



Towards More Realistic Membership Inference Attacks on Large Diffusion Models

Antoni Kowalczyk^{1,2*}, Jan Dubiński^{1*}, Stanisław Pawlak¹,
Przemysław Rokita¹, Tomasz Trzcíński^{1,3,4}, Paweł Morawiecki⁵

*equal contribution

1



2



3

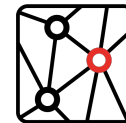


4

IDEAS
NCBR ○ ○ ○

5





Motivation

- Data Privacy
- Copyright issues
- Novel research topic



World ▾ Business ▾ Markets ▾ Sustainability ▾ More ▾



R



My View



Following



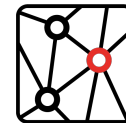
Saved

Litigation | Copyright | Litigation | Technology | Intellectual Property

Getty Images lawsuit says Stability AI misused photos to train AI

By **Blake Brittain**

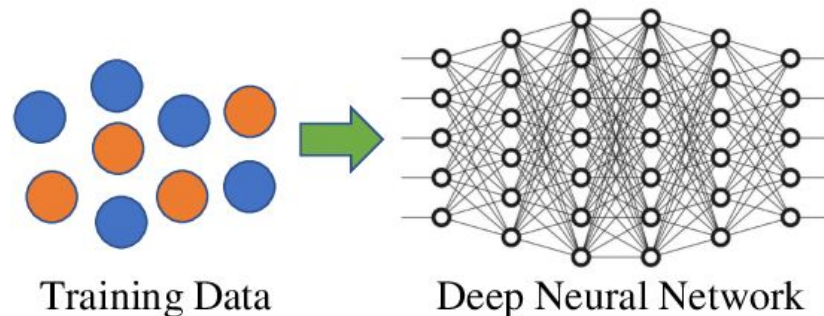
February 6, 2023 6:32 PM GMT+1 · Updated 9 months ago



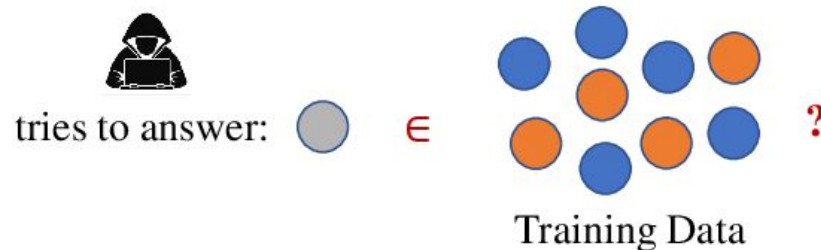
Membership Inference Attacks

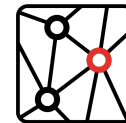
Was this example in the training set?

Training of Target Model



Membership Inference Attack on Target Model





Data

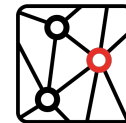
Members

- Train set
- Identifying them is **the goal** of a MIA
- Potentially: copyrighted artwork

Nonmembers

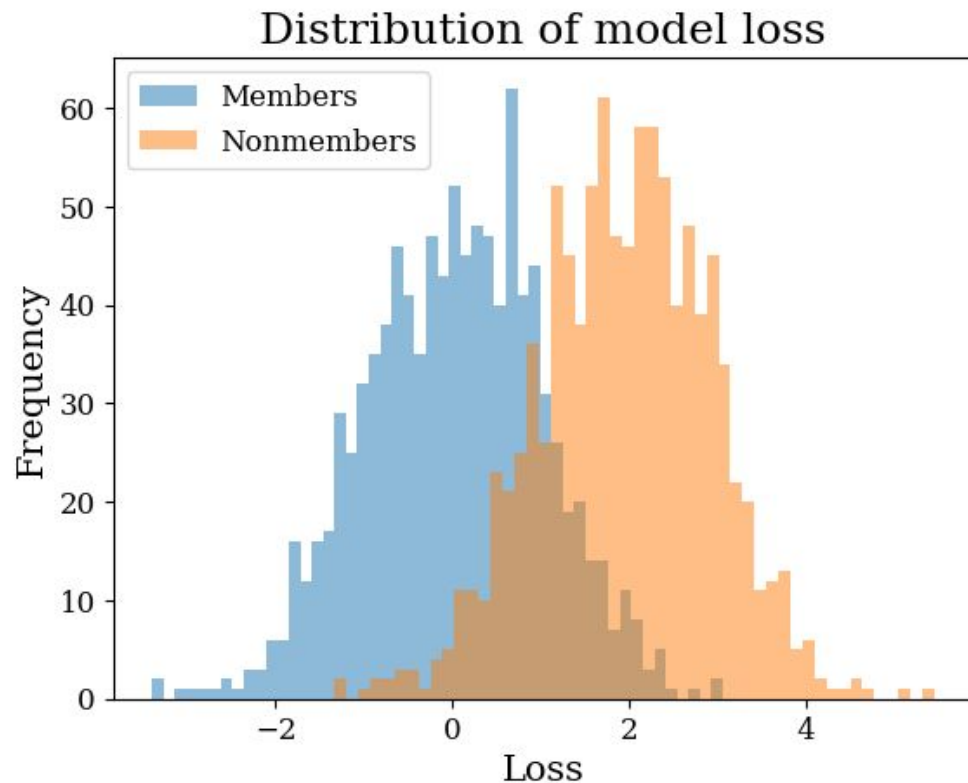
- Not used during training
- **Ideally: validation/test set**

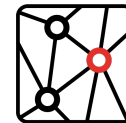
**MIA is a *classification*
task**



Loss Threshold Attack

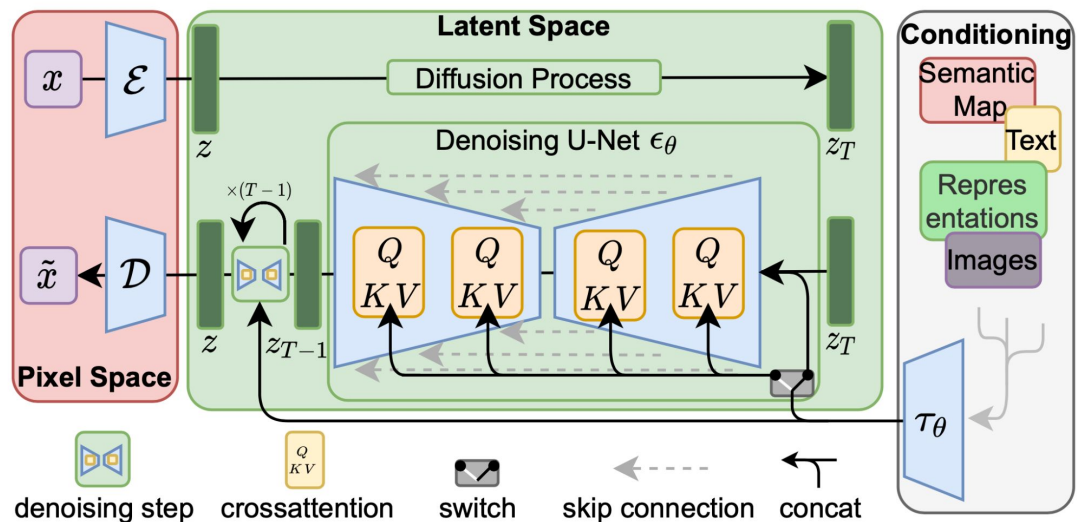
- If $\text{loss} < \text{threshold}$ then sample is a member
- Due to overfitting
- $\text{TPR@FPR}=1\%$

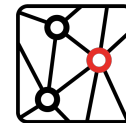




Stable Diffusion

- Large Latent Diffusion Model
- LAION-5B
- Fully open-source*
- SOTA Text2Image**



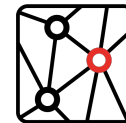


Issues

- Current cost of training: **\$100k**
- **150 000** GPU hours
- **No validation set!**



**Problem: We *do not* have
a natural nonmembers
set!**



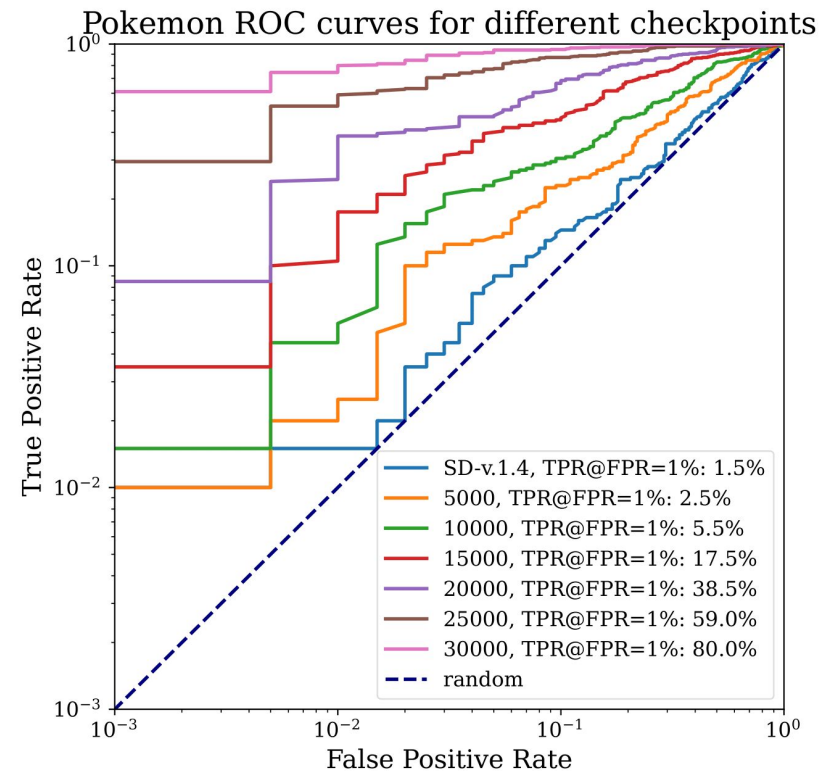
Option 1: Fine-tune Stable Diffusion on a smaller dataset

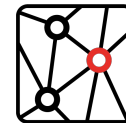
Pros

- cheap
- nonmembers easy to get
- easy to benchmark

Cons

- flawed
- trivial problem
- not applicable to real life scenarios





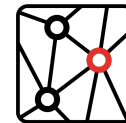
Option 2: Train Stable Diffusion from scratch

Pros

- nonmembers easy to get
- correct experimental setup

Cons

- extremely expensive
- impractical



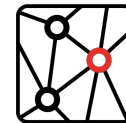
Option 3: Obtain nonmembers from a different source

Pros

- cheap
- data easy to collect in our case

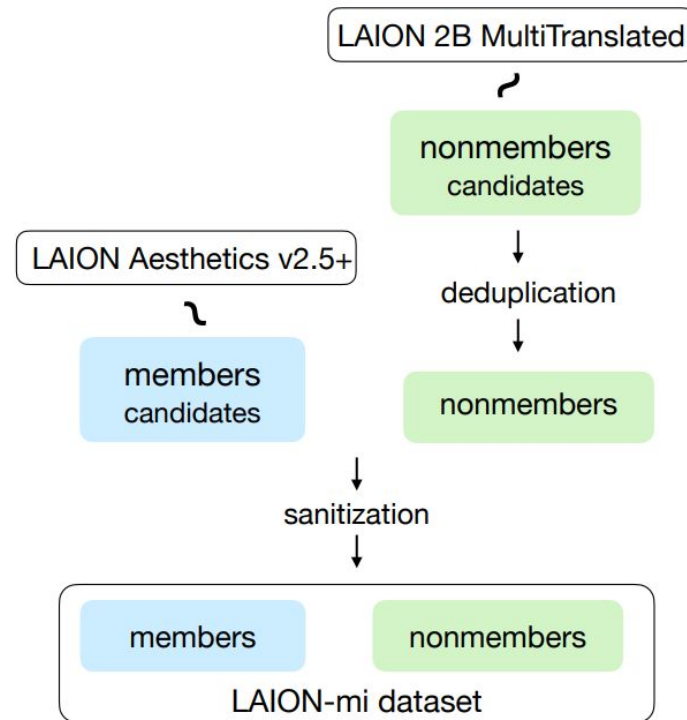
Cons

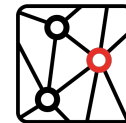
- possible distribution mismatch
- in effect: could lead to incorrect results



Solution: LAION-mi dataset

- Do not modify the original Stable Diffusion model
- Obtain the nonmembers set from other source
- Alleviate the distribution mismatch problem

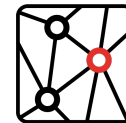




Challenge: Duplicates

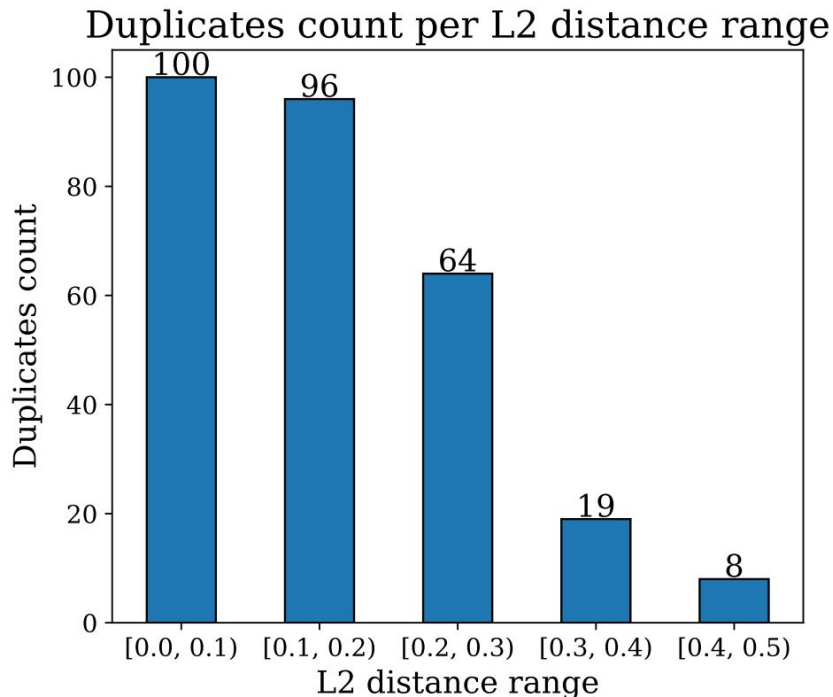
- 30% in LAION-2B EN
- Effect: nonmembers set contaminated with member samples

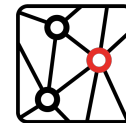




Solution: Deduplication

1. Query LAION-5B KNN-Index
2. Get duplicate candidates
3. Compute distances
4. Apply threshold to filter out duplicates

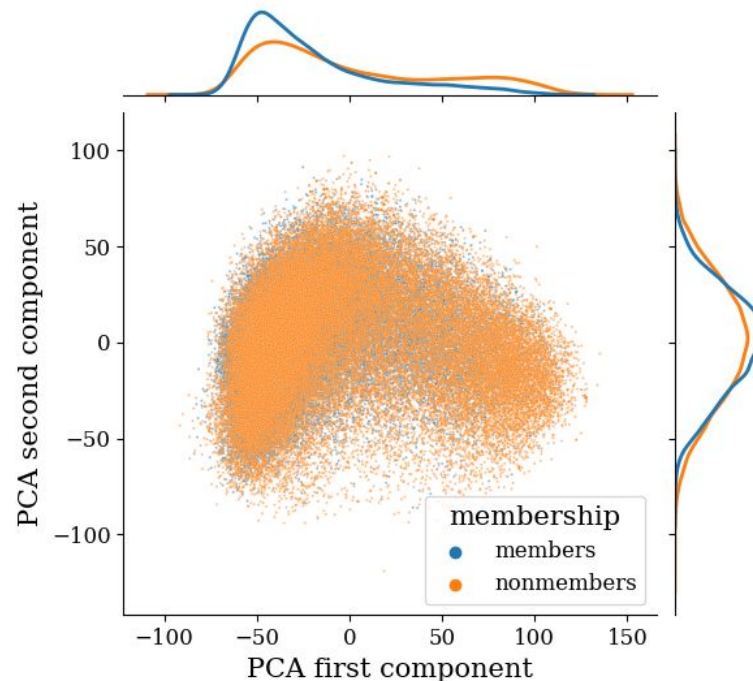


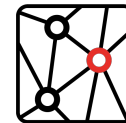


Challenge: Distribution Mismatch

- Images and descriptions
- Evaluation:
 - Visual (PCA)
 - FID
 - Training a classifier
- **Main focus on descriptions**

2D PCA plot - prompts before sanitization



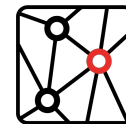


Solution: Sanitization

Start with a train set with half members and half nonmembers

Until almost random accuracy on the train set, repeat:

1. Train a classifier on a train set
2. Add a new classifier to all classifiers
3. Create a new train set by filtering out members, for which any of the classifiers classifies them as member, and use all nonmembers

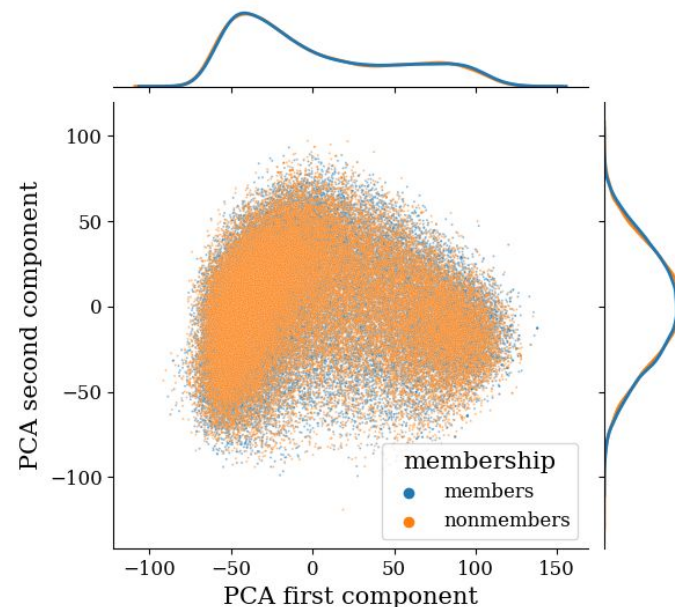


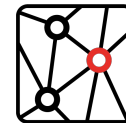
Results

LAION-mi : 40k members, 40k nonmembers

DATA SUBSET	FID	
	TEXT	IMAGES
MEMBERS INTERNAL - RANDOM	9.84	7.00
MEMBERS INTERNAL - SANITIZED	9.77	7.06
NONMEMBERS INTERNAL	9.73	7.01
COMPARATIVE - RANDOM	66.43	13.90
COMPARATIVE - SANITIZED	13.54	8.87

2D PCA plot - prompts iteration 3





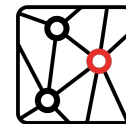
Evaluation on Stable Diffusion: Setup

Loss Threshold Attack

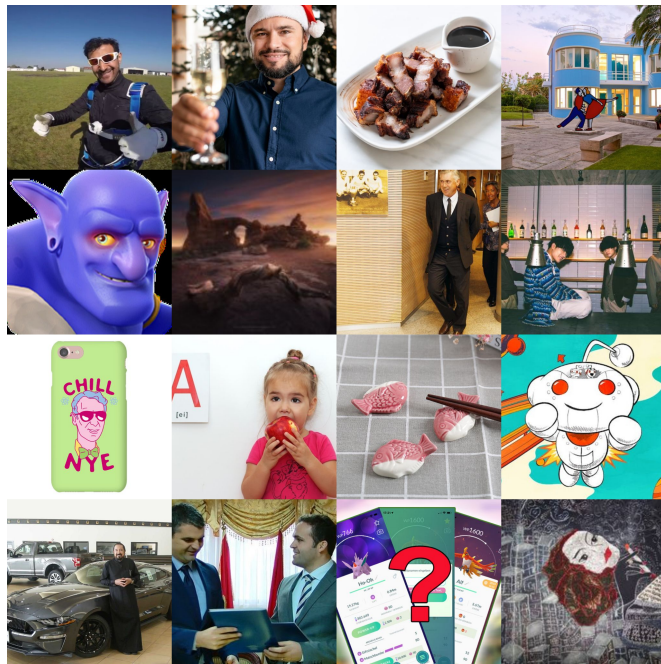
- Model Loss
- Pixel Error
- Latent Error

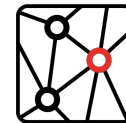
Datasets + Models

- LAION-mi + Stable Diffusion
- POKEMON + fine-tuned Stable Diffusion on POKEMON dataset



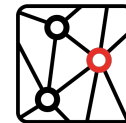
Datasets





Results

SCENARIO	LOSS	METHOD	TPR@FPR=1%. \uparrow	
			LAION-MI	POKEMON
WHITE-BOX		BASELINE LOSS THR.	1.92% \pm 0.59	80.9% \pm 2.27
	MODEL LOSS	REVERSED NOISING	2.51% \pm 0.73	97.3% \pm 0.93
		PARTIAL DENOISING	2.31% \pm 0.61	94.5% \pm 1.34
		REVERSED DENOISING	2.25% \pm 0.64	91.5% \pm 1.63
	LATENT ERROR	REVERSED NOISING	1.26% \pm 0.62	11.5% \pm 1.84
		PARTIAL DENOISING	2.42% \pm 0.62	99.5% \pm 0.4
		REVERSED DENOISING	2.17% \pm 0.64	61.1% \pm 2.74
	PIXEL ERROR	REVERSED NOISING	1.90% \pm 0.51	8.36% \pm 1.66
		REVERSED DENOISING	2.03% \pm 0.55	12.0% \pm 1.97
PARTIAL DENOISING		1.75% \pm 0.68	25.38% \pm 2.55	
GREY-BOX	LATENT ERROR	GENERATION FROM PROMPT	0.93% \pm 0.41	7.15% \pm 1.5
BLACK-BOX	PIXEL ERROR	GENERATION FROM PROMPT	0.35% \pm 0.19	12.0% \pm 1.9



Summary

- MIAs are still **hard**, or impractically **expensive**
- We point out **flawed methodology**
- Our contribution: **LAION-mi dataset & evaluation protocol**