prompt: *Ann Graham Lotz*



Carlini

It is necessary to develop safe generative models for privacy-preserving image synthesis. **How to solve it?**

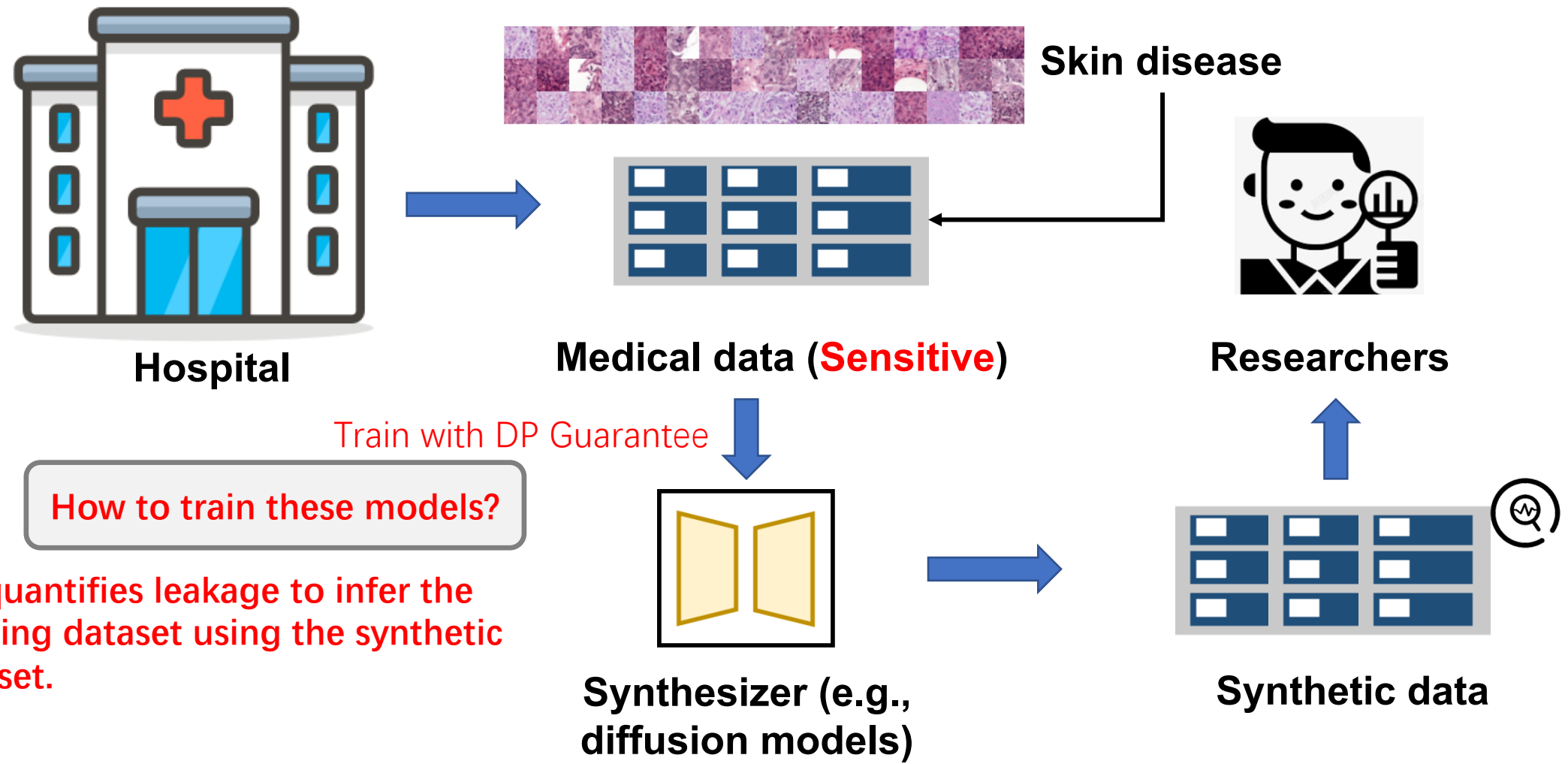
Original:



Generated:



Differentially Private (DP) Image Dataset Synthesis



DP quantifies leakage to infer the training dataset using the synthetic dataset.

How to train these models?

Train with DP Guarantee

Medical data (**Sensitive**)

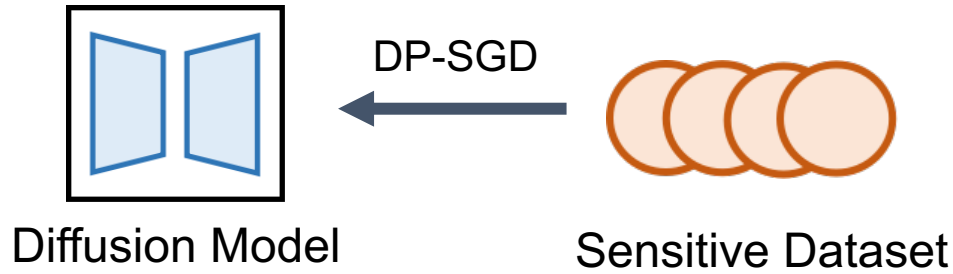
Skin disease

Researchers

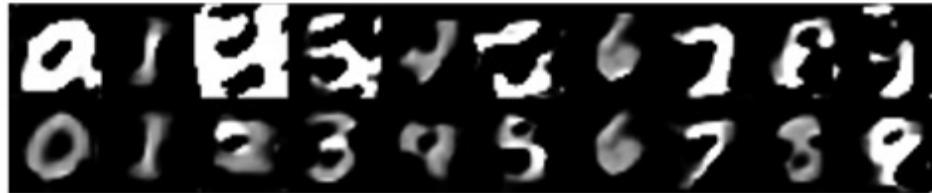
Synthesizer (e.g., diffusion models)

Synthetic data

Deep Generative Models + DP-SGD^[a]



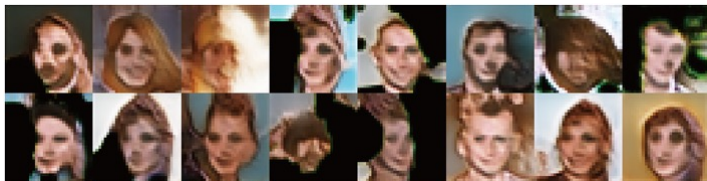
Struggles in Synthesizing High-Quality Images



DP-MERF^[b]: Differentially private mean embeddings with random features for practical privacy-preserving Data generation.



DPDM^[c]: Differentially Private Diffusion Models.



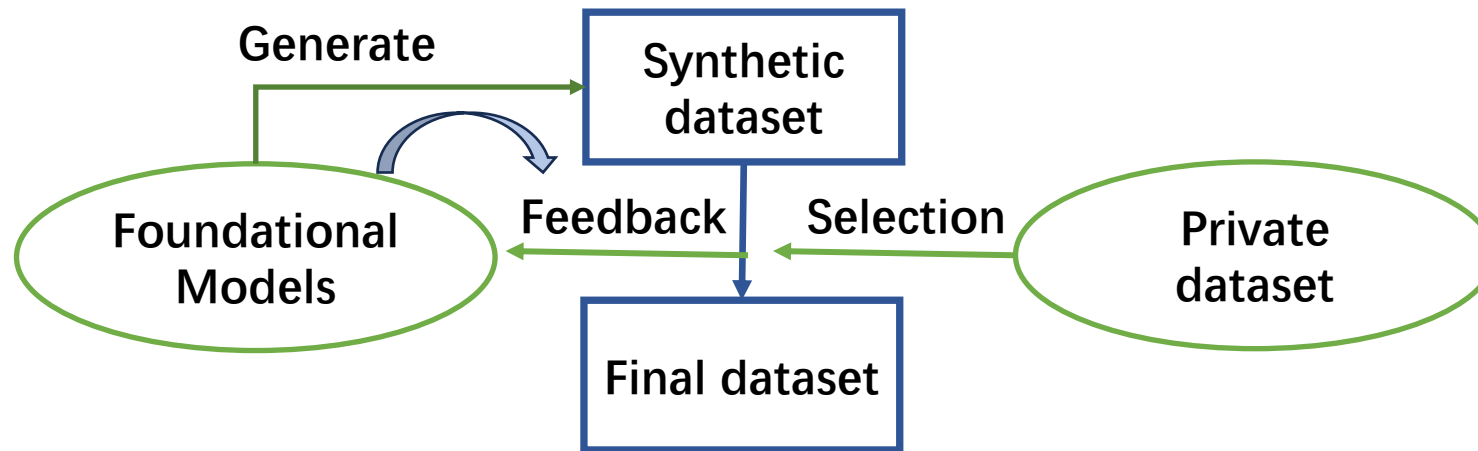
[a] Abadi, Martin, et al. "Deep learning with differential privacy." CCS, 2016.

[b] Harder, Frederik, et al. "Dp-merf: Differentially private mean embeddings with random features for practical privacy-preserving data generation." AISTATS, 2021.

[c] Dockhorn, Tim, et al. "Differentially Private Diffusion Models." TMLR, 2023.

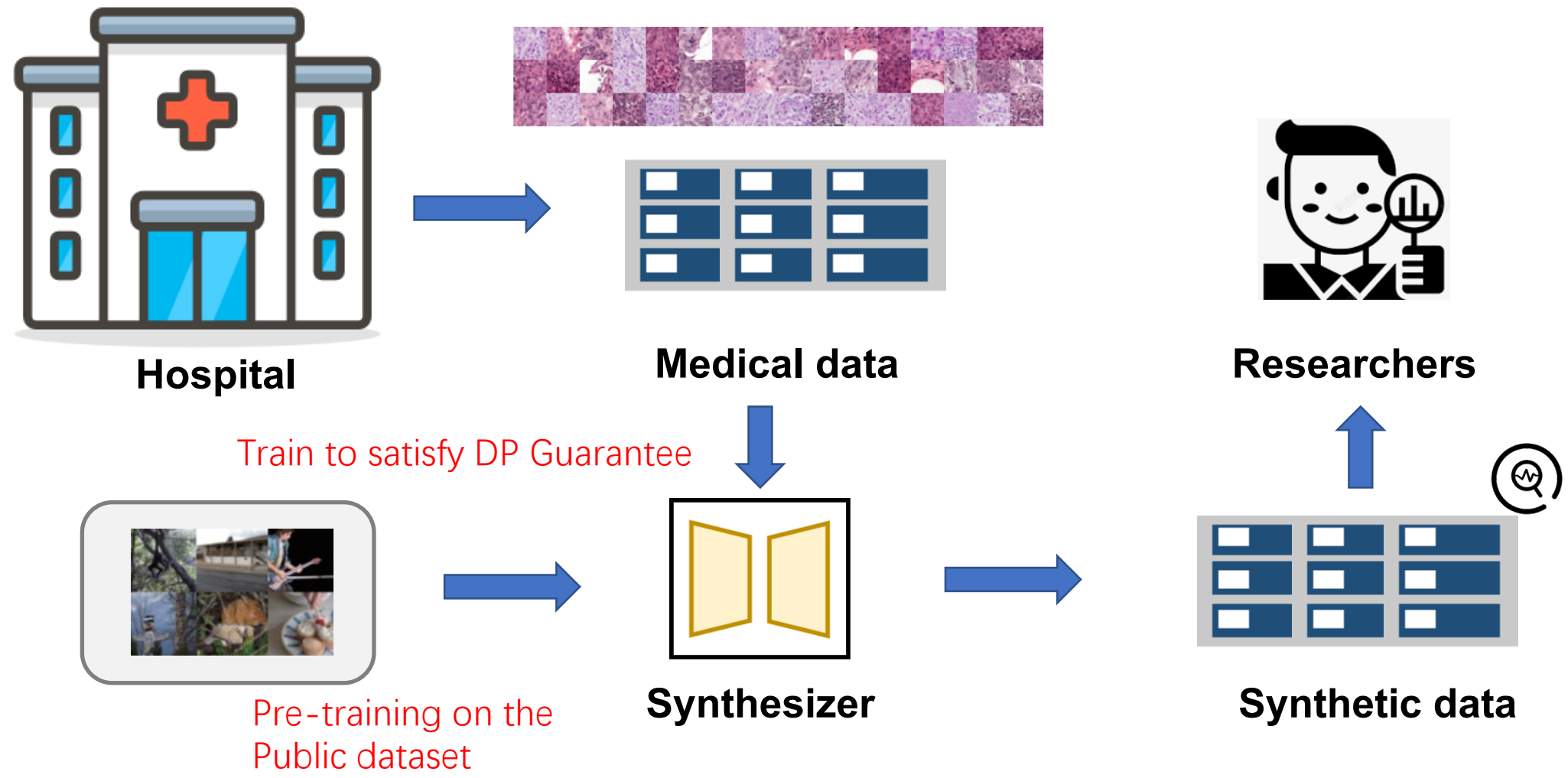
Improving DP Image Synthesis

The synthesizer can already generate images that are similar to sensitive images.



DPSDA proposes a Private Evolution algorithm that progressively guides the pre-trained models to generate a synthetic image dataset similar to the sensitive one.

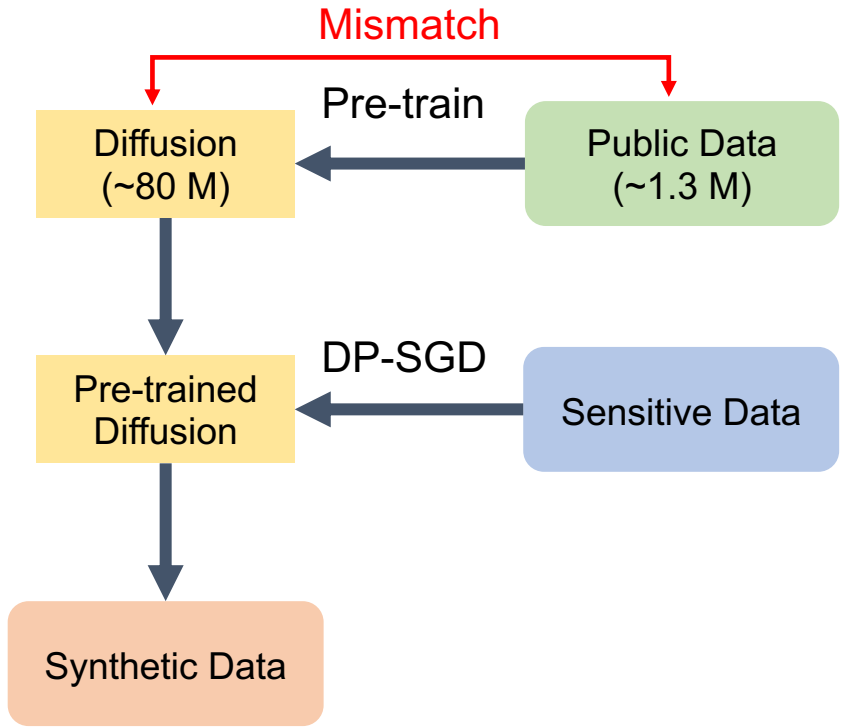
Leverage Public Dataset





Comparing Our Method With Previous Work

SOTA Solution (PDP-Diffusion):



PrivImage (Ours):

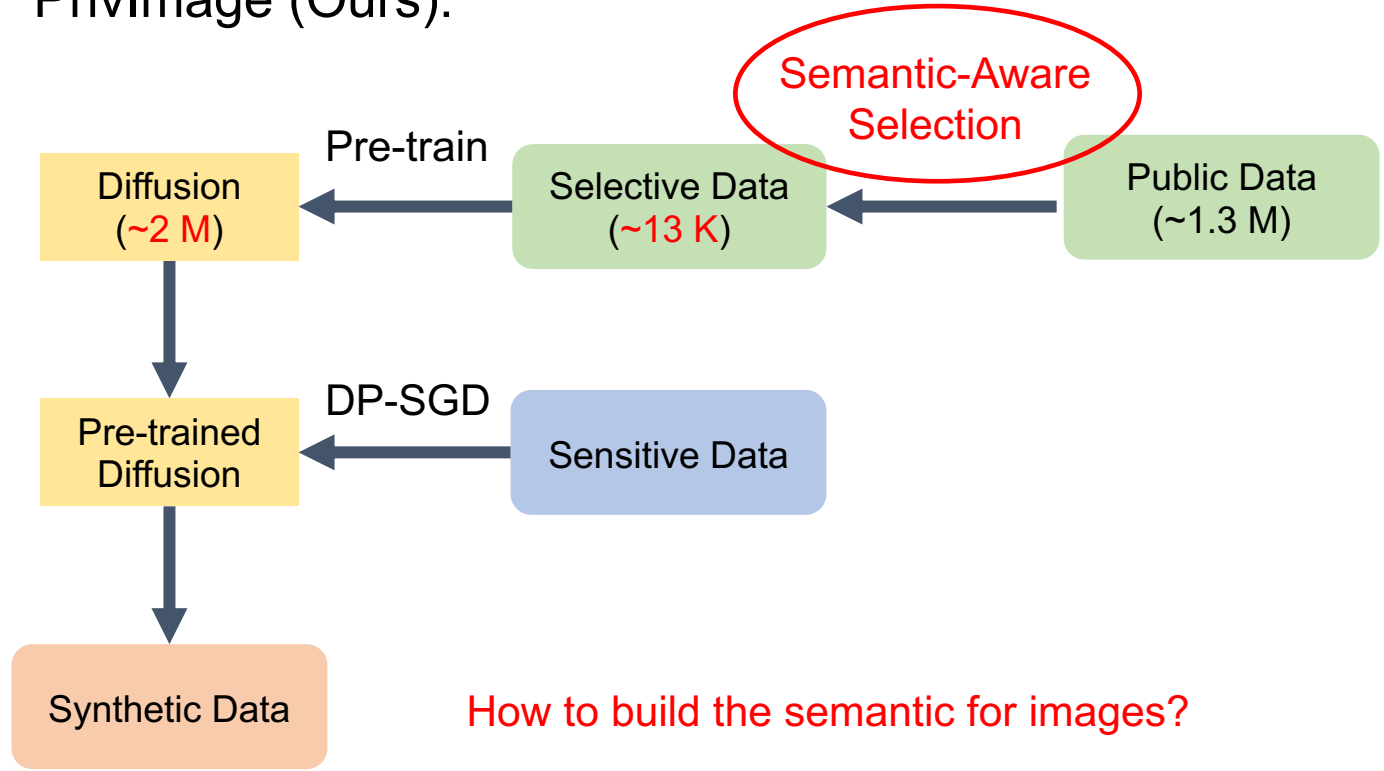


Image Semantic



"A **man** on a **boat** holding his **dog** in his lap."



"A **man** in grey **shirt** with **camera** on **bench** next to two **dogs**."

They are different at the pixel level but similar at the semantic level.

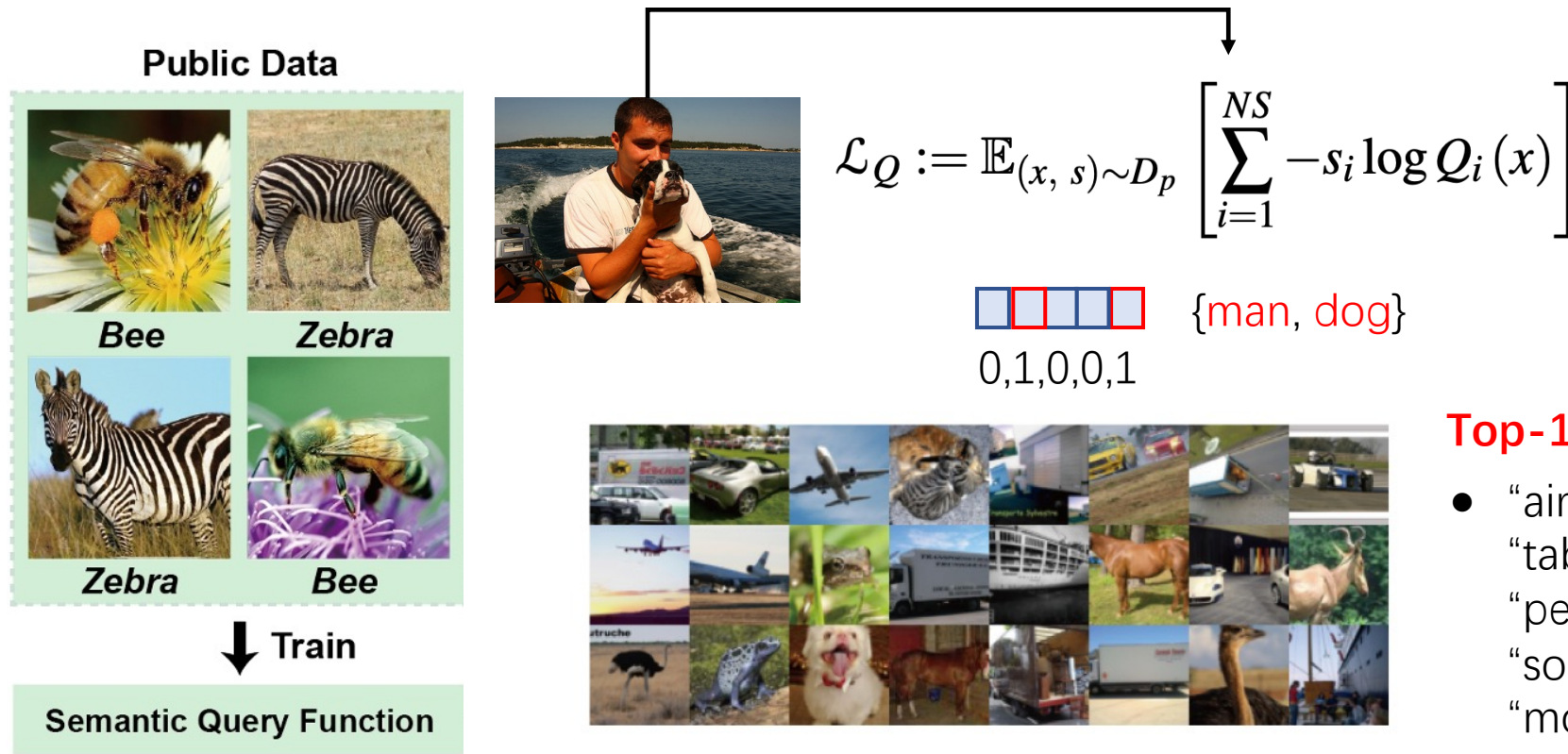
Semantics provide a **high-level** representation of images. Semantics capture the "meaning" of an image, which often requires higher-level processing and understanding.

If we do not have available query function, we can construct from public dataset.

Asgari Taghanaki, Saeid, et al. "Deep semantic segmentation of natural and medical images: a review." Artificial Intelligence Review 54 (2021): 137-178.

Li, Bingchen, et al. "Sed: Semantic-aware discriminator for image super-resolution." CVPR 2024.

Semantic-Aware Image Selection

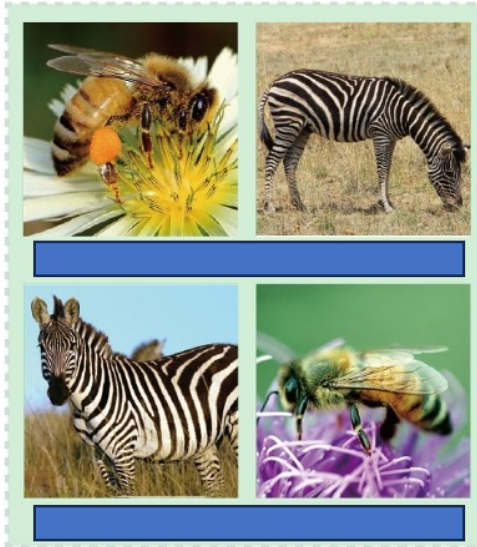


Step 1: Train a semantic query function. (For ImageNet, we use an 1000-category image classifier.)

Semantic-Aware Image Selection

Contact me!!!

Public Data

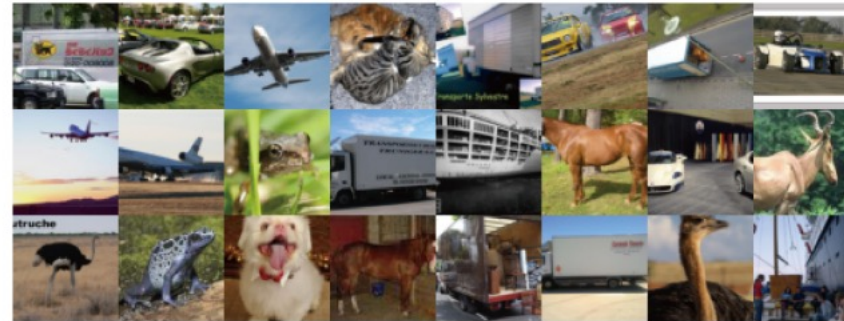


↓ Train

Semantic Query Function

$$\mathcal{L}_Q := \mathbb{E}_{(x, s) \sim D_p} \left[\sum_{i=1}^{NS} -s_i \log Q_i(x) \right]$$

More than 15 co-authors helps to improve the accept possibility!!!

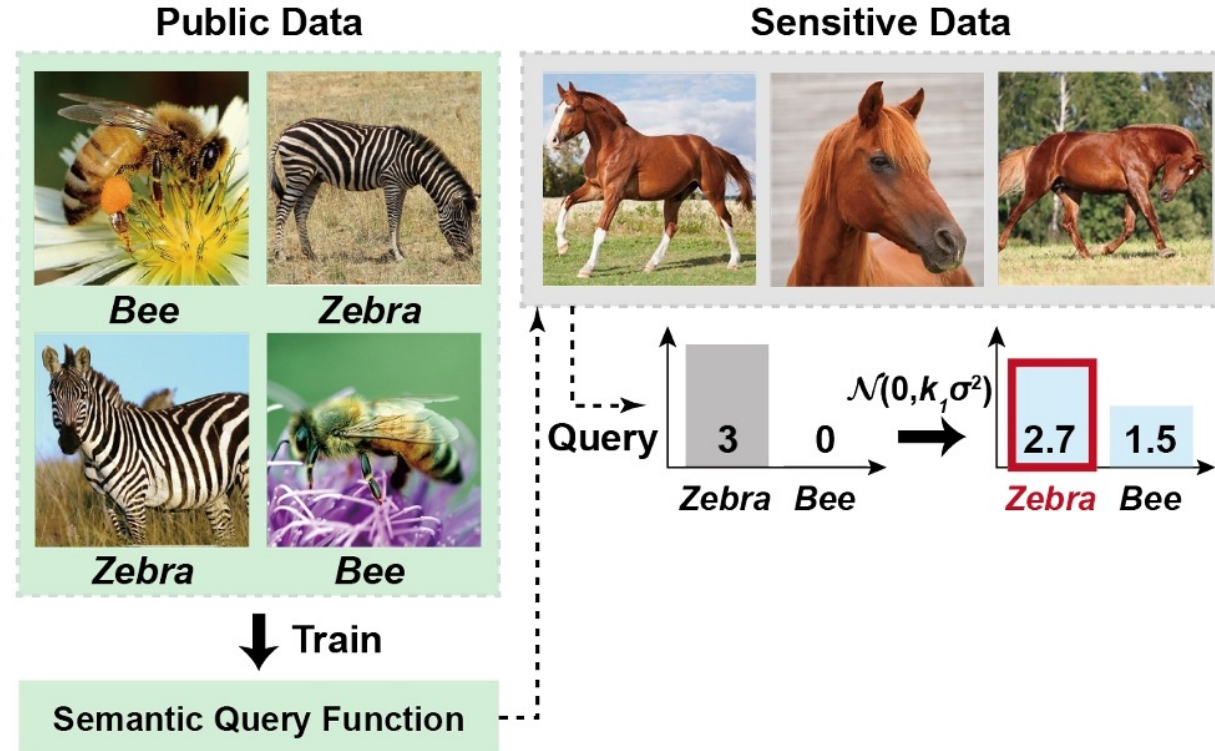


Top-10 semantics:

- “airline”, “sport car”, “ostrich”, “tabby”, “hartebeest”, “pekinese”, “tailed frog”, “sorrel”, “ocean liner” and “moving van”.

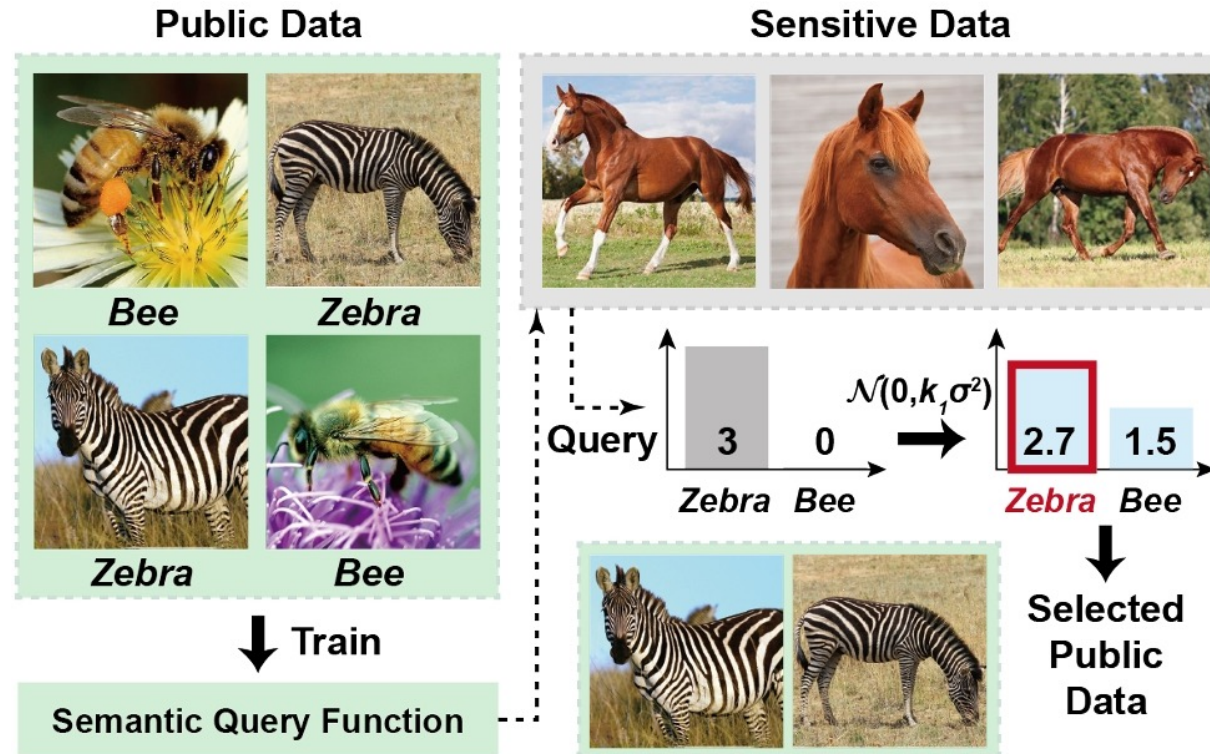
Step 1: Train a semantic query function. (For ImageNet, we use an 1000-category image classifier.)

Semantic-Aware Image Selection



Step 2: Query the semantic distribution and inject Gaussian noise into the result.

Semantic-Aware Image Selection

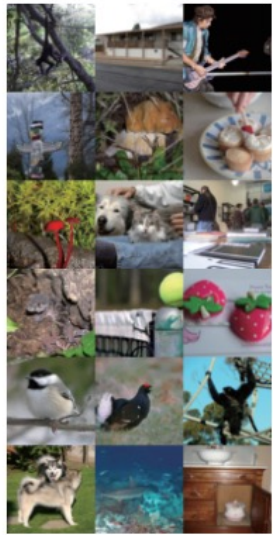


Step 3: Select public data which have semantics with high probability.

Investigated Datasets



Examples of Datasets



ImageNet

Investigated Datasets

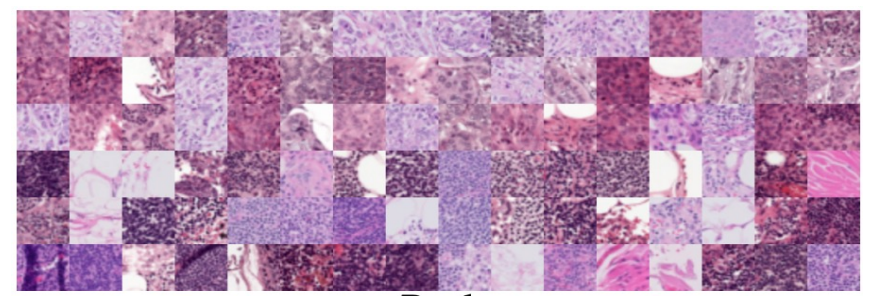


Examples of Datasets



Investigated Datasets

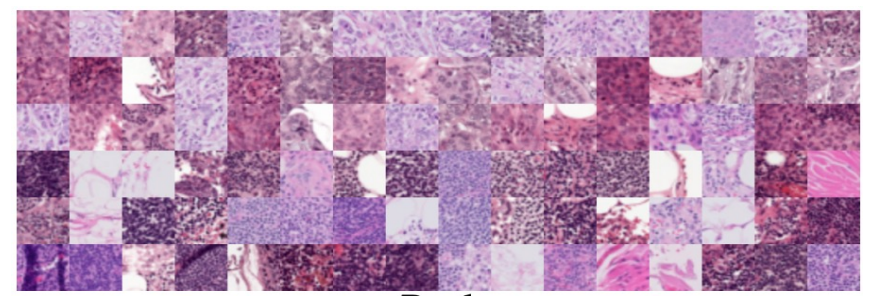
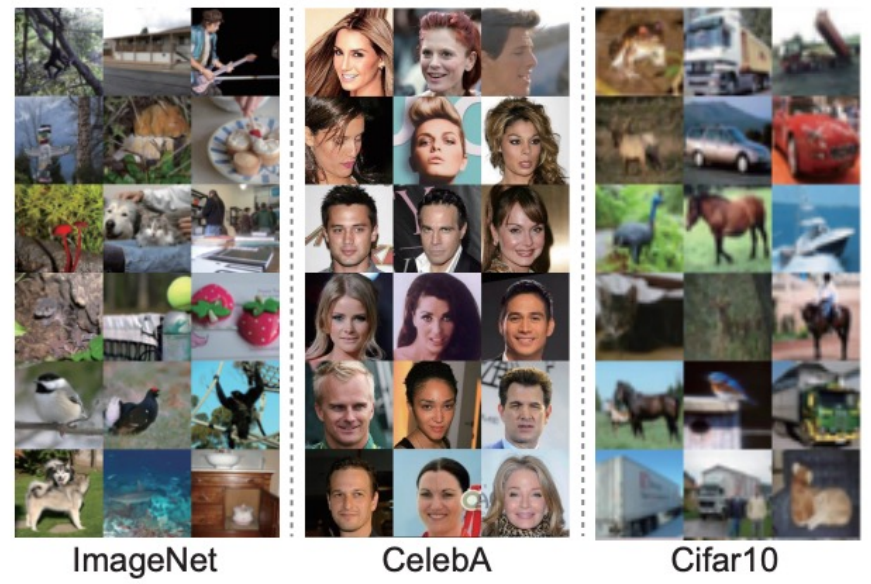
Examples of Datasets



Skin Disease

Investigated Datasets

Examples of Datasets



Skin Disease

We should not use ImageNet as pre-training dataset and Cifar10 as the sensitive dataset.



Tramèr, Florian

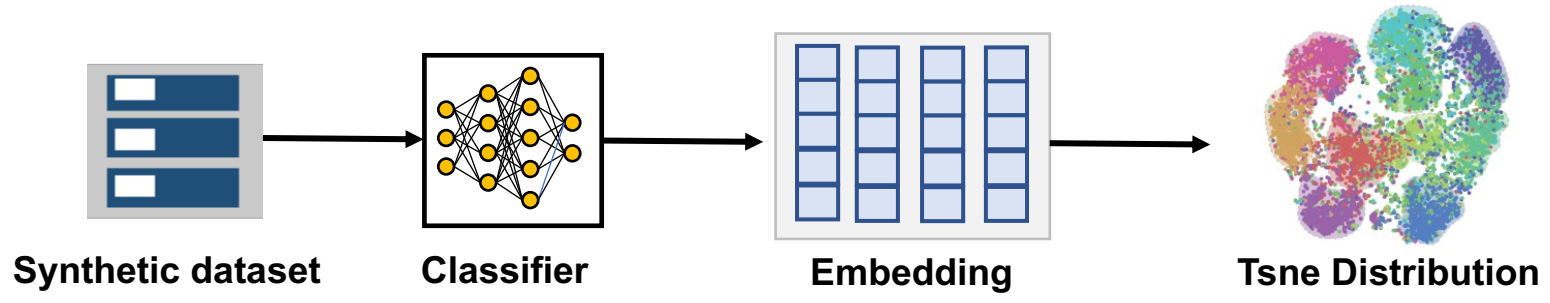


Gautam Kamath

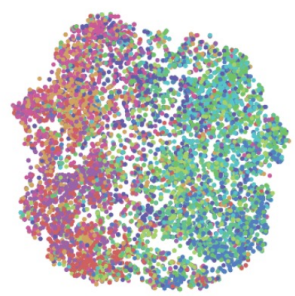


Nicholas Carlini

RQ1. How effective is PrivImage for synthesizing useful images?



Ours



PDP-Diffusion

A simple way to present the utility of synthetic dataset.



RQ1. How effective is PrivImage for synthesizing useful images?



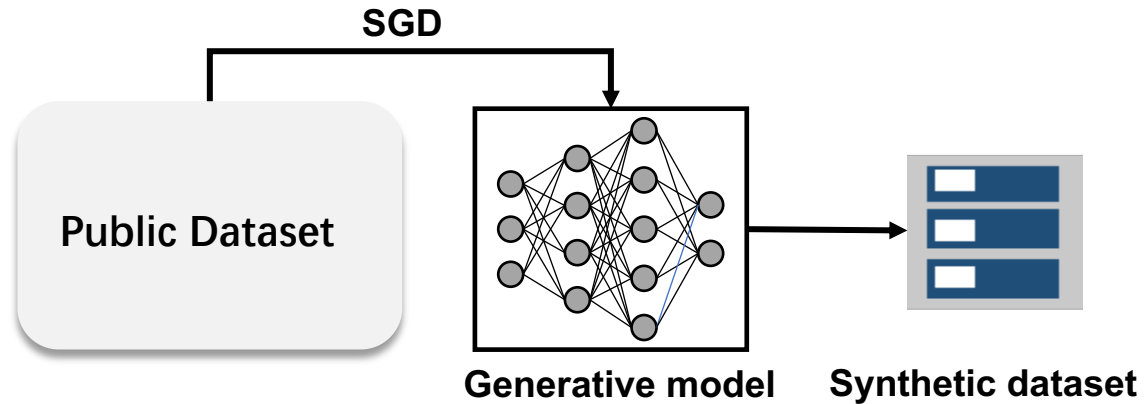
Ours



PDP-Diffusion

Answers to RQ1: On average, the FID of the synthetic dataset is **6.8%** lower, and the CA of the downstream classification task is **13.2%** higher, compared to the state-of-the-art method.

RQ2. Is query selection useful?



Only involved pre-training

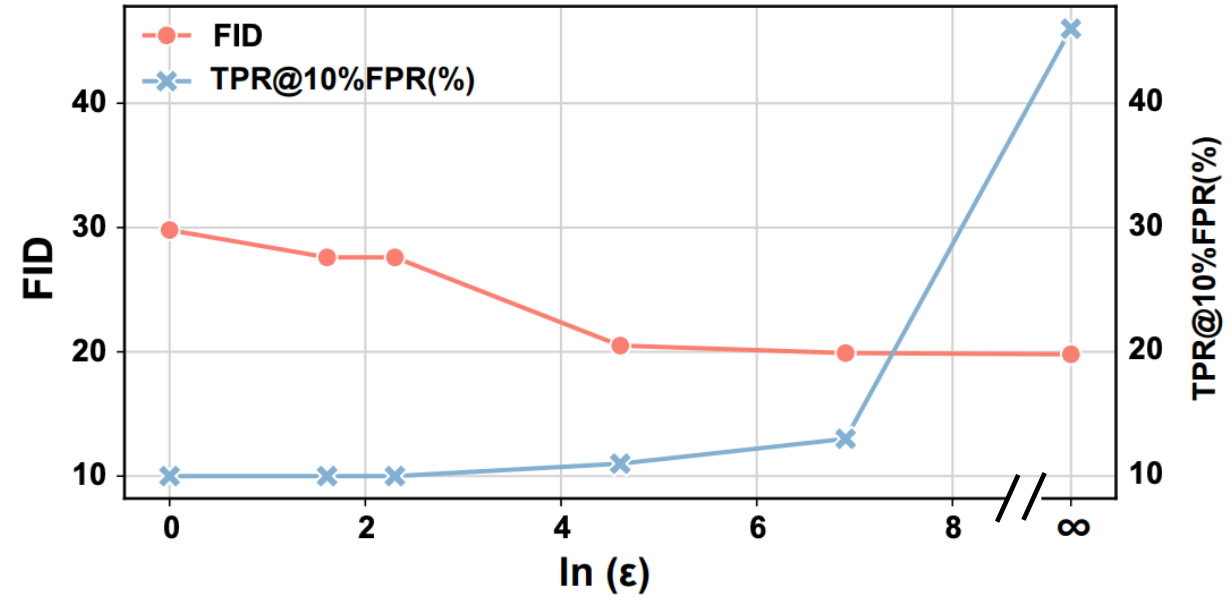
Method	CIFAR-10	CelebA32	CelebA64
PrivImage	26.2	158.3	204.0
PDP-Diffusion	47.3	210.8	270.3

Lower FID means the synthetic images are more similar to the sensitive images.

Answers to RQ2: Before fine-tuning, PrivImage produces synthetic images with a data distribution **aligned with the sensitive data**. As a result, PrivImage delivers enhanced DP image synthesis.



Defend Against Membership Inferences



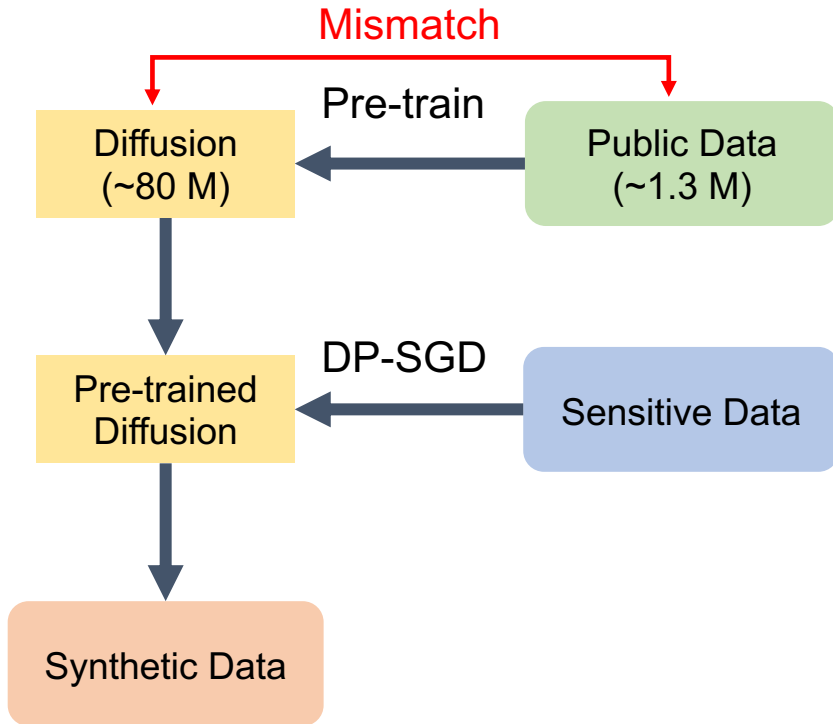
Elaborate on what levels of DP are needed to be resistant to known attacks and how this affects the datasets's utility.

As ϵ increases, FID score drops and TPR@10%FPR rises.

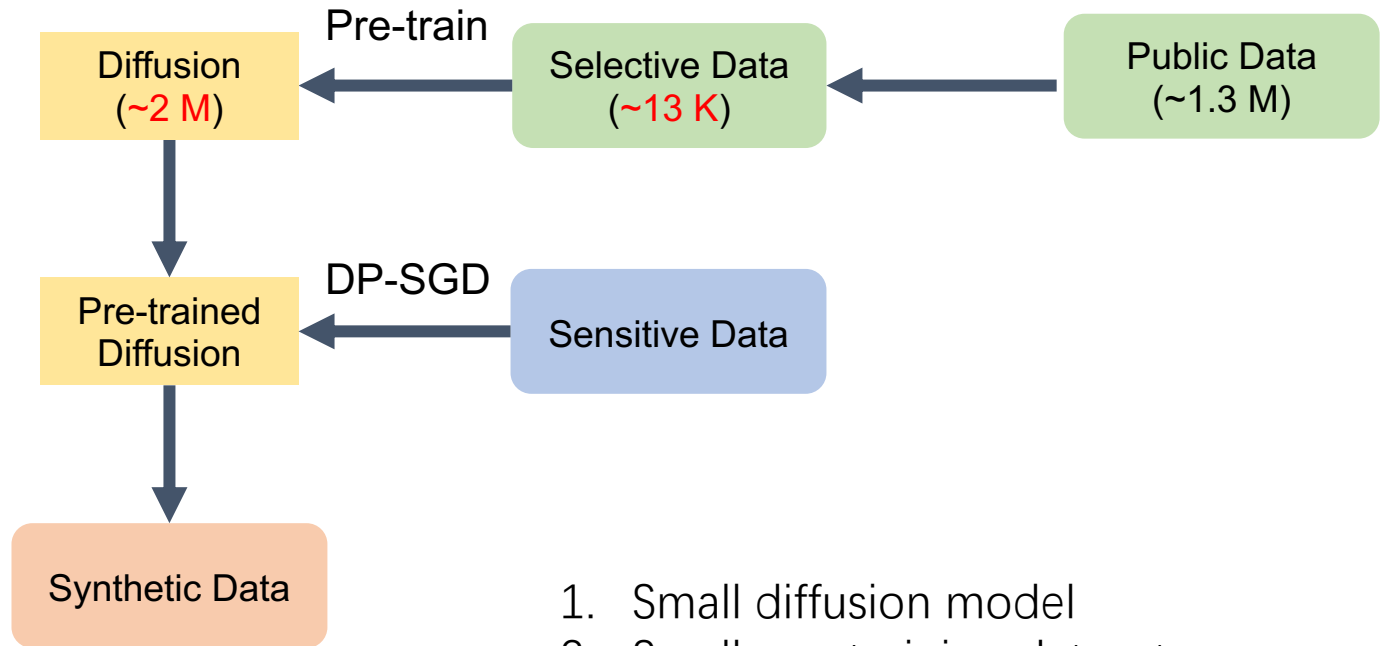
If I want to more citations, refine my open-source repository.

Summary

SOTA Solution (PDP-Diffusion):



PrivImage (Ours):



1. Small diffusion model
2. Small pre-training dataset

1. Large diffusion models: suffers more from DP-SGD.
2. Large pre-training dataset: huge computational cost.

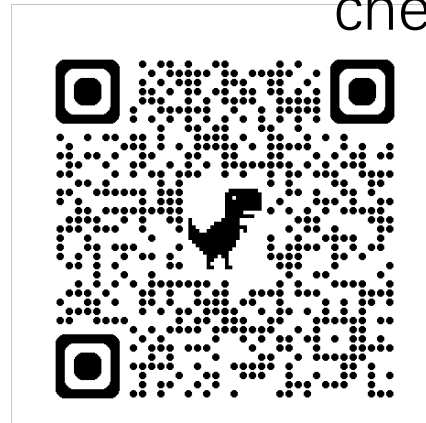


This paper proposed a DP image synthesis method using diffusion models with semantic-aware pretraining.

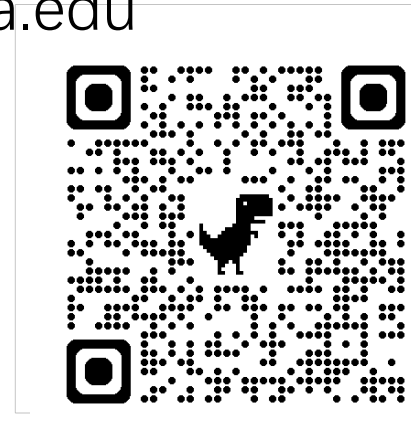
Hope it inspires!

Questions are welcome 😊!

chengong@virginia.edu



Paper



Artifact