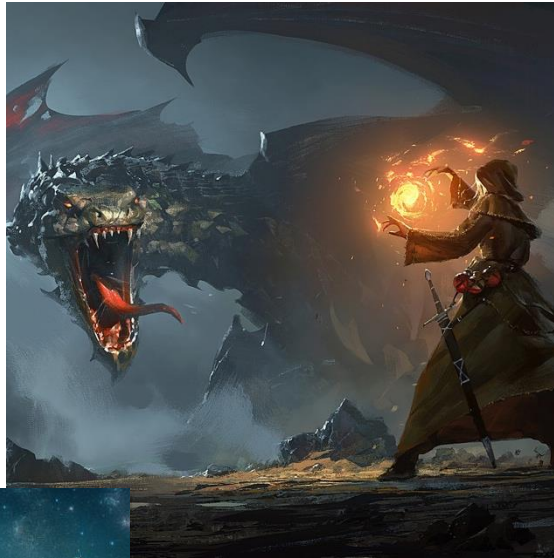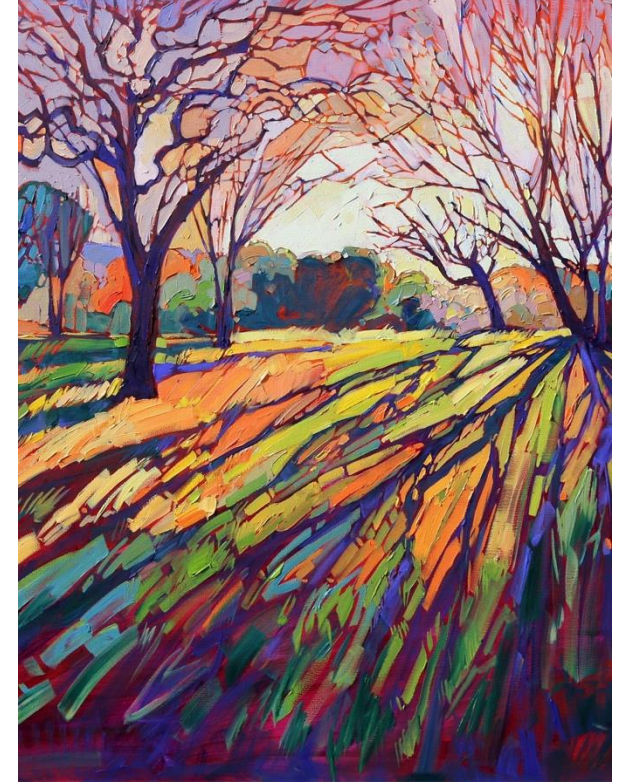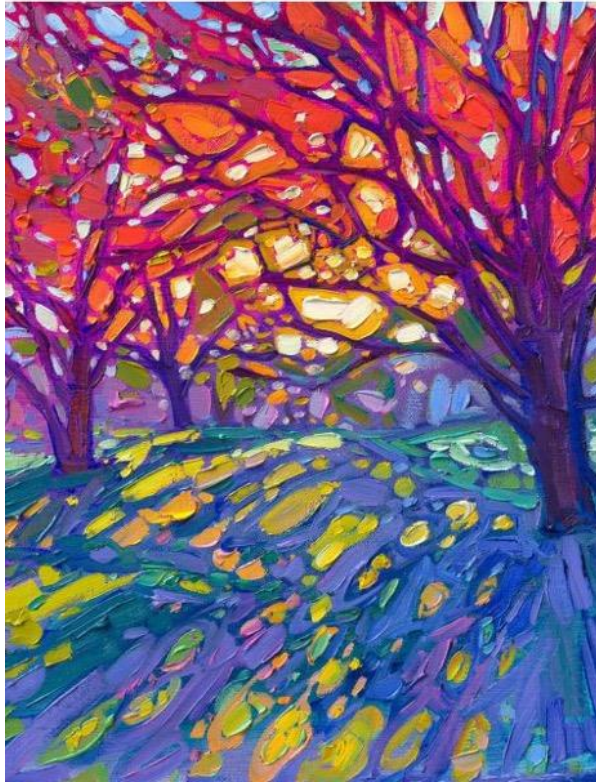# CDI: Copyrighted Data Identification in Diffusion Models

Jan Dubiński, Antoni Kowalczuk, Franziska Boenisch, Adam Dziedzic

# Diffusion Models are amazing image generators

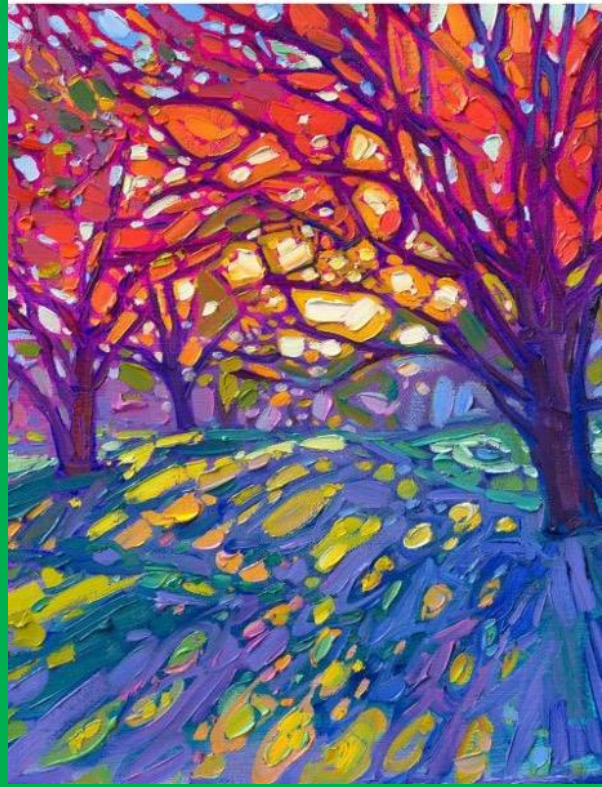# DMs are trained on billions of images



Which image is generated as
*„a forest painting in the style of Erin Hanson"*?

# DMs are trained on billions of images



*„a forest painting in the style of Erin Hanson"* by Stable Diffusion

# DMs are trained on billions of images



Real paintings by Erin Hanson

Artists and Illustrators Are Suing Three A.I. Art Generators for Scraping and 'Collaging' Their Work Without Consent

The plaintiffs claim the A.I. tools have unlawfully scraped and used their artwork in training datasets.

Artists and Illustrators Are Suing

ARTIFICIAL INTELLIGENCE / TECH / LAW

## Getty Images sues AI art generator Stable Diffusion in the US for copyright infringement

/ Getty Images has filed a case against Stability AI, alleging that the company copied 12 million images to train its AI model 'without permission ... or compensation.'

By **James Vincent**, a senior reporter who has covered AI, robotics, and more for eight years at The Verge.

Feb 6, 2023, 5:56 PM GMT+1

16  Comments (16 New)

*An illustration from Getty Images' lawsuit, showing an original photograph and a similar image (complete with Getty Images watermark) generated by Stable Diffusion.* Image: Getty Images

Artists and Illustrators Are Suing

ARTIFICIAL INTELLIGENCE / TECH /

**Getty Images Stable Diffus infringement**

An illustration from Getty Images' l
similar image (complete with Getty I
Diffusion. Image: Getty Images

**Generative AI Lawsuits Timeline: Legal Cases vs. OpenAI, Microsoft, Anthropic, Nvidia and More**

March 13, 2024 by Joe Panettieri

# Was the sample used to train the Diffusion Model?
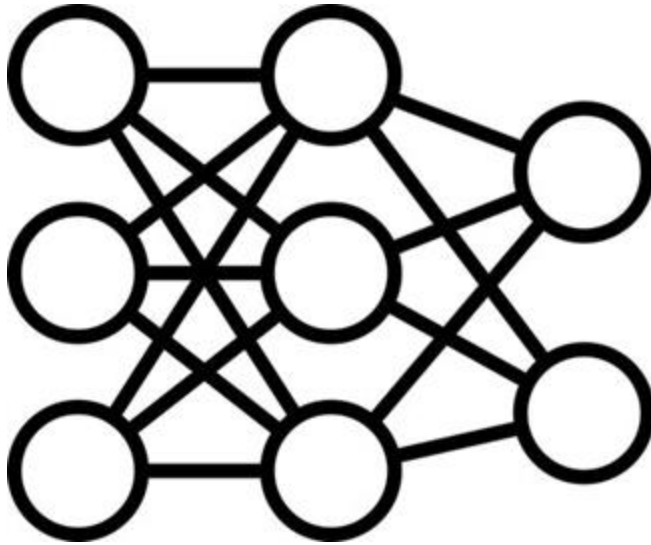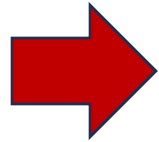


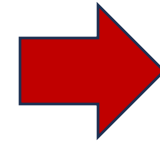Sample $s$

$\in$

Training Data $D$

# Membership Inference Attack



1. Choose sample *s*

2. Query the model *M*
with sample *s*

3. Decision:
was *s* in the train set?

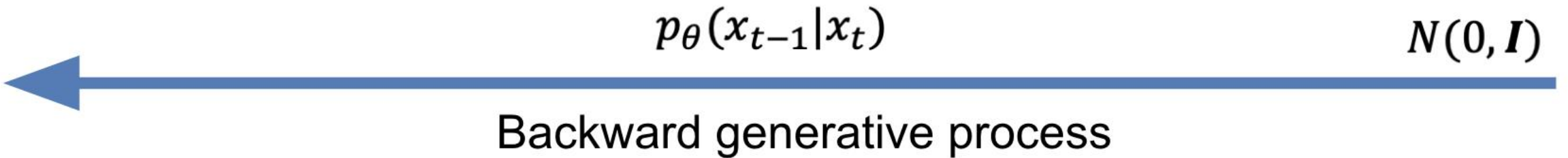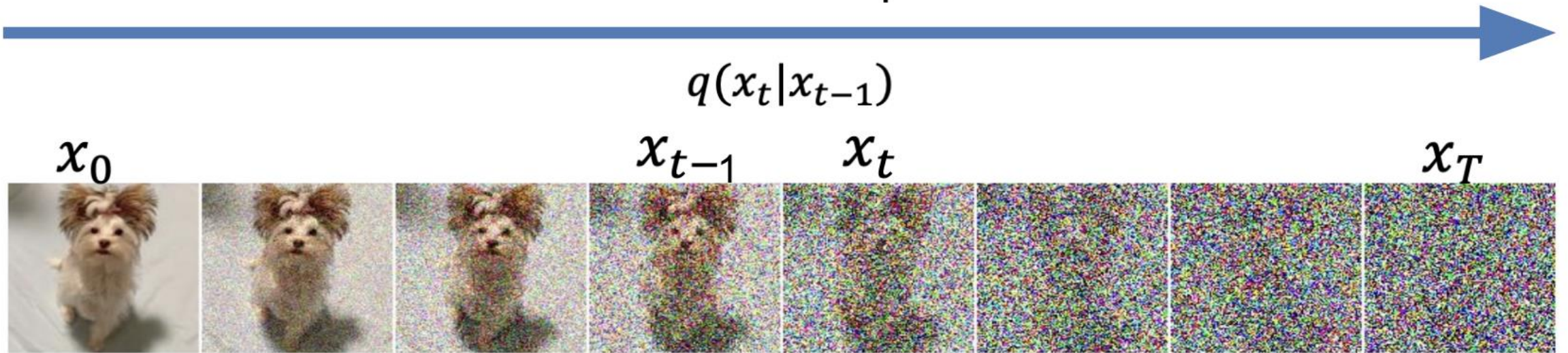# Loss threshold membership inference attack

Machine learning models minimize loss
on the training (*members*) set

$$IF\ Loss(sample) < Threshold$$
$$THEN\ Member$$
$$ELSE\ Nonmember$$

# Diffusion Models

Forward diffusion process

$$q(x_t|x_{t-1})$$



$x_0$  $x_{t-1}$  $x_t$  $x_T$

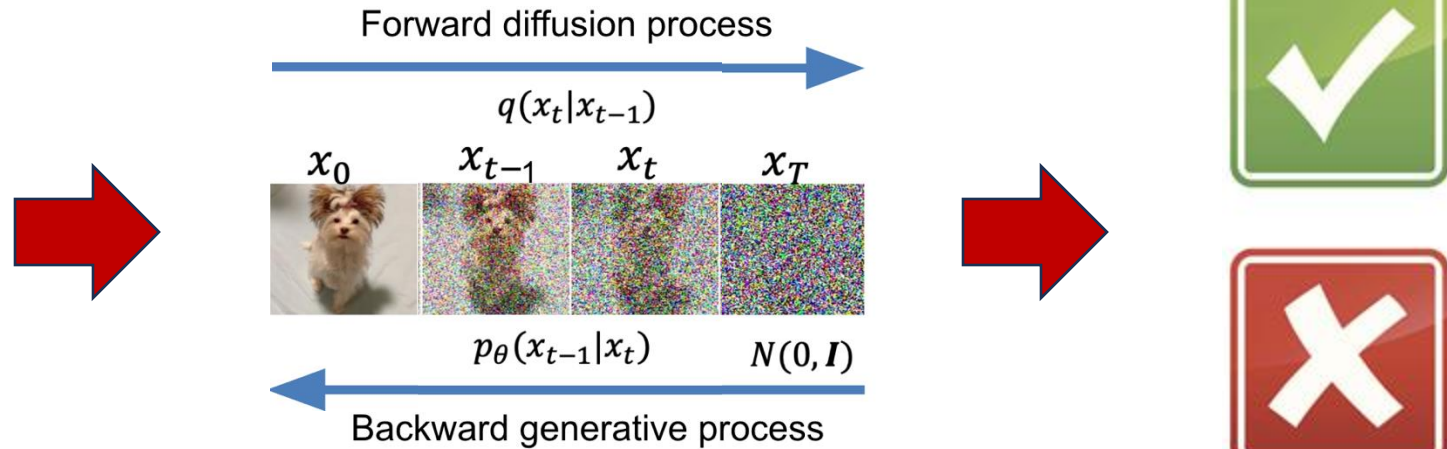$$p_\theta(x_{t-1}|x_t)$$

$$N(0, \boldsymbol{I})$$

Backward generative process

# Membership Inference Attack on Diffusion Models



1. Choose sample $s$

2. Query the model $M$ with sample $s$

3. Decision: was $s$ in the train set?

*IF loss(sample) < threshold THEN member ELSE nonmember*
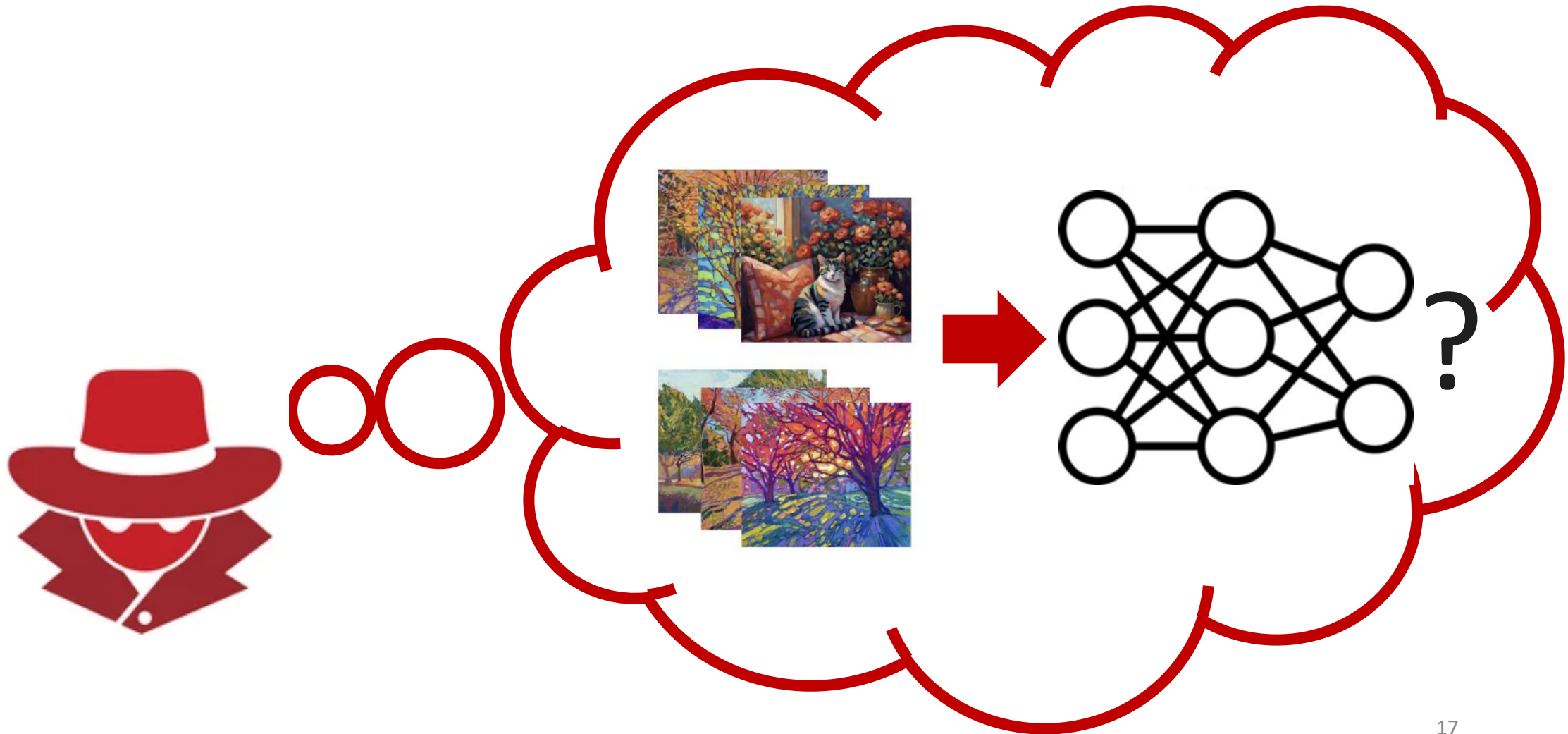
# Membership Inference for Large Diffusion Models
## Fails

| MIA | Max. TPR@FPR=1% |
|---|---|
| Denoising Loss | 2.24% |
| Secmi | 2.44% |
| PIA | 5.57% |
| PIAN | 1.53% |

# Membership Inference for Large Diffusion Models
## Fails

| MIA | Max. TPR@FPR=1% |
|---|---|
| Denoising Loss | 2.24% |
| Secmi | 2.44% |
| PIA | 5.57% |
| PIAN | 1.53% |

How can we do better?

# Dataset Inference

# Dataset Inference for Diffusion Models

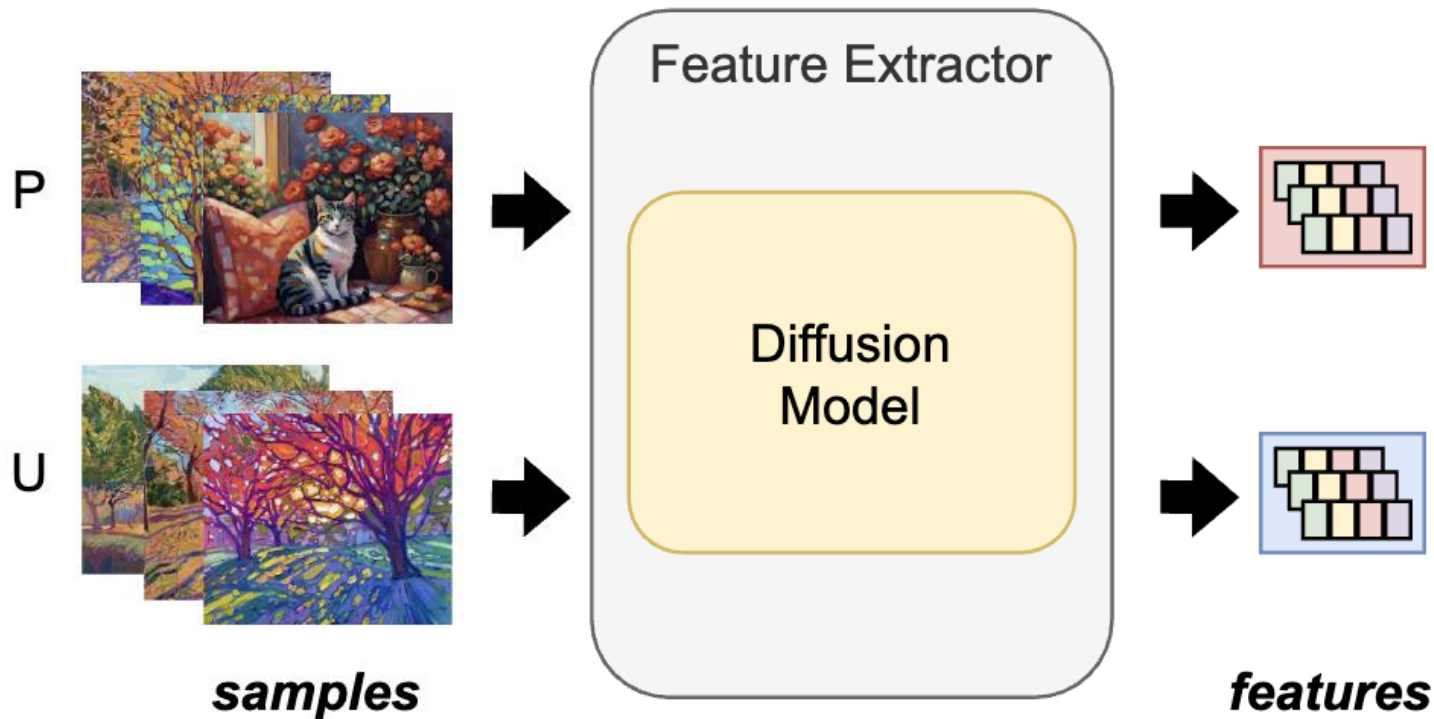# Step 1: Prepare Data and Model



$P$ - published pictures – used for training?
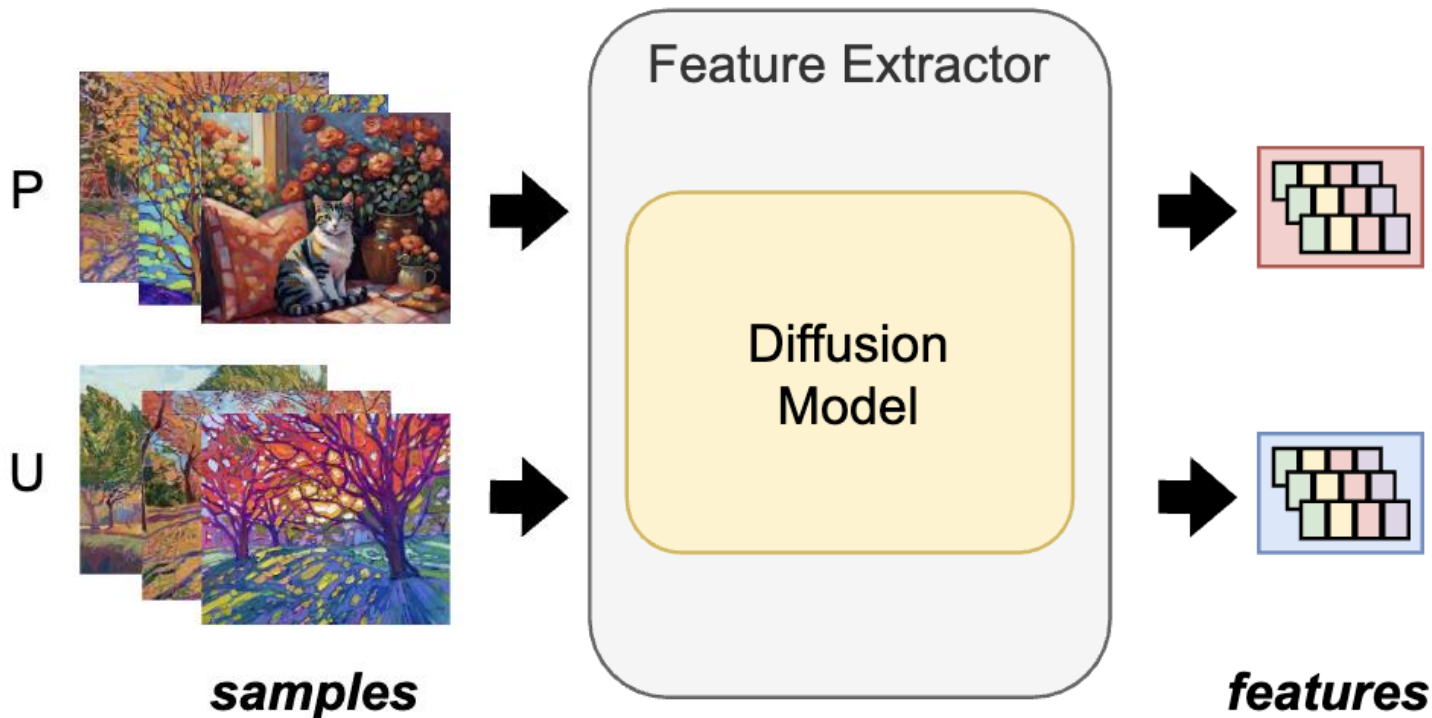
$U$ - unpublished pictures

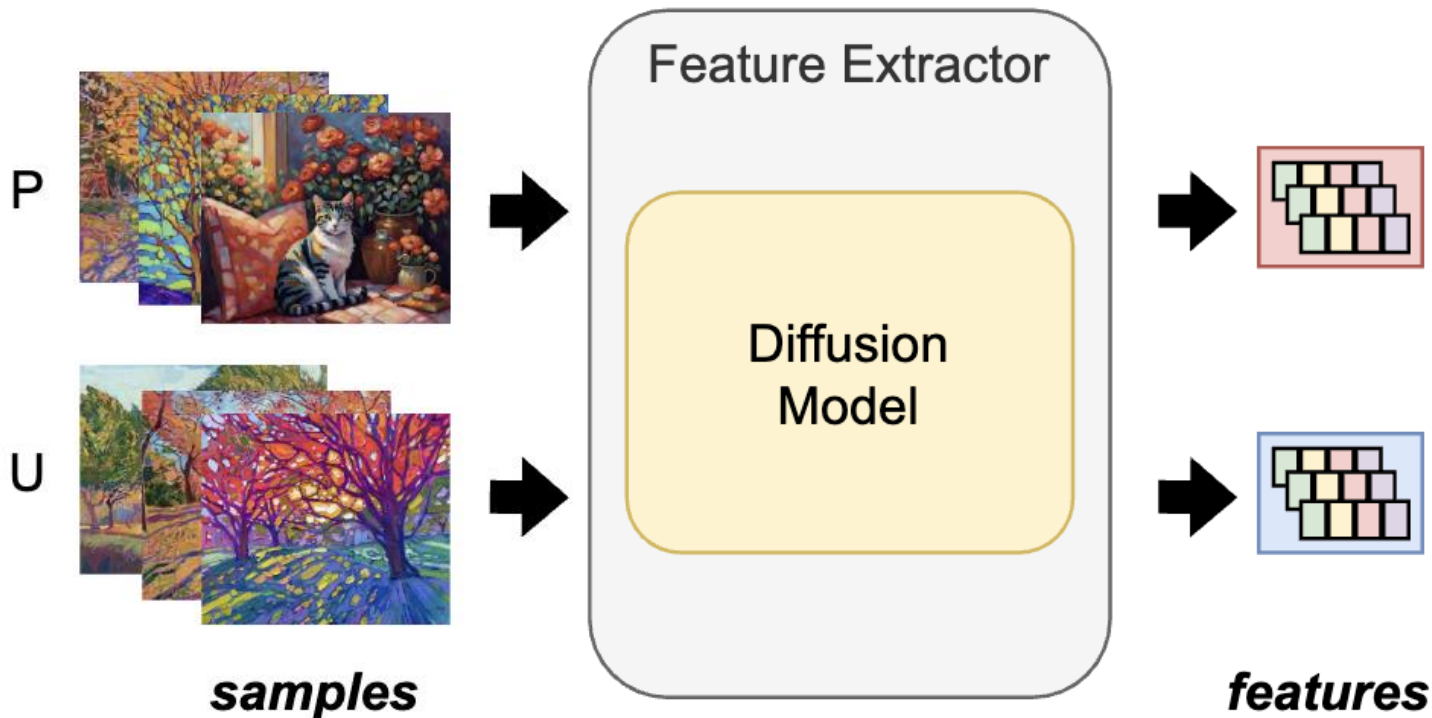**samples**

# Step 2: Feature Extraction



Existing MIA Features:

1. Denoising Loss

# Step 2: Feature Extraction



Existing MIA Features:
1. Denoising Loss
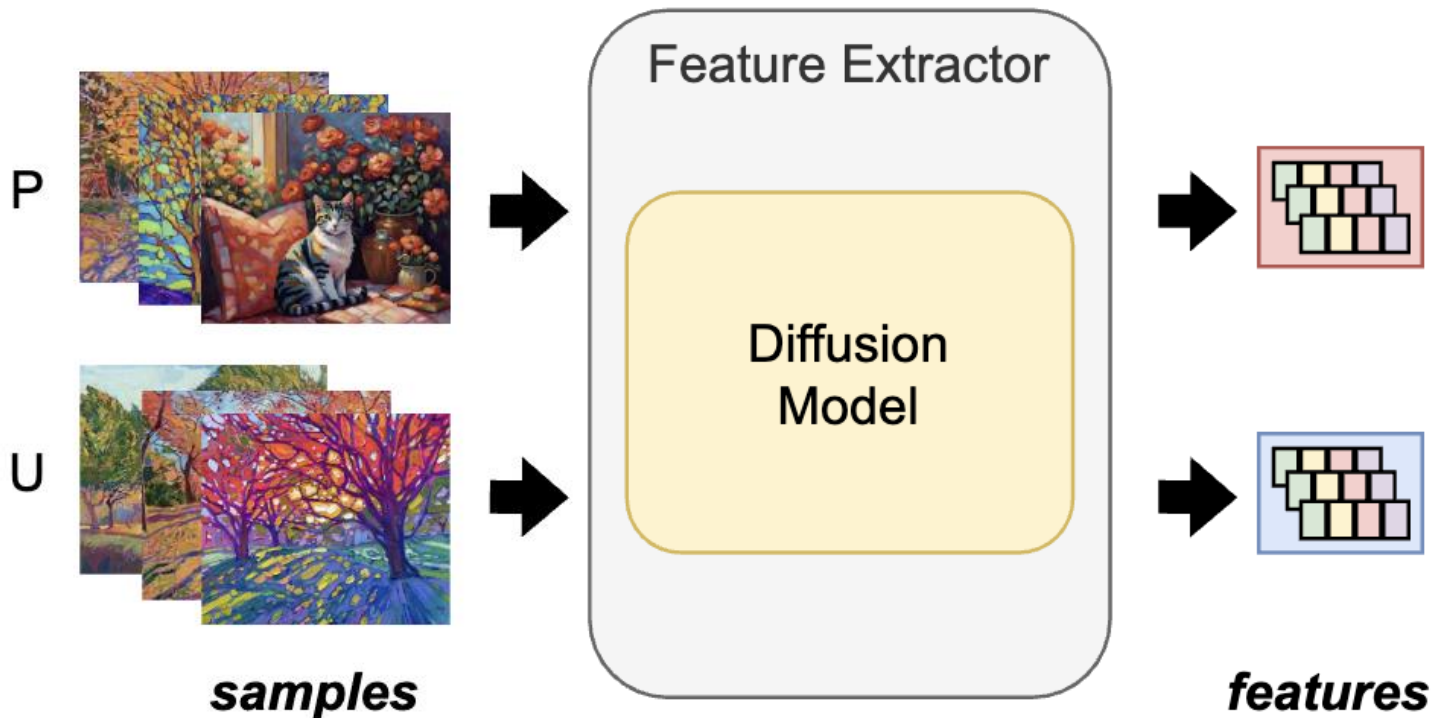2. SecMI

# Step 2: Feature Extraction



Existing MIA Features:
1. Denoising Loss
2. SecMI
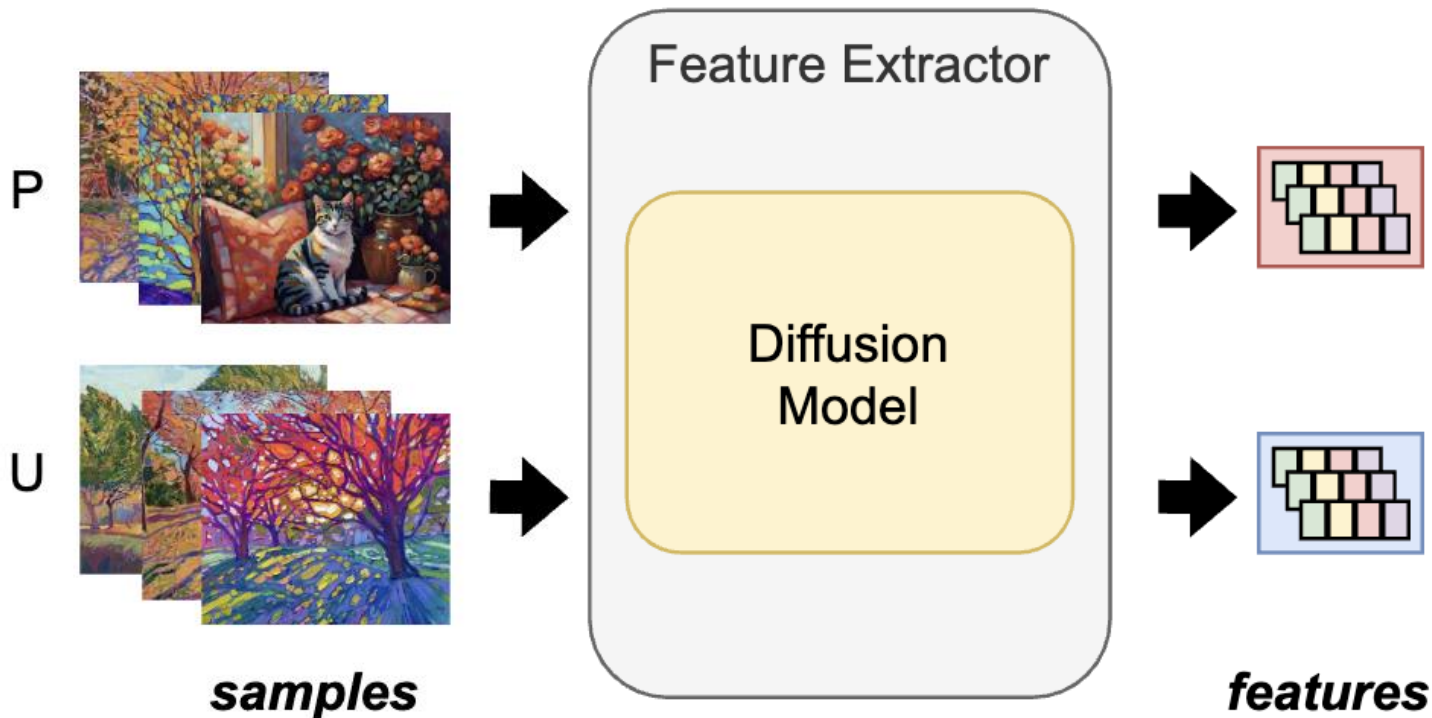3. PIA/PIAN

# Step 2: Feature Extraction



Existing MIA Features:
1. Denoising Loss
2. SecMI
3. PIA/PIAN

**Our new features**
1. Multiple Loss
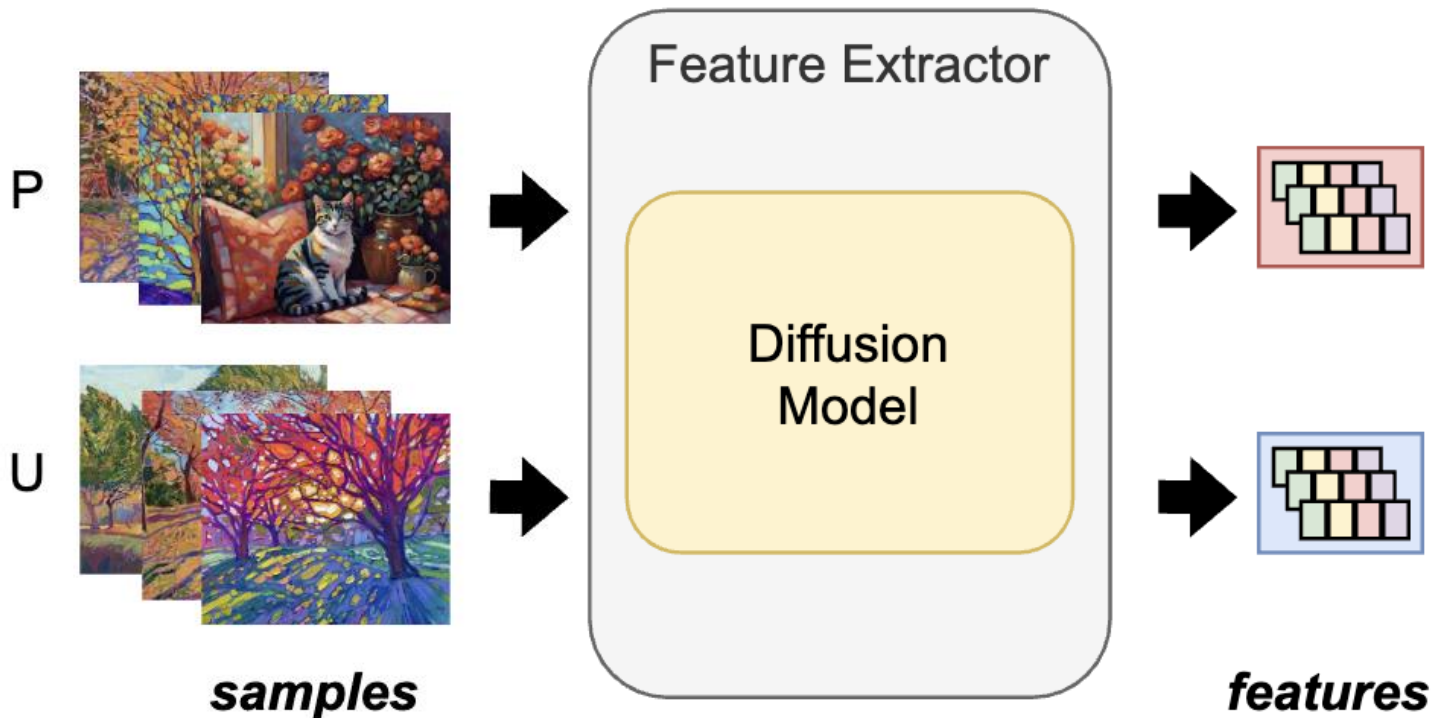
# Step 2: Feature Extraction



Existing MIA Features:
1. Denoising Loss
2. SecMI
3. PIA/PIAN

**Our new features**
1. Multiple Loss
2. Noise Optimisation

# Step 2: Feature Extraction



Existing MIA Features:
1. Denoising Loss
2. SecMI
3. PIA/PIAN

**Our new features**
1. Multiple Loss
2. Noise Optimisation
3. Gradient Masking
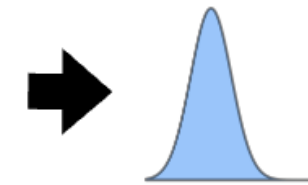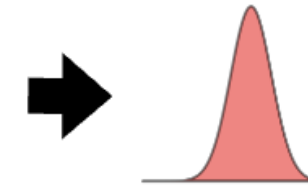
# Step 3: Scoring model

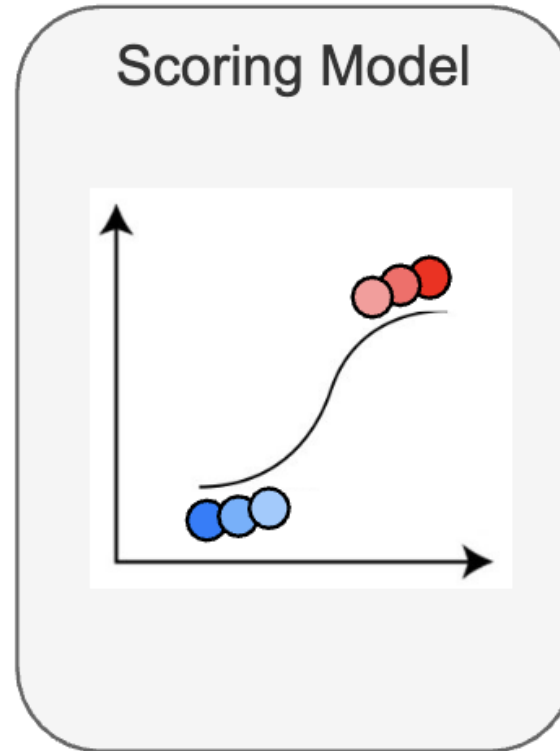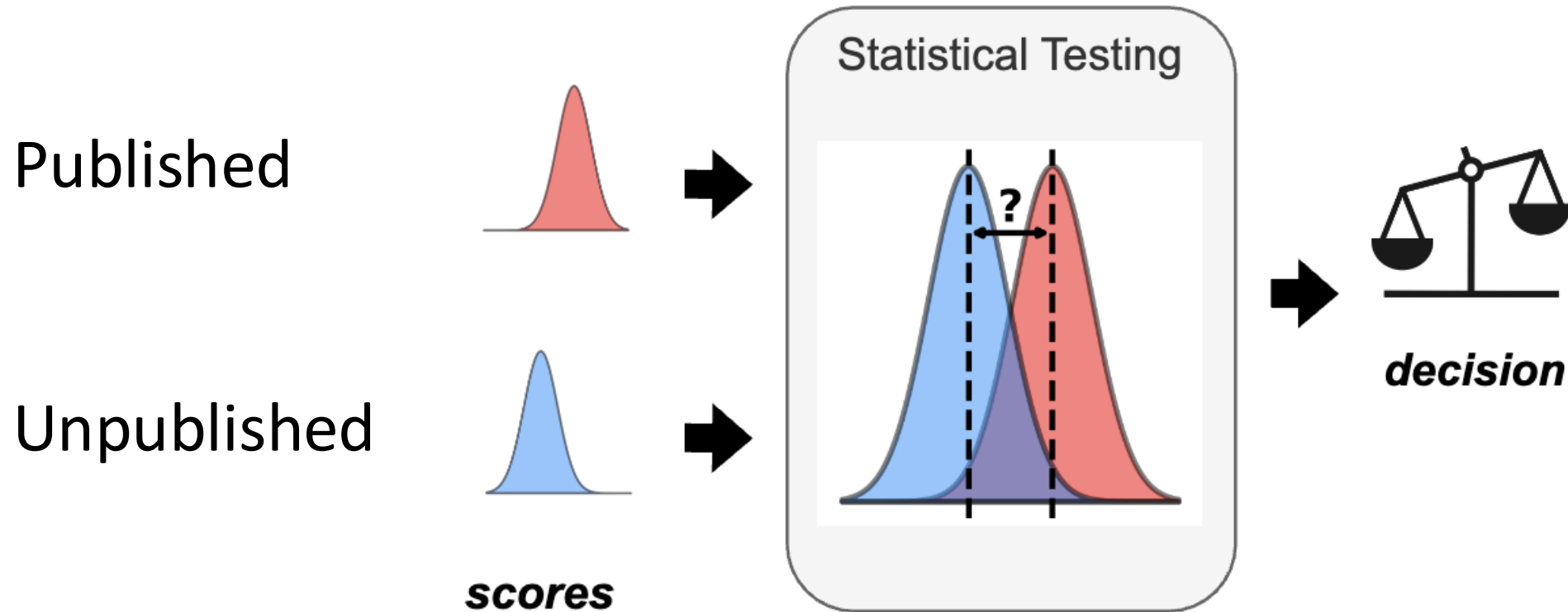$$s: Rk \rightarrow [0,1]$$



Published

Unpublished

fit and score

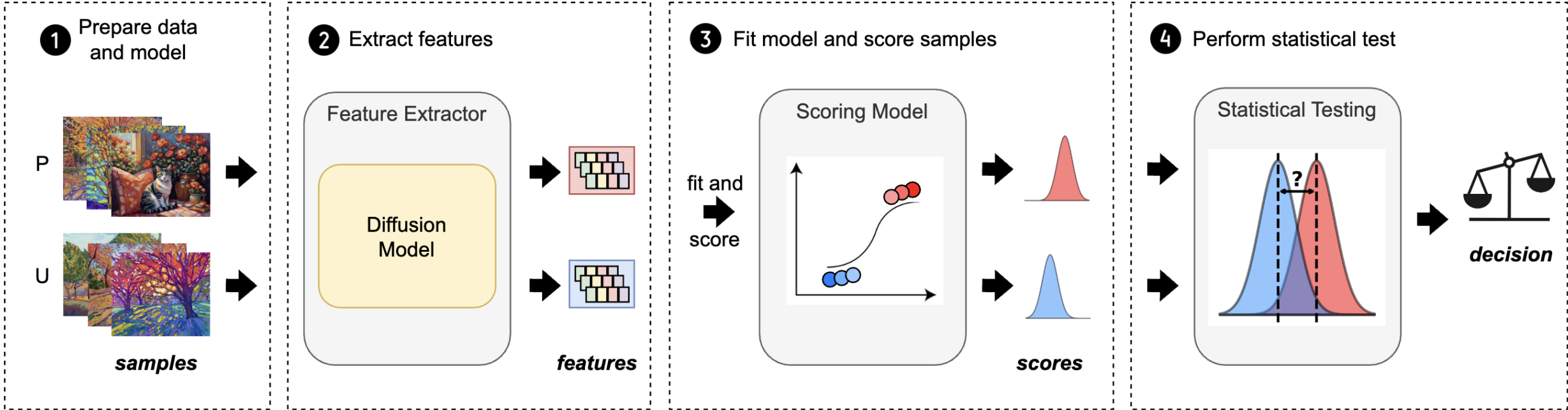Scoring Model

features

scores

# Step 4: Statistical testing

$$H_0: scores(Published) <= scores(Unpublished)$$



Published

Unpublished

scores

Statistical Testing

?

decision

# End-to-End solution



1 Prepare data and model
*samples*

2 Extract features
Feature Extractor
Diffusion Model
*features*

3 Fit model and score samples
Scoring Model
fit and score
*scores*

4 Perform statistical test
Statistical Testing
*decision*

# Experimental set-up



Diffusion Models

Trained on ImageNet

Class conditioned 256x256 images

Class conditioned 256x256 images

LDM
UViT-256
DiT-256

UViT-512
DiT-512

# Experimental set-up

# CDI works



In some case we need <100 samples

# Our new features lower the number of samples

# No False Positives

| | LDM | DiT-256 | UViT-512 | UViT-T2i |
|---|---|---|---|---|
| Members | $10^{-6}$ | $10^{-59}$ | $10^{-31}$ | ~0.0 |
| Nonmembers | 0.4 | 0.39 | 0.39 | 0.39 |

# Works if only part of data was used in training

# Key findings



1. Choose sample **s**

3. Decision: was **s** in the train set?

**State of the art membership inference methods fail on large diffusion models!**

# Key findings



1. Choose sample **s**

2. ... **M**

3. Decision: was **s** in the train se...

Forward diffusion process
$q(x_t|x_{t-1})$
$x_{t-1}$

**We shift to Dataset Inference to protect the Intellectual Property in data collections**

Forward diffusion process
$q(x_t|x_{t-1})$
$x_0 \quad x_{t-1} \quad x_t \quad x_T$
$p_\theta(x_{t-1}|x_t) \qquad N(0, I)$
Backward generative process

?

# Key findings



1. Choose sample *s*

2. Query the model *M* with sample *s*

3. Decision: was *s* in the train set?



① Prepare data and model

② Extract features

Feature Extractor

Diffusion Model

③ Fit model and score samples

Scoring Model

fit and score

④ Perform statistical test

Statistical Testing

decision

P

U

samples

features

scores

**CDI successfully identifies data collections used in training of large diffusion models**