



# GAN-Leaks: A Taxonomy of Membership Inference Attack against Generative Models



Dingfan Chen<sup>1</sup>



Ning Yu<sup>2,3</sup>



Yang Zhang<sup>1</sup>



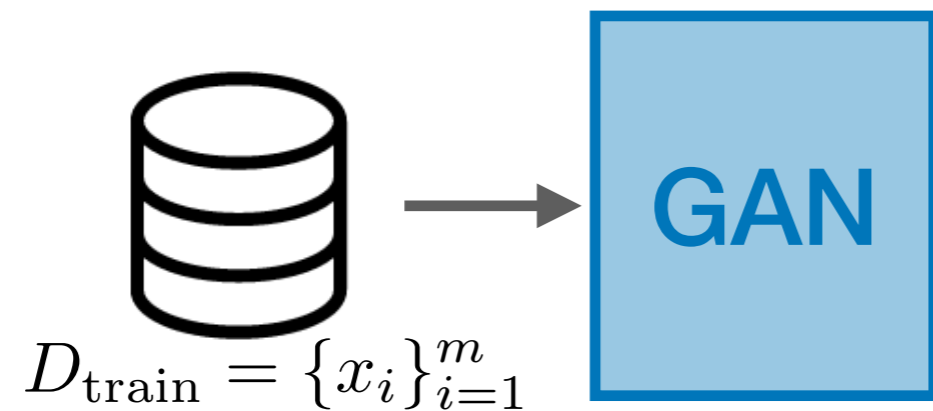
Mario Fritz<sup>1</sup>

<sup>1</sup>CISPA Helmholtz Center for Information Security, Germany

<sup>2</sup>Max Planck Institute for Informatics, Germany

<sup>3</sup>University of Maryland, College Park

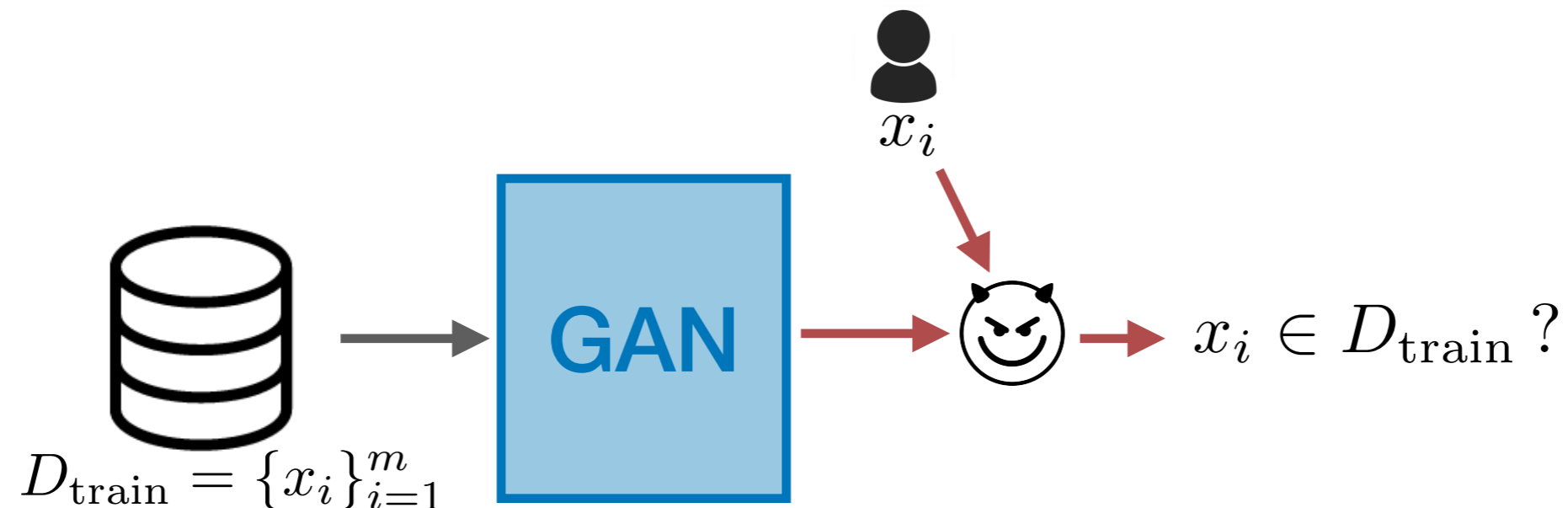
# Motivation



- Generative adversarial Networks (GANs)<sup>1</sup> have been largely used on privacy sensitive datasets, e.g., face images and medical records.

<sup>1</sup> Goodfellow et al., “Generative Adversarial Nets”, NIPS 2014

# Motivation



- Generative adversarial Networks (GANs)<sup>1</sup> have been largely used on privacy sensitive datasets, e.g., face images and medical records.
- **Our work: Membership Inference Attack** against GANs (whether a query sample  $x_i$  has been used to train a GAN model?)
- Crucial to understand and control privacy leakage; provides insights for privacy-preserving data sharing

<sup>1</sup> Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

# Contributions

- **Taxonomy**
  - categorize attack scenarios against generative models
  - benchmark future research
- **Novel attack models**
  - generic; easy-to-implement; effective; theoretically grounded
- **Extensive evaluation**
  - 3 datasets with diverse data modalities, 5 victim models, 4 attack scenarios ...

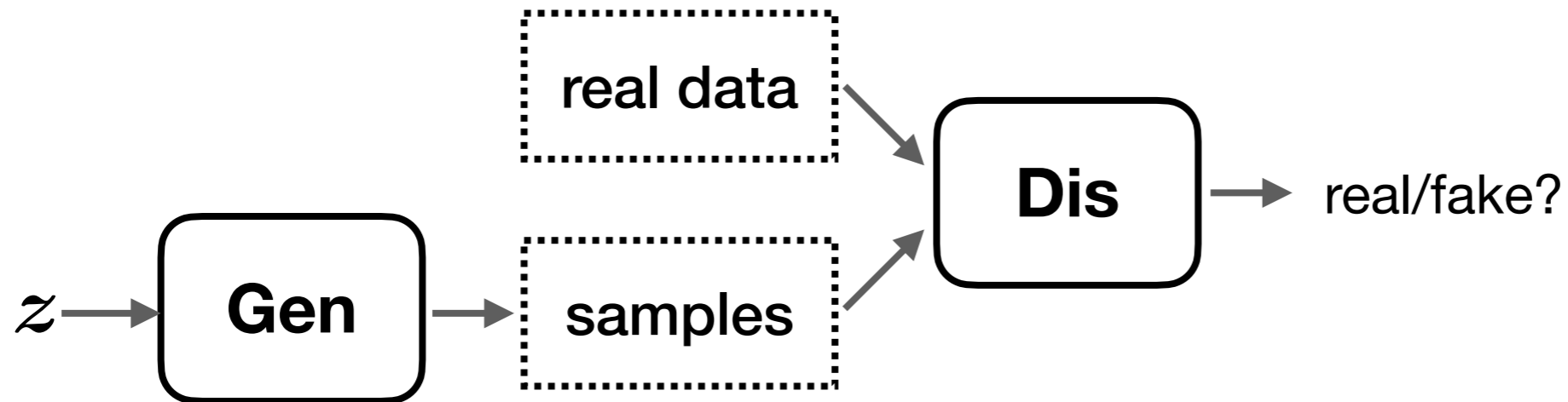
# Taxonomy

# Taxonomy

- white-box □/black-box ■?
- which GANs' components are accessible?

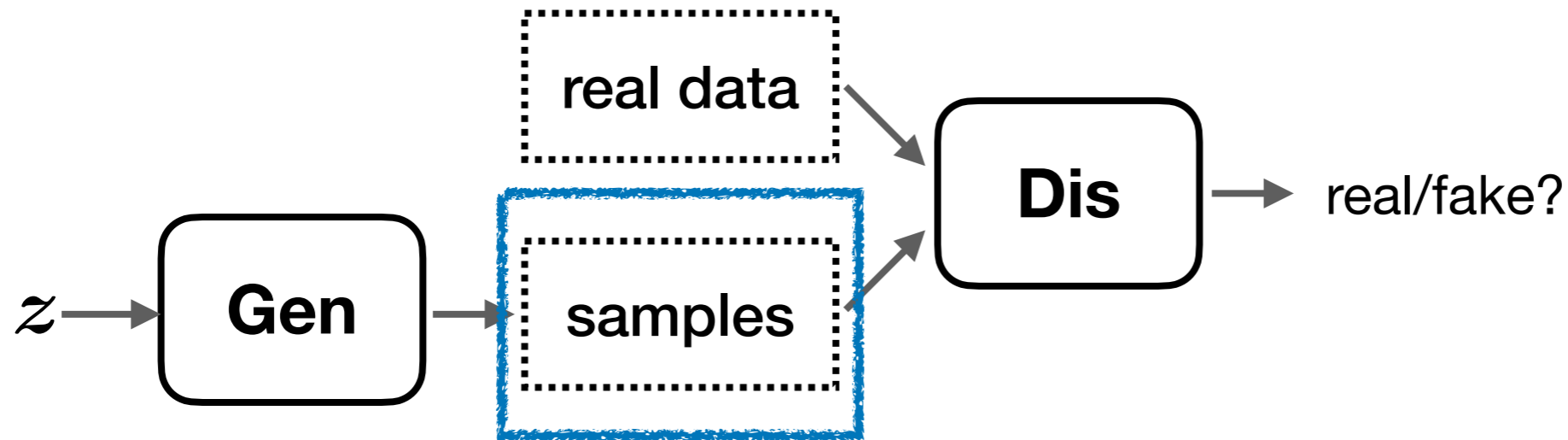
# Taxonomy

- white-box □/black-box ■?
- which GANs' components are accessible?  
( $z$ : latent code; Gen: Generator; Dis: Discriminator)



# Taxonomy

- white-box □/black-box ■?
- which GANs' components are accessible?  
( $z$ : latent code; Gen: Generator; Dis: Discriminator)



|                                   | Latent code | Generator | Discriminator |
|-----------------------------------|-------------|-----------|---------------|
| (1) Full black-box <sup>1,2</sup> | ×           | ■         | ×             |

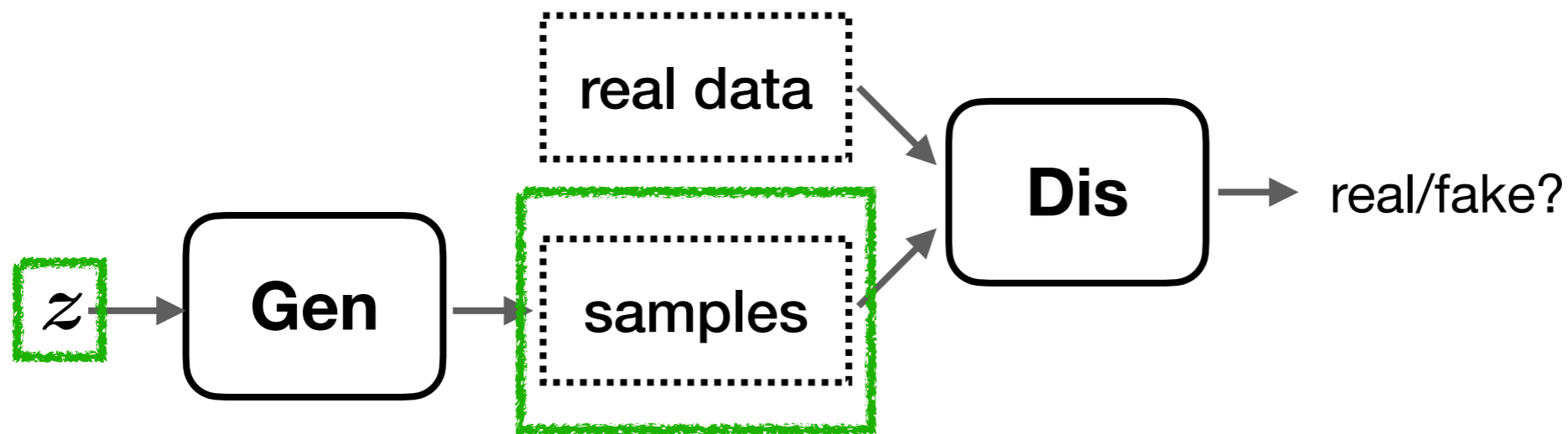
<sup>1</sup> Hayes et al., "LOGAN: Evaluating Privacy Leakage of Generative Models Using Generative Adversarial Networks", PoPETs 2019

<sup>2</sup> Hilprecht et al., "Monte Carlo and Reconstruction Membership Inference Attacks against Generative Models", PoPETs 2019



# Taxonomy

- white-box □/black-box ■?
- which GANs' components are accessible?  
( $z$ : latent code; Gen: Generator; Dis: Discriminator)



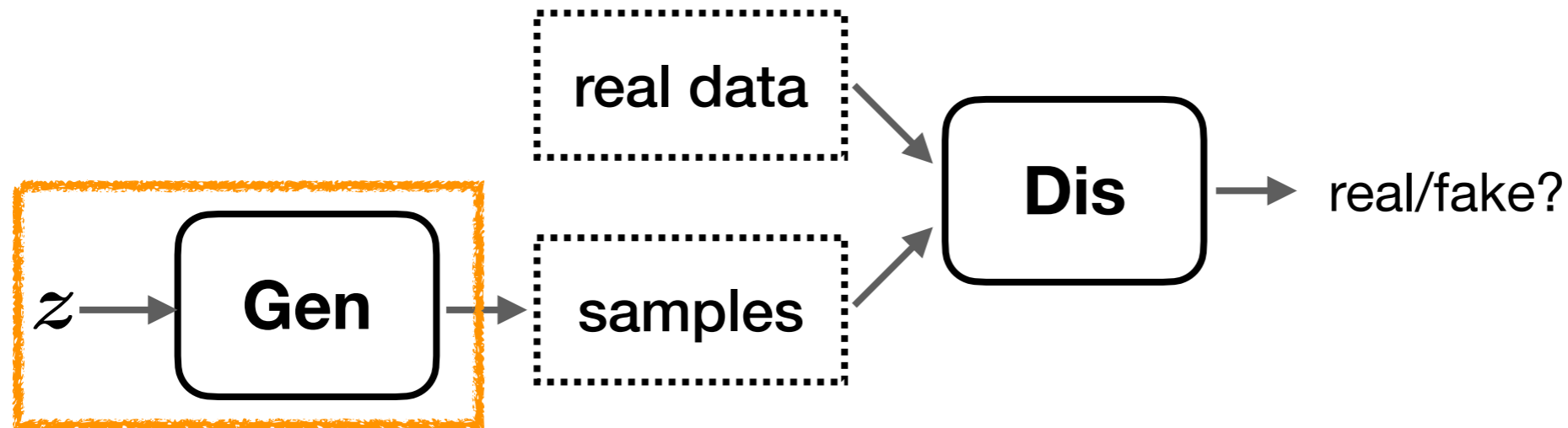
|                                   | Latent code | Generator | Discriminator |
|-----------------------------------|-------------|-----------|---------------|
| (1) Full black-box <sup>1,2</sup> | ✗           | ■         | ✗             |
| (2) Partial black-box             | ✓           | ■         | ✗             |

<sup>1</sup> Hayes et al., “LOGAN: Evaluating Privacy Leakage of Generative Models Using Generative Adversarial Networks”, PoPETs 2019

<sup>2</sup> Hilprecht et al., “Monte Carlo and Reconstruction Membership Inference Attacks against Generative Models”, PoPETs 2019

# Taxonomy

- white-box /black-box  ?
- which GANs' components are accessible?  
( $z$ : latent code; Gen: Generator; Dis: Discriminator)



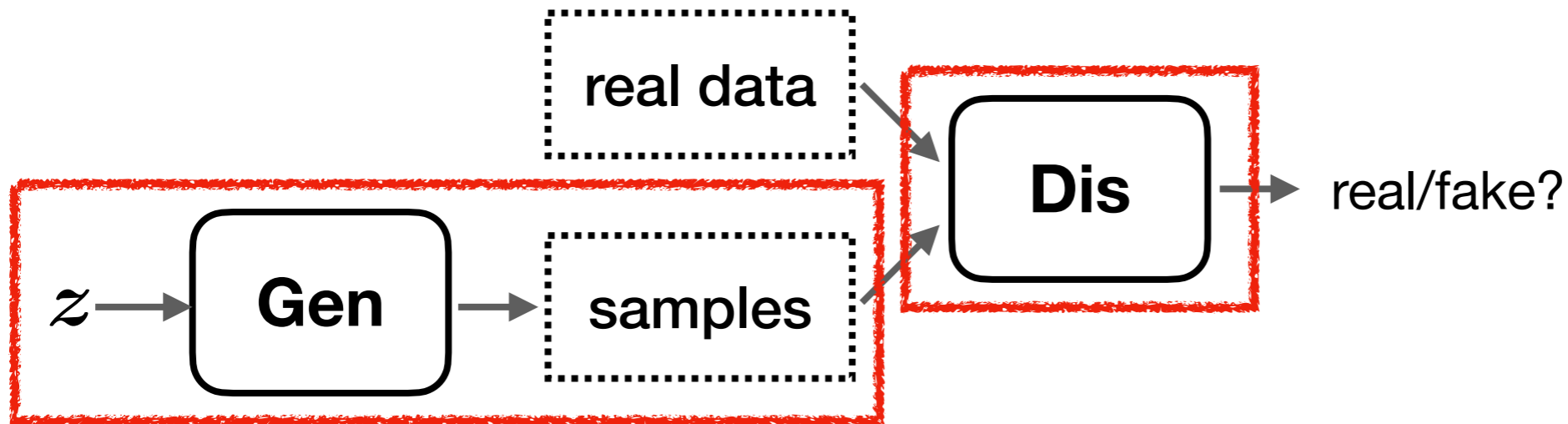
|                                   | Latent code | Generator | Discriminator |
|-----------------------------------|-------------|-----------|---------------|
| (1) Full black-box <sup>1,2</sup> | ✗           | ■         | ✗             |
| (2) Partial black-box             | ✓           | ■         | ✗             |
| (3) White-box                     | ✓           | □         | ✗             |

<sup>1</sup> Hayes et al., “LOGAN: Evaluating Privacy Leakage of Generative Models Using Generative Adversarial Networks”, PoPETs 2019

<sup>2</sup> Hilprecht et al., “Monte Carlo and Reconstruction Membership Inference Attacks against Generative Models”, PoPETs 2019

# Taxonomy

- white-box /black-box  ?
- which GANs' components are accessible?  
( $z$ : latent code; Gen: Generator; Dis: Discriminator)

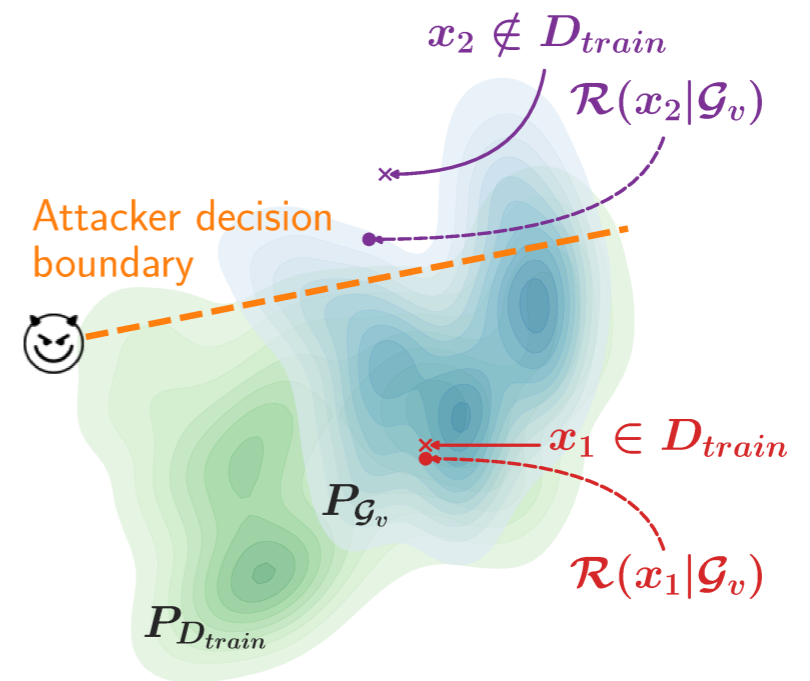


|  | Latent code | Generator | Discriminator |
|--|-------------|-----------|---------------|
| (1) Full black-box <sup>1,2</sup>      | ✗           | ■         | ✗             |
| (2) Partial black-box                  | ✓           | ■         | ✗             |
| (3) White-box                          | ✓           | □         | ✗             |
| (4) Accessible full model <sup>1</sup> | ✓           | □         | ✓             |

<sup>1</sup> Hayes et al., "LOGAN: Evaluating Privacy Leakage of Generative Models Using Generative Adversarial Networks", PoPETs 2019

<sup>2</sup> Hilprecht et al., "Monte Carlo and Reconstruction Membership Inference Attacks against Generative Models", PoPETs 2019

# Method



# Method

- **Insight:**  
Smaller reconstruction error for training set data.

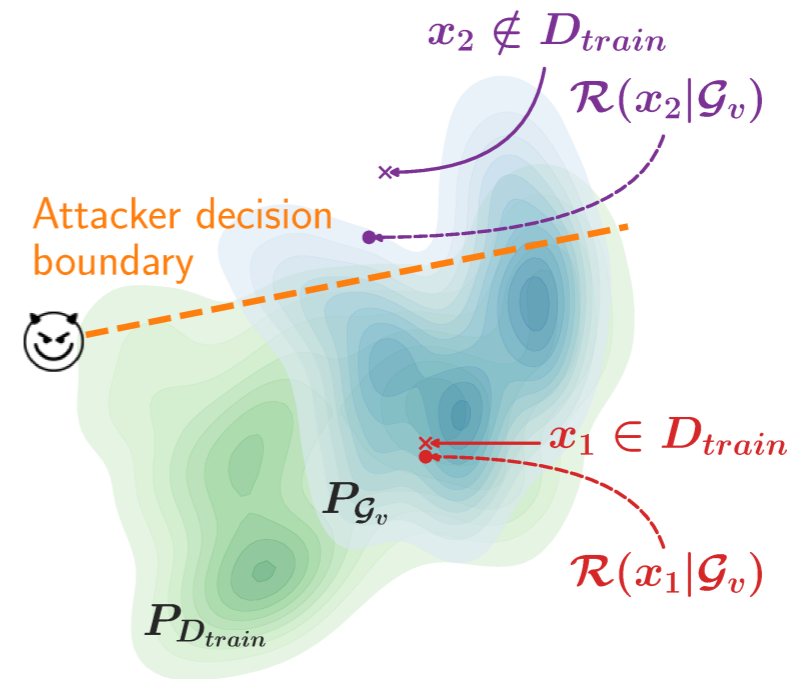
- **Generic Model:**  
optimization problem

$$\mathcal{R}(x|\mathcal{G}_v) = \mathcal{G}_v(z^*)$$

$$z^* = \operatorname{argmin}_z L(x, \mathcal{G}_v(z))$$

- **Objective:**

$$\operatorname{minimize}_z L(x, \mathcal{G}_v(z)) = \lambda_1 L_2(x, \mathcal{G}_v(z)) + \lambda_2 L_{\text{lpips}}(x, \mathcal{G}_v(z)) + \lambda_3 L_{\text{reg}}(z)$$



# Method

- **Insight:**  
Smaller reconstruction error for training set data.

- **Generic Model:**  
optimization problem

$$\mathcal{R}(x|\mathcal{G}_v) = \mathcal{G}_v(z^*)$$

$$z^* = \operatorname{argmin}_z L(x, \mathcal{G}_v(z))$$

- **Objective:**

$$\operatorname{minimize}_z L(x, \mathcal{G}_v(z)) = \lambda_1 L_2(x, \mathcal{G}_v(z)) + \lambda_2 L_{\text{lips}}(x, \mathcal{G}_v(z)) + \lambda_3 L_{\text{reg}}(z)$$

- **Different settings:**

(1) Full black-box

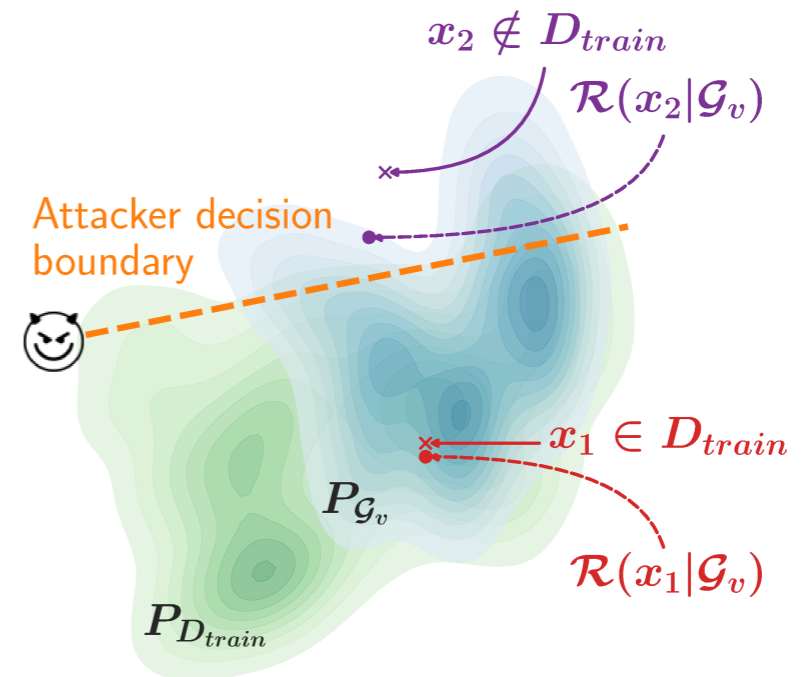
KNN search

(2) Partial black-box

Powell's conjugate direction method

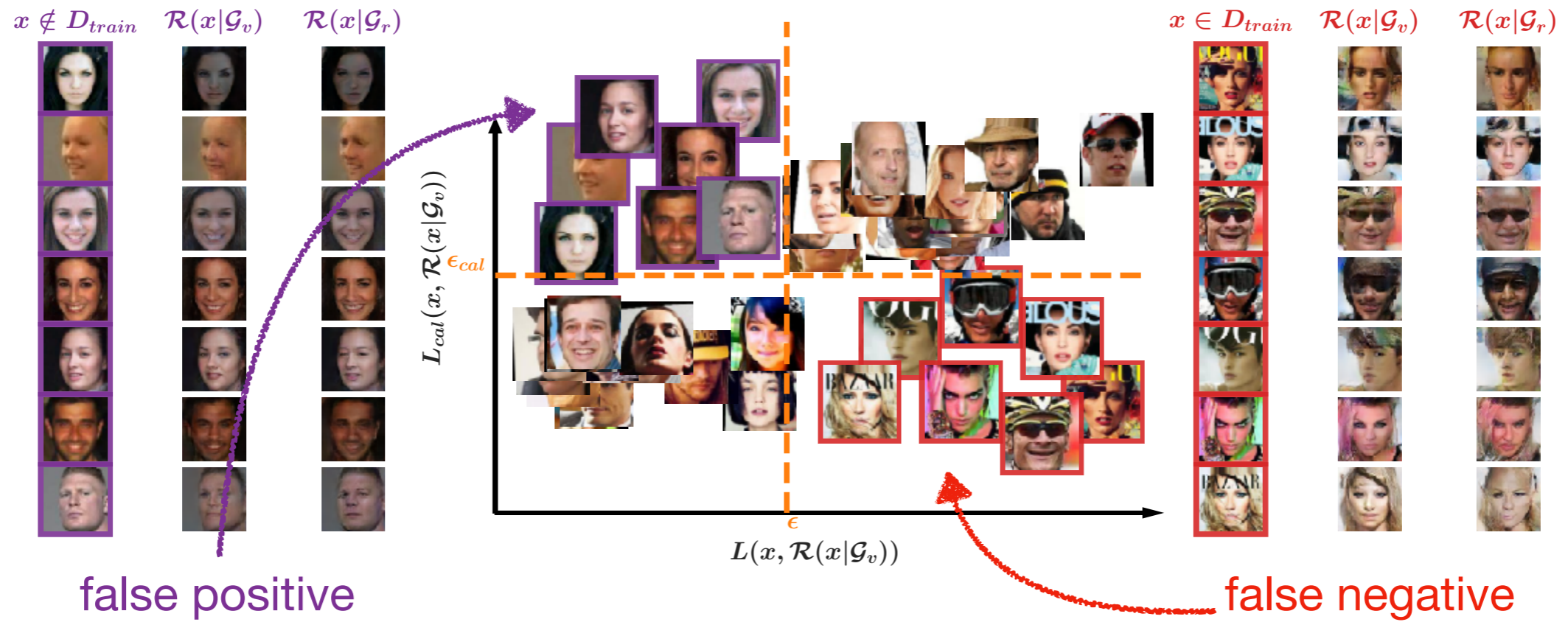
(3) White-box

L-BFGS quasi-Newton method



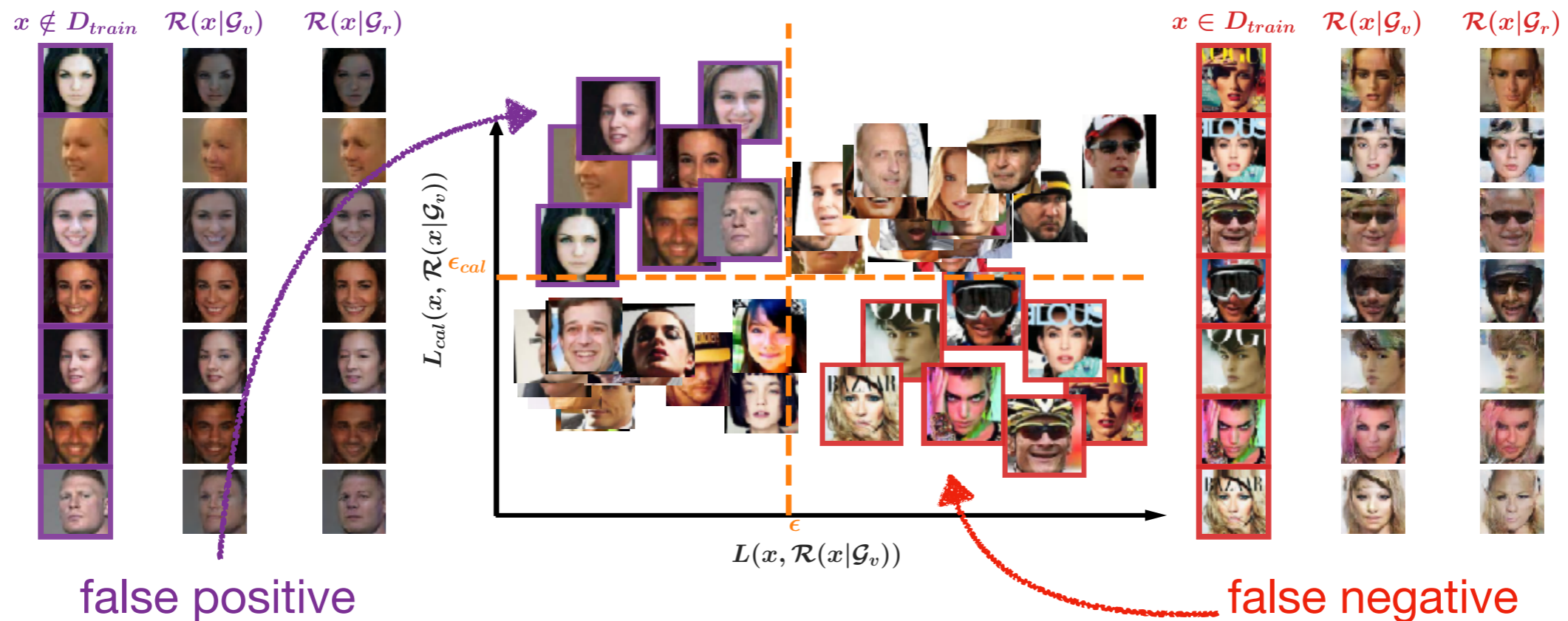
# Method

- Observe:  
Reconstruction error affected by the appearance



# Method

- Observe:  
Reconstruction error affected by the appearance



- Solution:

## Attack Calibration

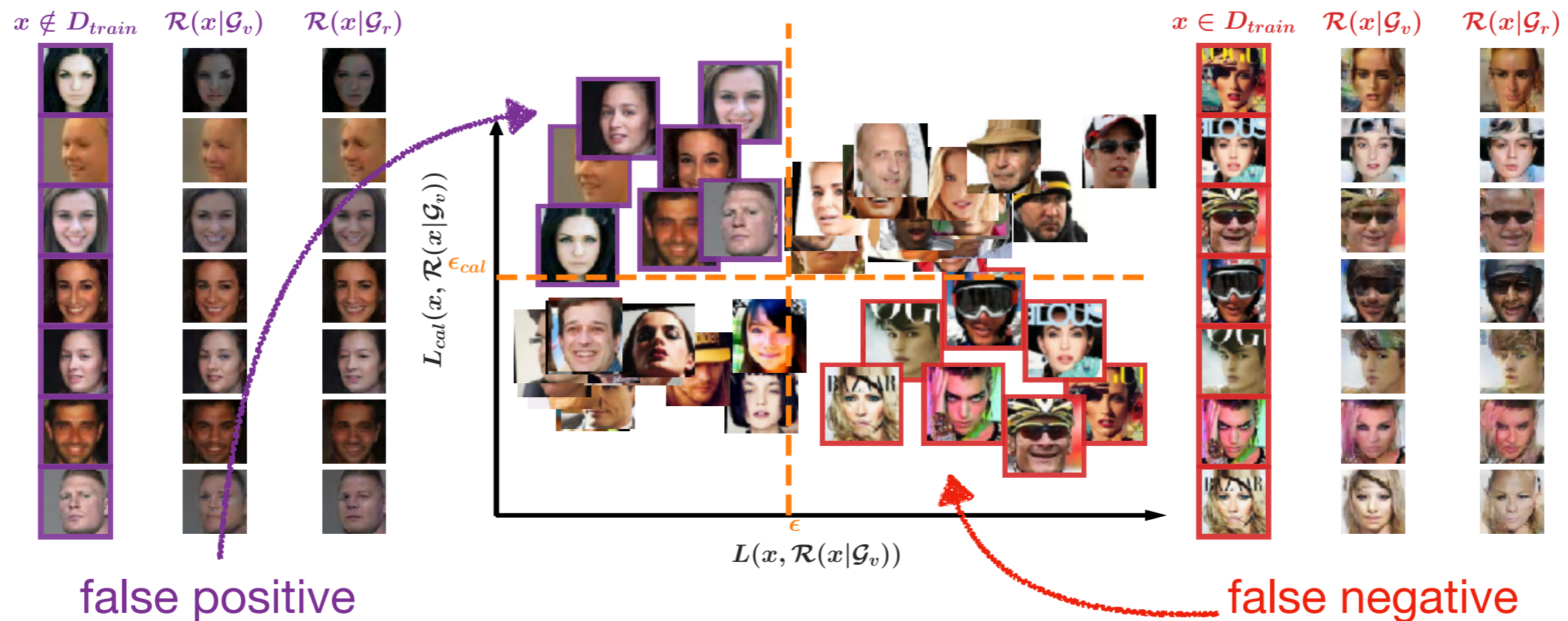
$$L_{cal}(x, \mathcal{R}(x|\mathcal{G}_v)) = L(x, \mathcal{R}(x|\mathcal{G}_v)) - L(x, \mathcal{R}(x|\mathcal{G}_r))$$

victim model          reference model



# Method

- **Observe:**  
Reconstruction error affected by the appearance



- **Solution:** **Attack Calibration**  
$$L_{cal}(x, \mathcal{R}(x|\mathcal{G}_v)) = L(x, \mathcal{R}(x|\mathcal{G}_v)) - L(x, \mathcal{R}(x|\mathcal{G}_r))$$

victim model          reference model

- **Theory:** near-optimal under a Bayesian perspective<sup>1</sup>

<sup>1</sup> Sablayrolles et al., “White-box vs Black-box: Bayes Optimal Strategies for Membership Inference”, ICML 2019

# Experiments

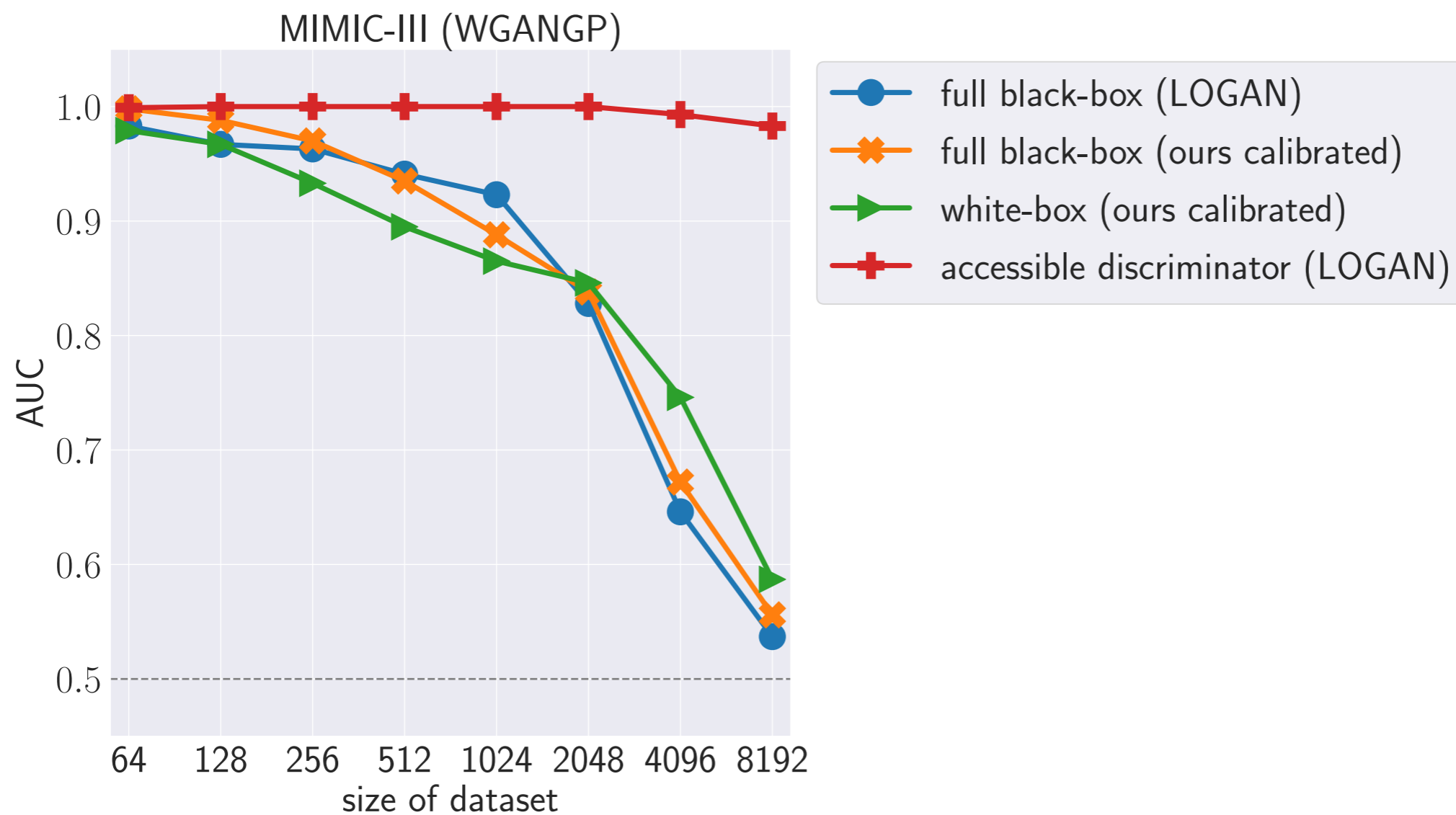
- **3 Datasets**
  - Face: CelebA
  - Location: Instagram
  - Medical: MIMIC III
- **5 GAN models**
  - PGGAN, WGANGP, DCGAN, VAEGAN, MedGAN
- **2 Baselines**
  - LOGAN<sup>1</sup>, MC<sup>2</sup>
- **Systematic Analysis**
  - dataset size, model architectures, attack settings, defense...
- **Metric**
  - AUC (Area Under the ROC Curve)
  - larger AUC → better attacker**

<sup>1</sup> Hayes et al., “LOGAN: Evaluating Privacy Leakage of Generative Models Using Generative Adversarial Networks”, PoPETs 2019

<sup>2</sup> Hilprecht et al., “Monte Carlo and Reconstruction Membership Inference Attacks against Generative Models”, PoPETs 2019

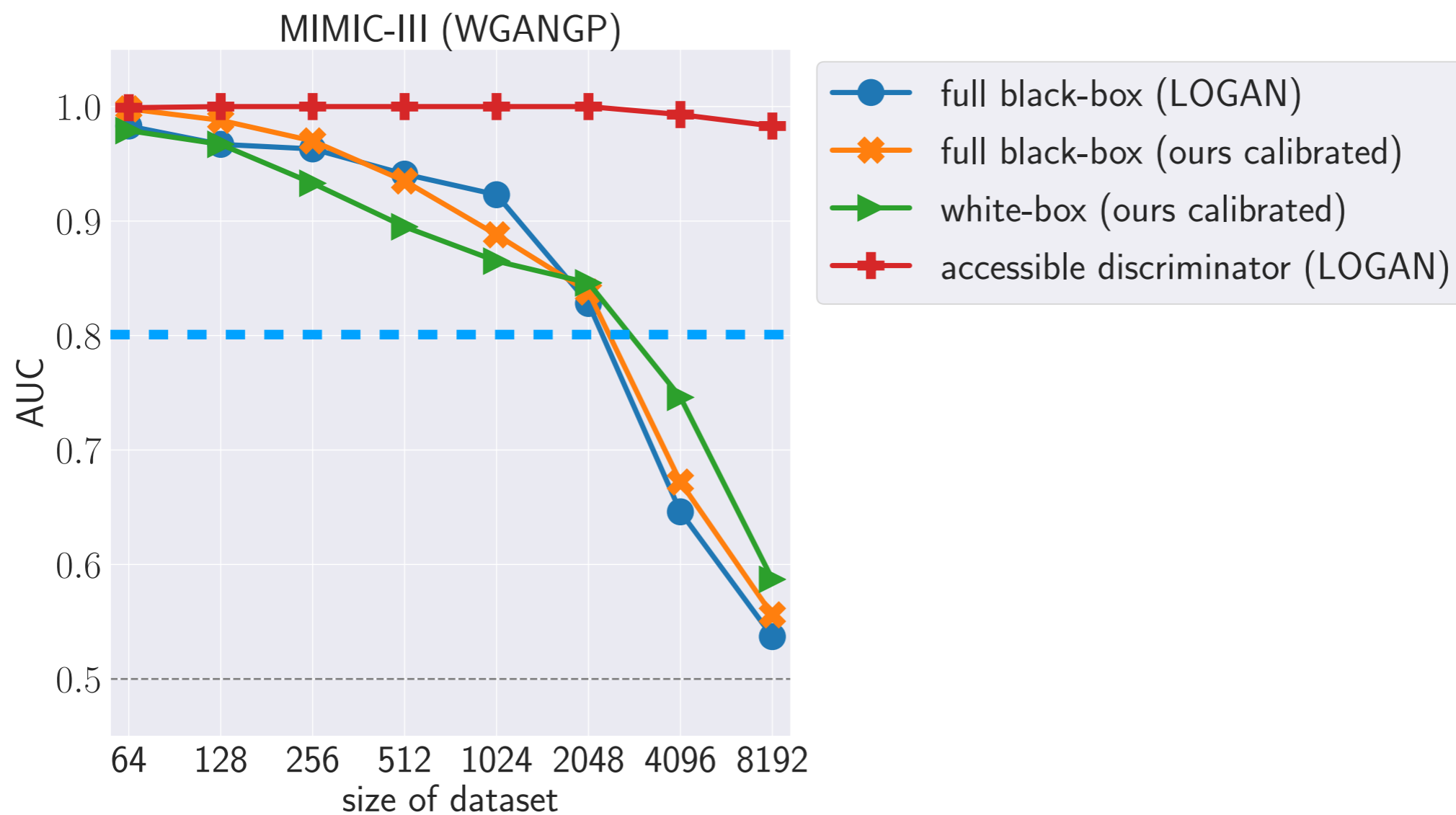
# Experiments

- Attack effectiveness — Dataset size



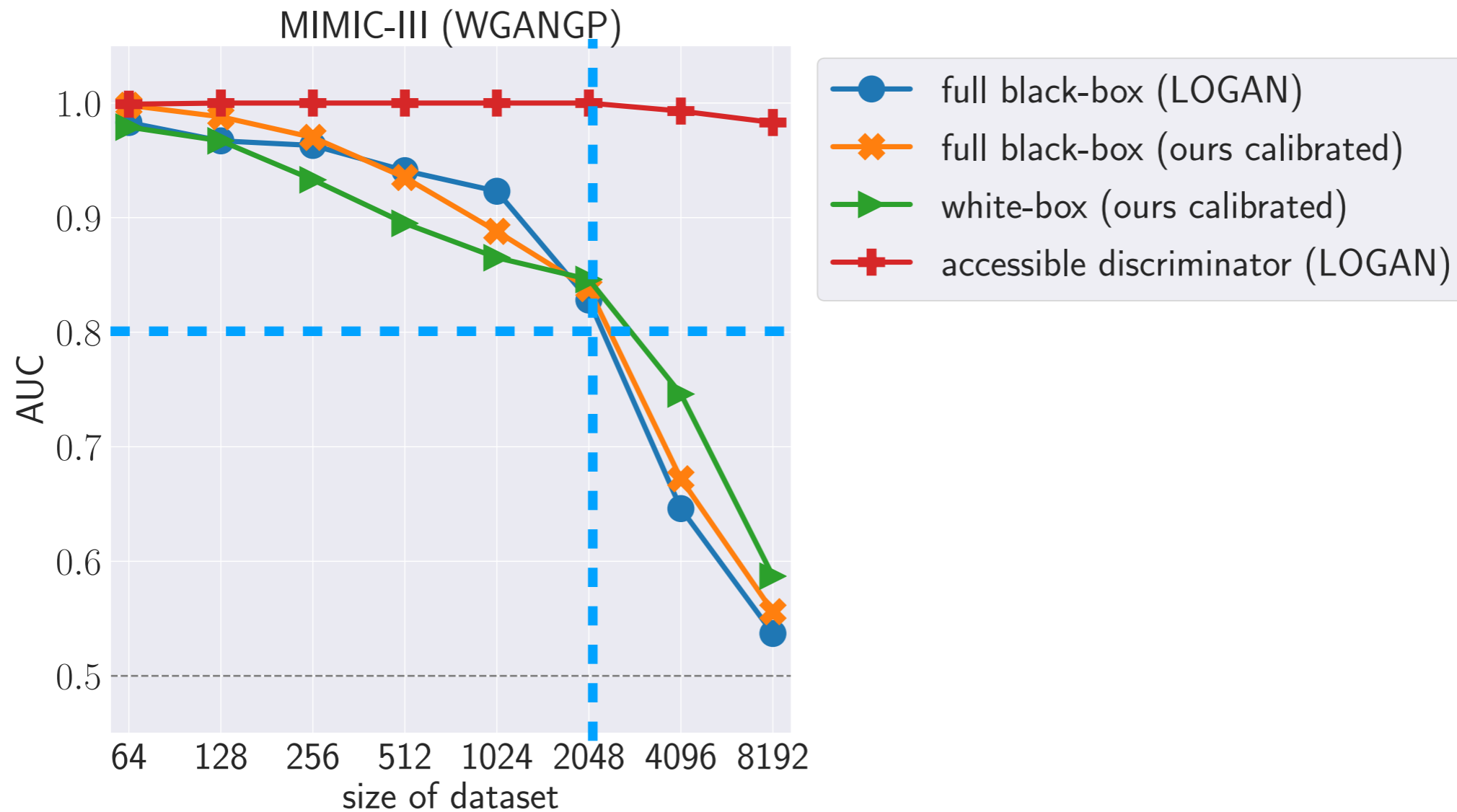
# Experiments

- Attack effectiveness — Dataset size



# Experiments

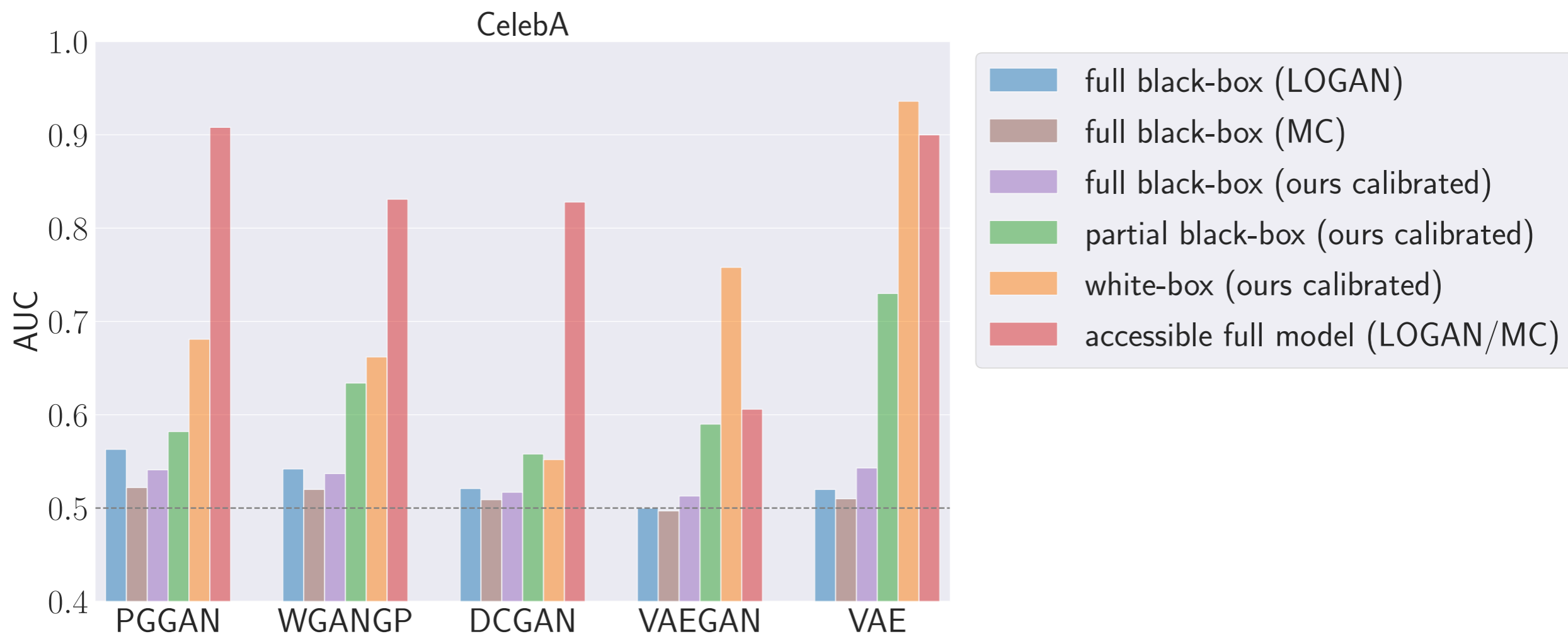
- Attack effectiveness — Dataset size



**medical dataset:** high privacy risk (AUC > 0.8) for ~2k training samples

# Experiments

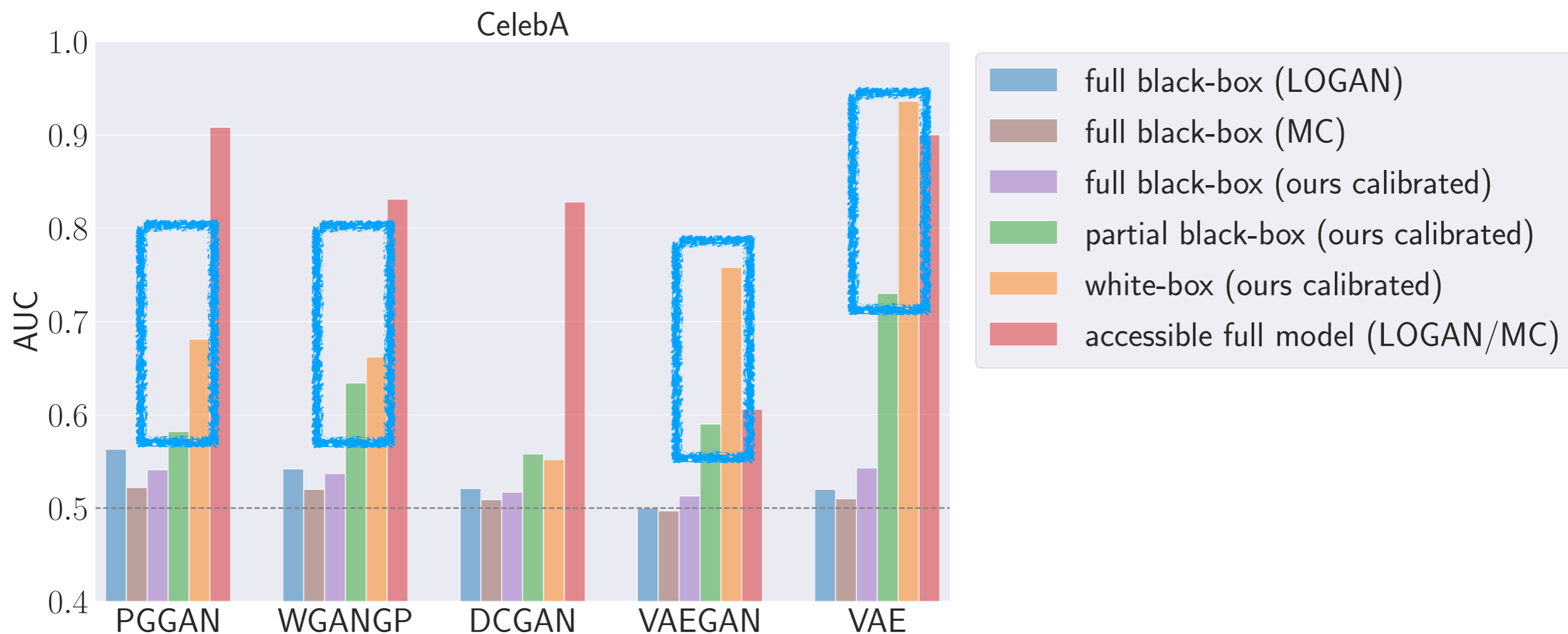
- Attack effectiveness — Attack settings, model architecture



**face dataset:** 20k training samples

# Experiments

- Attack effectiveness — Attack settings, model architecture



**face dataset:** 20k training samples

attacks are effective in practical settings

# More Details in the paper

## GAN-Leaks: A Taxonomy of Membership Inference Attack against Generative Models

Dingfan Chen<sup>1</sup>

Ning Yu<sup>2,3</sup>

Yang Zhang<sup>1</sup>

Mario Fritz<sup>1</sup>

Code and Models are available on [Github](#)



<https://github.com/DingfanChen/GAN-Leaks>

<sup>1</sup>CISPA Helmholtz Center for Information Security, Germany

<sup>2</sup>Max Planck Institute for Informatics, Germany

<sup>3</sup>University of Maryland, College Park