

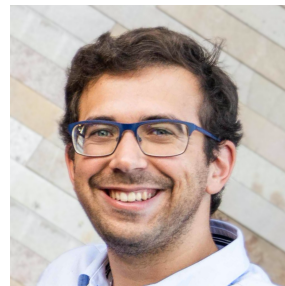
# Dataset Inference: Ownership Resolution in Machine Learning



Pratyush Maini <sup>1,2,3</sup>



Mohammad Yaghini <sup>2,3</sup>



Nicolas Papernot <sup>2,3</sup>

1



2



3



# Overview

- Why is model privacy important?
  - Primer on Model Extraction and Membership Inference
- Model Stealing – Threat Models
- Dataset Inference
  - Train-Test Prediction Margin
  - Blind Walk
  - Confidence Regressor
  - Ownership Resolution
  - Results

# Developing High-performing ML models is expensive

Computational  
Cost

Private Data

Intellectual  
Contribution

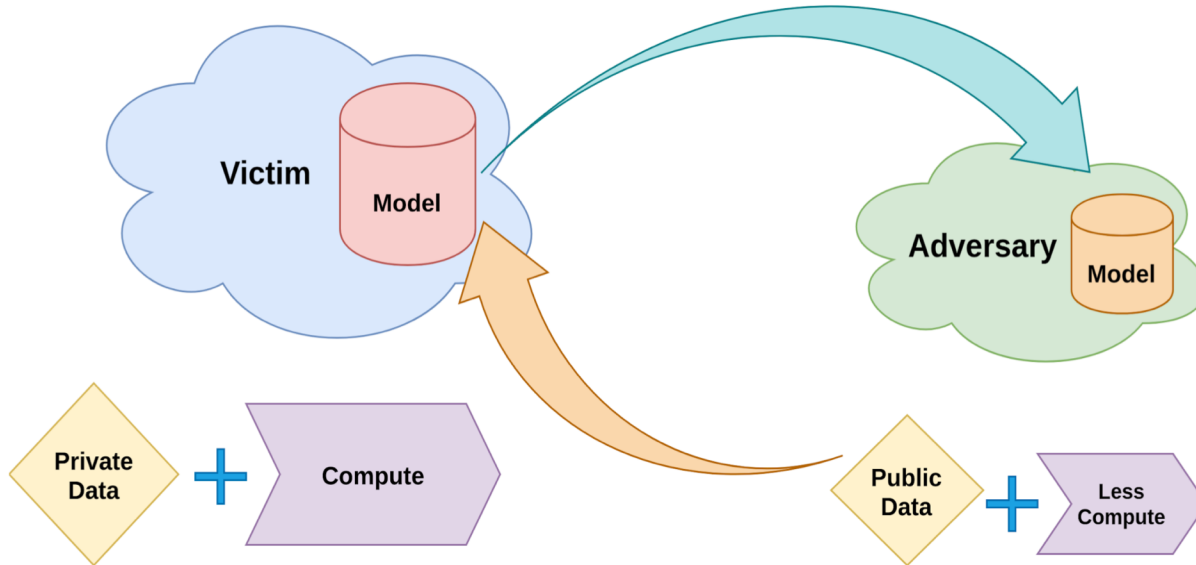
## Model Stealing Attacks are a realistic threat

Copying a model's predictions with significantly lesser cost at the adversary's end.



# Model Extraction

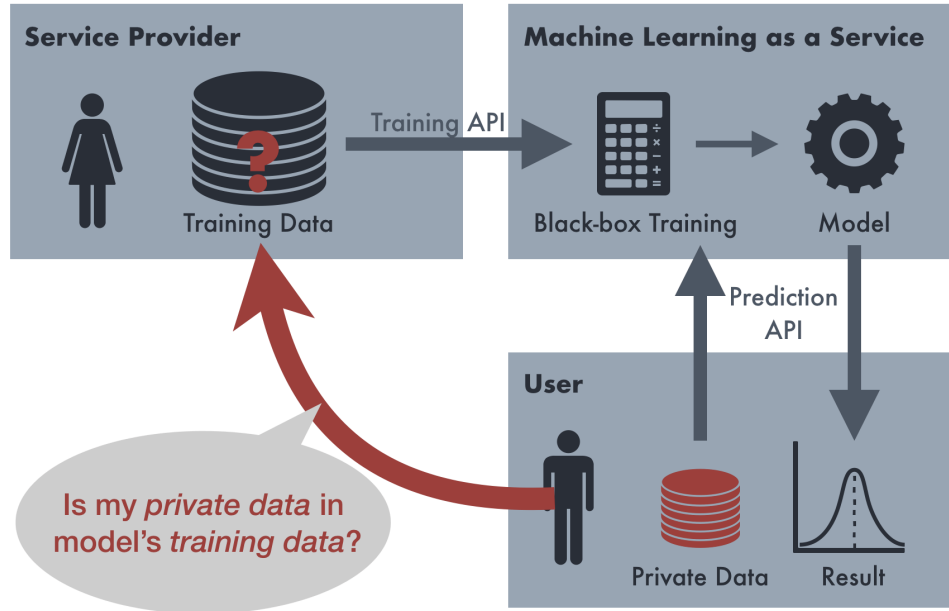
- Using predictions from an ML API (victim) to train a surrogate model using some publicly available dataset.





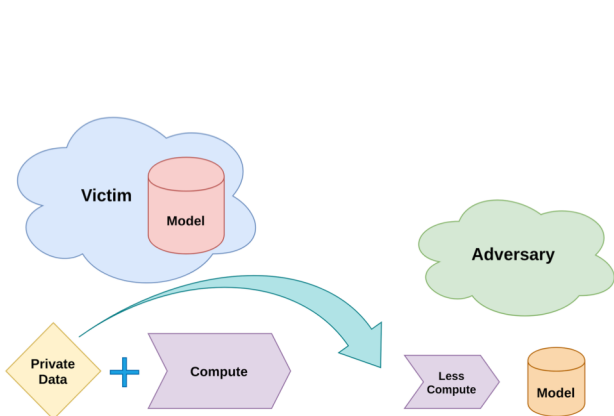
# Membership Inference

- Inferring the membership of a data-point in a model's training set.



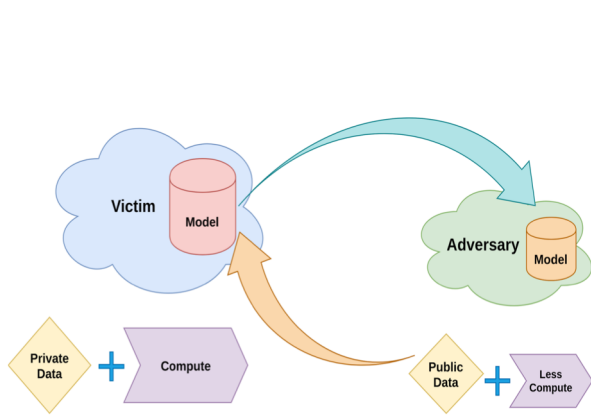
# Model Stealing Attacks: How?

An adversary may gain varying degrees of access to your 'Knowledge'



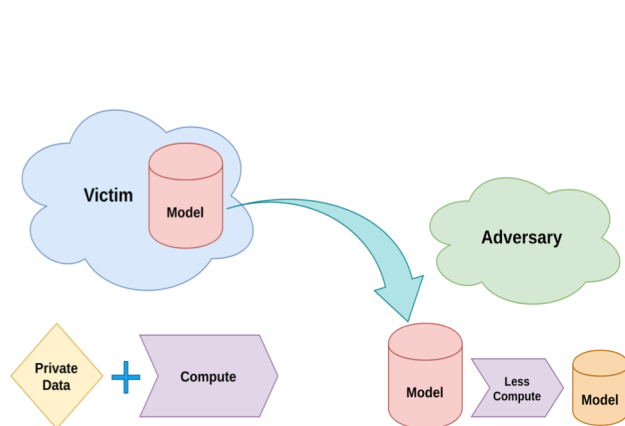
## Data Access: $A_D$

- Distillation
- Train on different architectures or hyperparameters



## Query Access: $A_Q$

- Label Access
- Logit Access

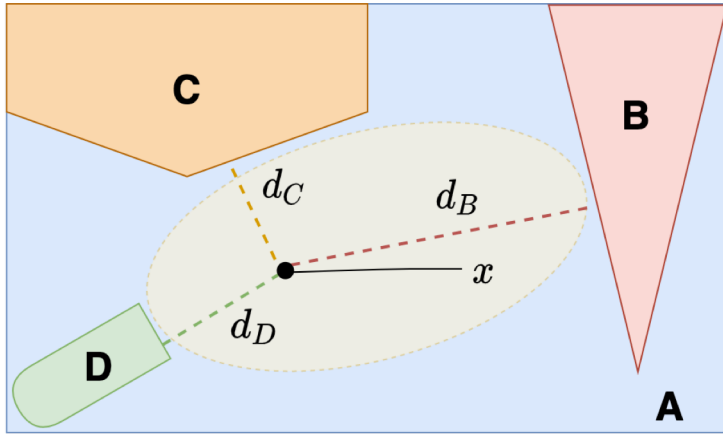


## Model Access: $A_M$

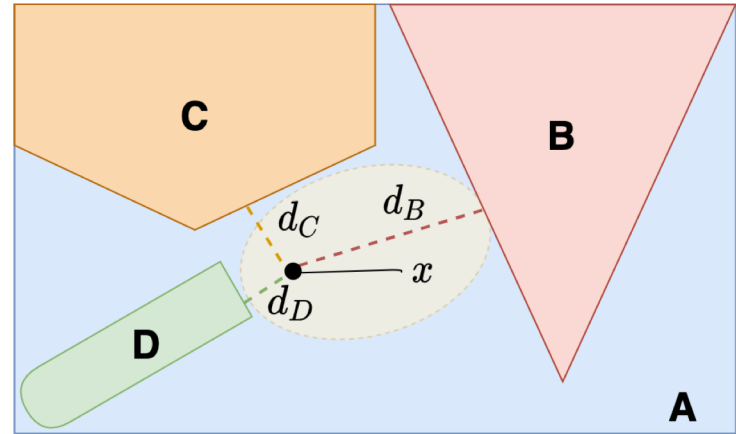
- Zero shot learning
- Fine tuning

# Dataset Inference Exploits Train-Test Prediction Certainty

Prediction Margin if  $x$  was in Train set



Prediction Margin if  $x$  was in Test set



# Analysis on a Linear Model

Binary Classification

$$y \sim \{-1, +1\};$$

Linearly Separable

$$\mathbf{x}_1 = y \cdot \mathbf{u} \in \mathbb{R}^K,$$

Gaussian Noise

$$\mathbf{x}_2 \sim \mathcal{N}(0, \sigma^2 I) \in \mathbb{R}^D$$

$$h(\mathbf{X}) = \mathbf{w}_1 \cdot \mathbf{x}_1 + \mathbf{w}_2 \cdot \mathbf{x}_2$$

Linear Classifier

**Theorem 1 (Train-Test Margin)** *Given a linear classifier  $h(\cdot)$  trained to classify inputs  $(\mathbf{X}, y) \in \mathcal{S} \subset \mathcal{D} \subset \mathbb{R}^{K+D}$ , the difference in the expected prediction margin for  $\mathbf{X}$  in  $\mathcal{S}$  and  $\mathcal{D} \setminus \mathcal{S}$ , given by  $\mathbb{E}_{\mathbf{X} \sim \mathcal{S}} [y \cdot h(\mathbf{X})] - \mathbb{E}_{\mathbf{X} \sim \mathcal{D} \setminus \mathcal{S}} [y \cdot h(\mathbf{X})] = D\sigma^2$ .*

$$\mathbf{w}_1 \leftarrow \mathbf{w}_1 + \alpha y^{(i)} \mathbf{x}_1^{(i)}$$

$$\mathbf{w}_2 \leftarrow \mathbf{w}_2 + \alpha y^{(i)} \mathbf{x}_2^{(i)}$$

# Dataset Inference Succeeds when Membership Inference Fails

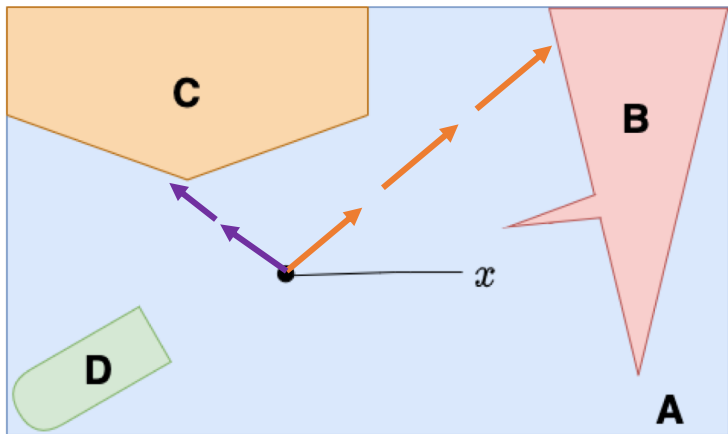
**Theorem 2 (Failure of MI)** Given a linear classifier  $h(\cdot)$  trained on  $\mathcal{S} \subset \mathcal{D} \subset \mathbb{R}^{D+K}$ , the probability that an adversary  $\mathcal{M}$  correctly predicts the membership of inputs randomly belonging to the training or test set,  $\mathbb{P}_{\mathbf{X} \sim \mathcal{R}} [\mathcal{M}(\mathbf{X}, h(\cdot)) = b] = 1 - \Phi\left(-\sqrt{\frac{D}{2m}}\right)$ , and decreases with  $|\mathcal{S}| = m$ . Moreover,  $\lim_{m \rightarrow \infty} \mathbb{P}_{\mathbf{X} \sim \mathcal{R}} [\mathcal{M}(\mathbf{X}, h(\cdot)) = b] = 0.5$ .

**Theorem 3 (Success of DI)** Choose  $b \leftarrow \{0, 1\}$  uniformly at random. Given an adversary's linear classifier  $h(\cdot)$  trained on  $\mathcal{D} \setminus \mathcal{S} \subset \mathbb{R}^{K+D}$ , if  $b = 0$ , and on  $\mathcal{S} \subset \mathcal{D}$  otherwise. The probability  $\mathcal{V}$  correctly decides if an adversary stole its knowledge  $\mathbb{P}[\psi(\mathcal{D}, h(\cdot)) = b] = 1 - \Phi\left(-\frac{\sqrt{D}}{2}\right)$ . Moreover,  $\lim_{D \rightarrow \infty} \mathbb{P}[\psi(\mathcal{D}, h(\cdot)) = b] = 1$ .

# How do you calculate the prediction certainty?

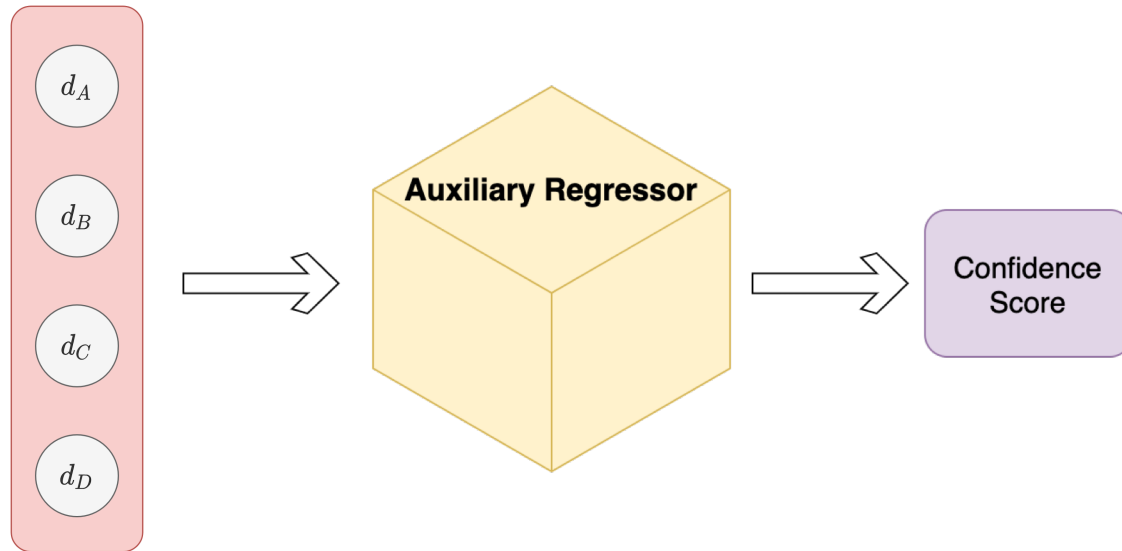
*Blind Walk* : Black-box method to estimate the prediction certainty

- Sample Random Noise*
- Take Small Steps in that direction till you reach class boundary*
- Aggregate the distance over many noise directions to create a feature embedding.*



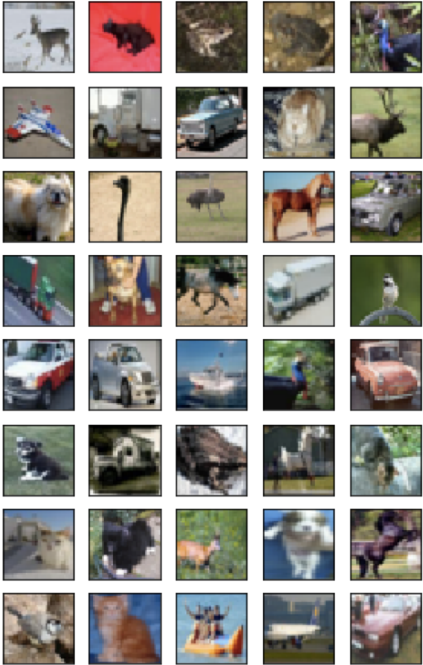
$$\text{emb}_{(\mathbf{x}, y)}^i(f) = \min_{k \in \mathbb{N}} d(x, x + k\delta_i);$$
$$s.t. f(x + k\delta_i) = t; t \neq y$$

# Training an Auxiliary Regressor

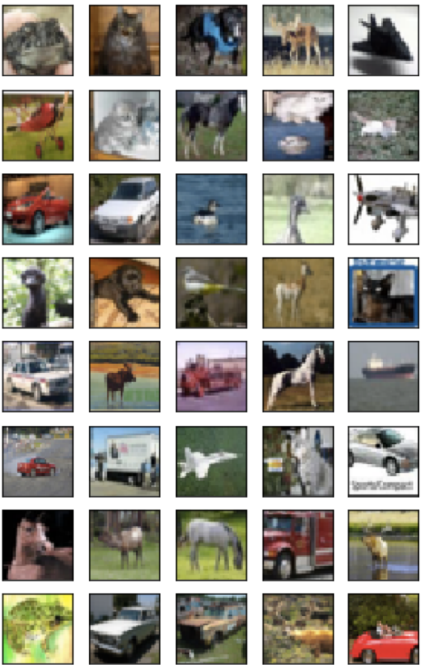


# Ownership Resolution by aggregation of Confidence Scores

Training Set



Test Set

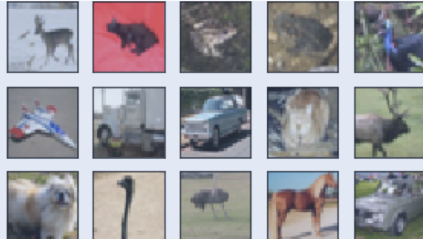


**Step 1:**  
Sample inputs from  
the train & test set

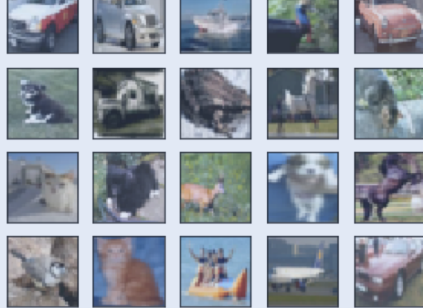


# Ownership Resolution by aggregation of Confidence Scores

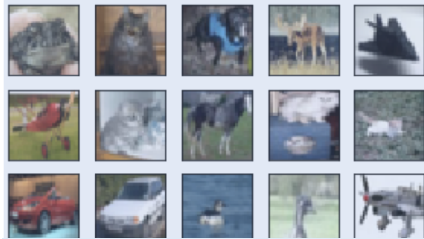
Training Set



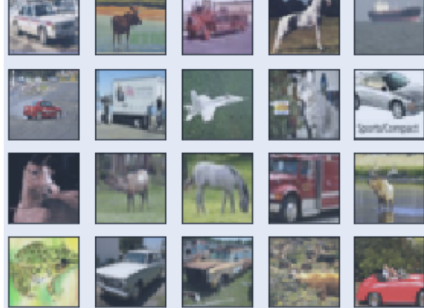
Distance embedding for each input



Test Set



Distance embedding for each input



**Step 2:**  
Generate embeddings for prediction margin

# Ownership Resolution by aggregation of Confidence Scores

Training Set



Confidence Scores for each Embedding



Test Set



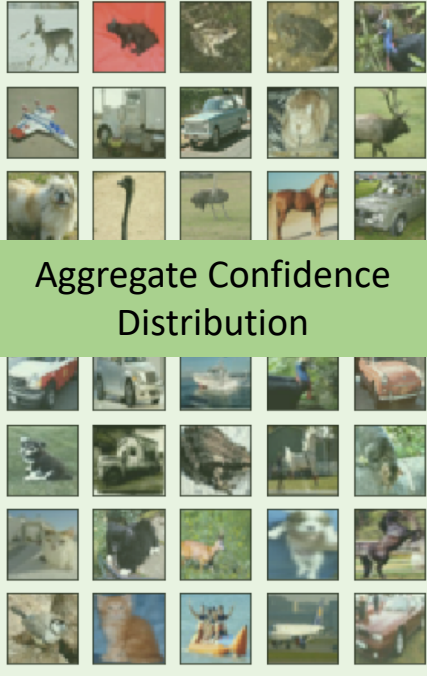
Confidence Scores for each Embedding



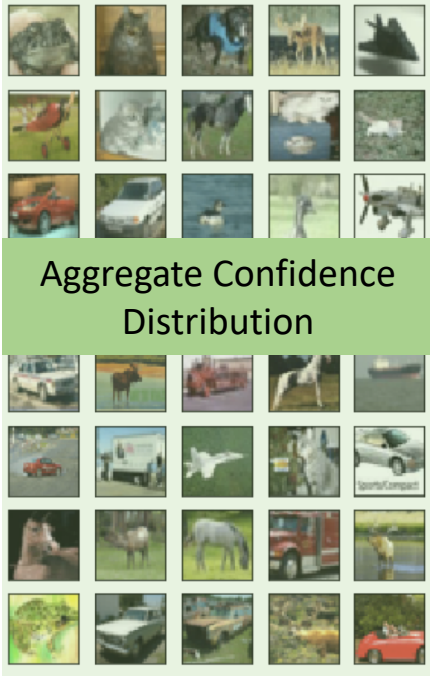
**Step 3:**  
Pass embeddings through auxiliary regressor

# Ownership Resolution by aggregation of Confidence Scores

Training Set



Test Set



**Step 4: One sided t-Test:**

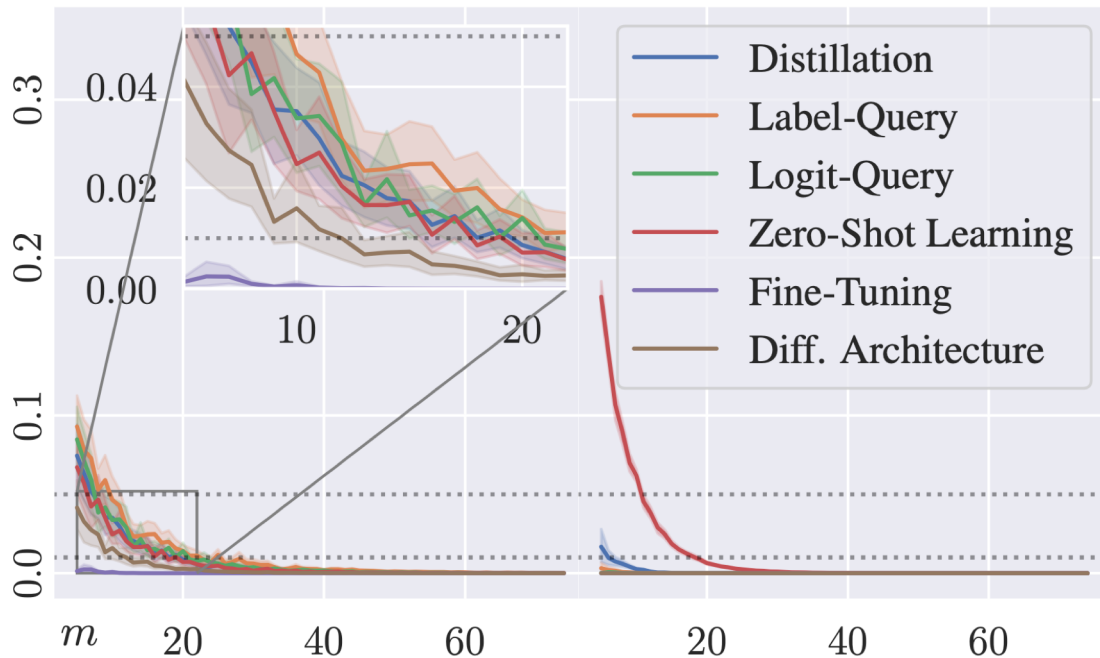
$$H_0: \mu_{test} \geq \mu_{train}$$

If stolen,  $H_0$  would be rejected.

# DI is successful across CIFAR10, SVHN, CIFAR100 and ImageNet

p-value CIFAR10 - White-Box

CIFAR10 - Black-Box



p-value against number of revealed samples ( $m$ )

Dataset Inference resolves ownership by revealing fewer than 60 private samples, with FPR < 1%

## Key Take-aways from Dataset Inference (DI)

1. Requires few private points to prove ownership.
2. Can be performed in less than 30,000 queries to the adversary.
3. White-box access is not essential to DI
4. Out-of-the-box solution that does not require overfitting or retraining.
5. Does not have a trade-off with task accuracy.