

# Auditing Data Privacy for Machine Learning

Reza Shokri

Data Privacy and Trustworthy ML Research Lab  
National University of Singapore (NUS)

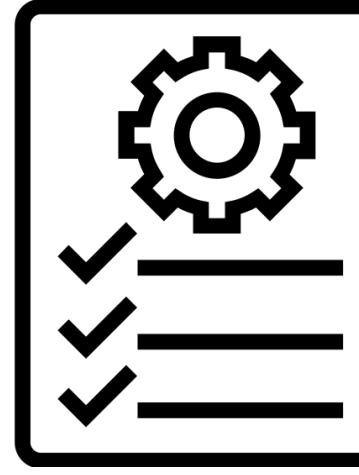
# Main Takeaways

- There is a difference between confidentiality and **privacy**
- Machine learning algorithms are extremely high **risk** algorithms, and are vulnerable to inference attacks
- We need to **audit** privacy in machine learning, but currently we are not doing it!
- There is a systematic and quantitative **method**, and an open source **tool**, for auditing privacy risks in machine learning

# Privacy Regulations



Systematic description of data collection, storage and processing



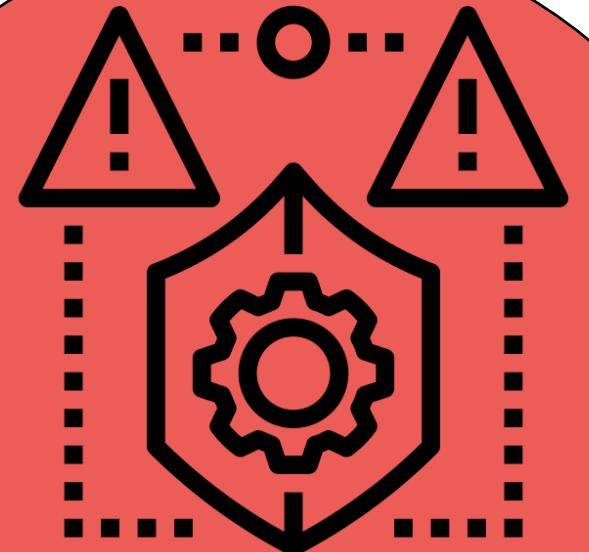
Assess necessity and proportionality



Likelihood and impact of the threats on individuals



Assess potential threats to the data



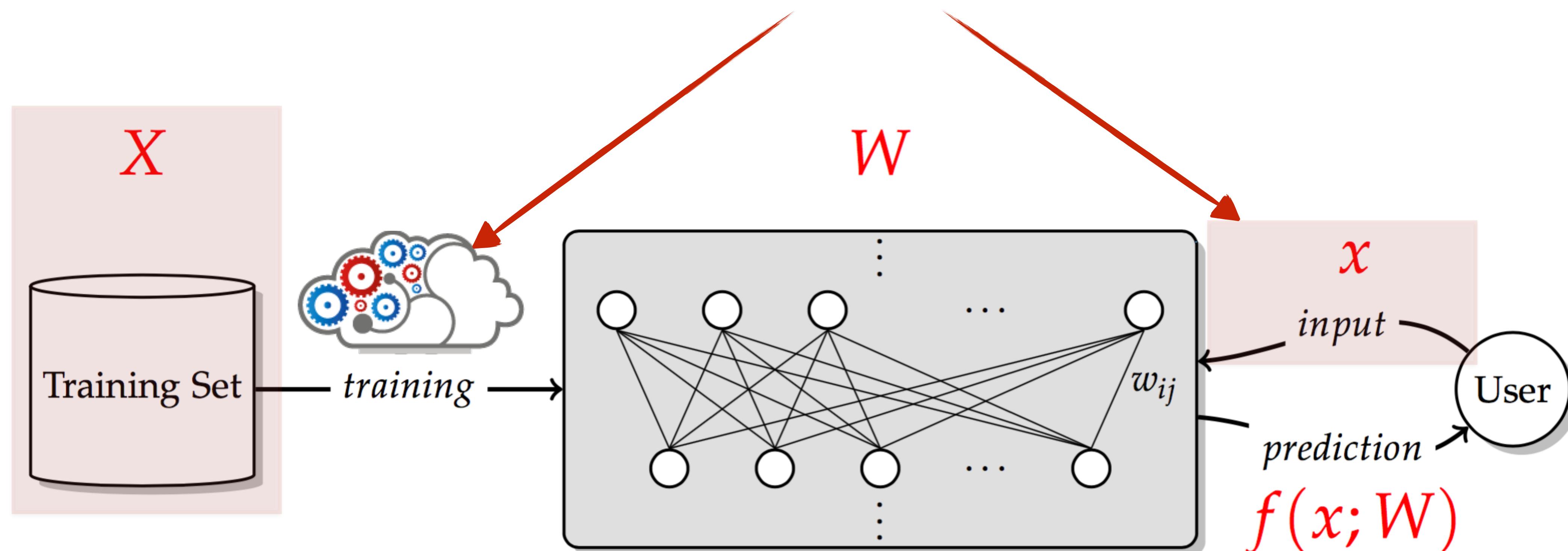
Identify and analyze possible risk mitigation measures

## GDPR — Data Protection Impact Assessment

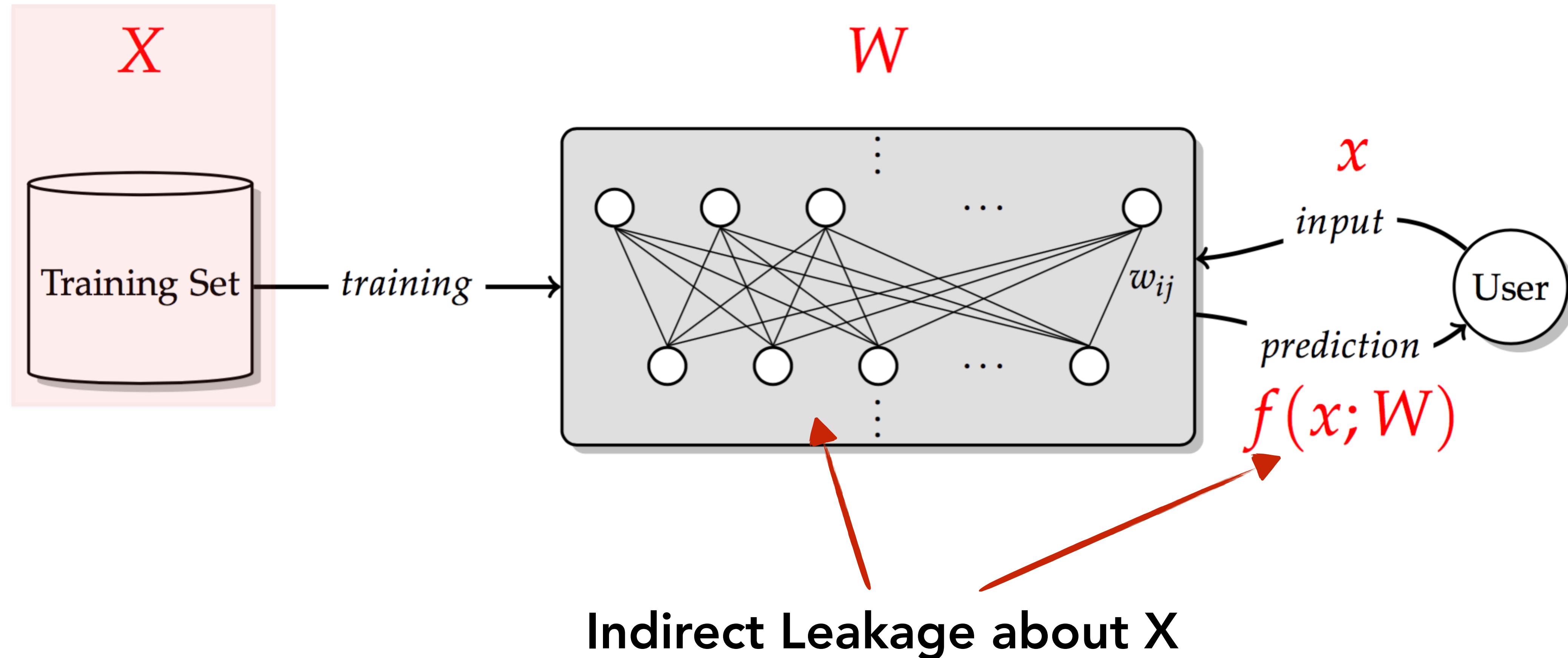
The focus is mostly on data collection, data sharing, access control, ...

# Direct Privacy Risks in Machine Learning

## Direct Access to Sensitive Data



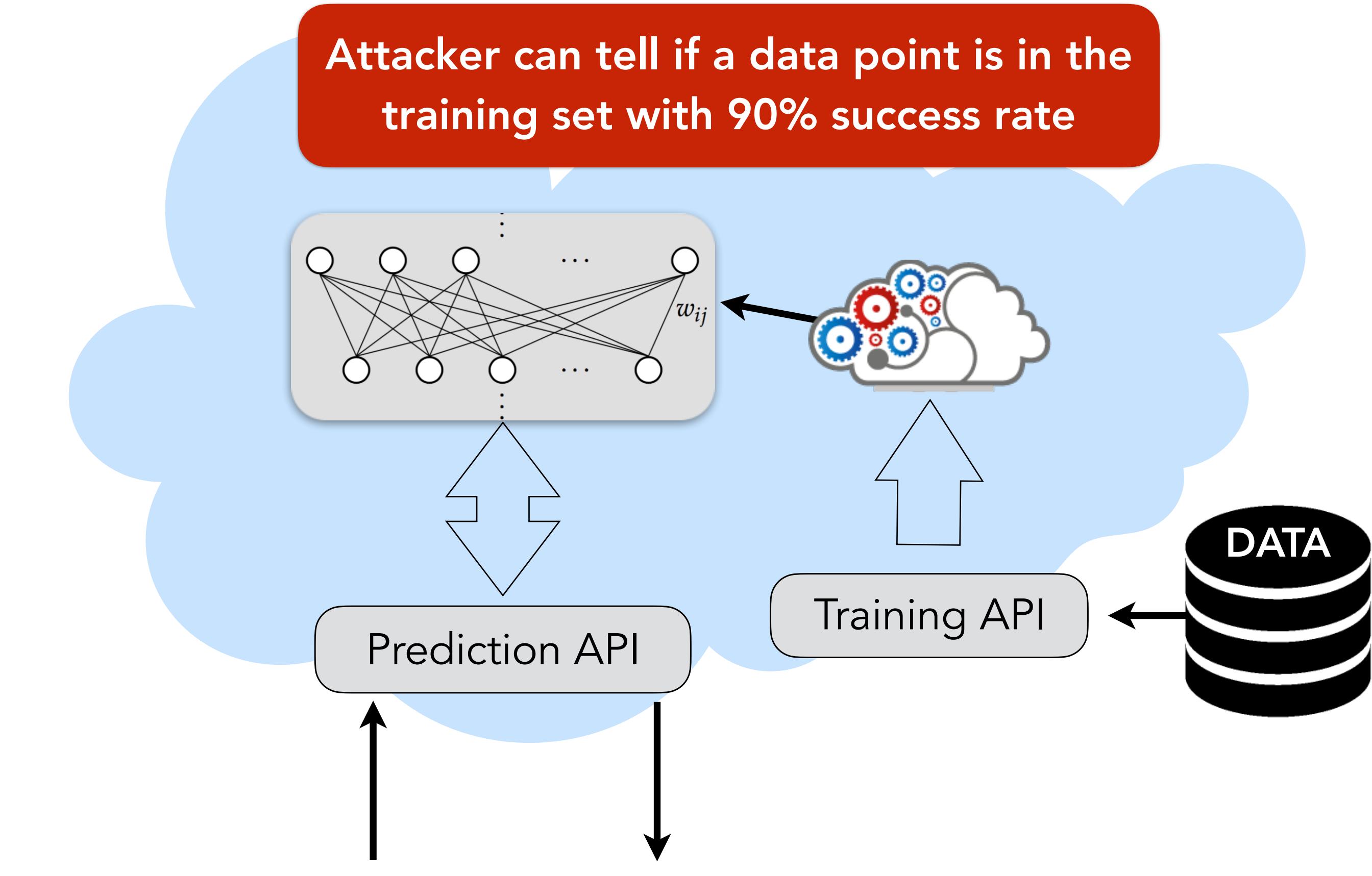
# Indirect Privacy Risks in Machine Learning



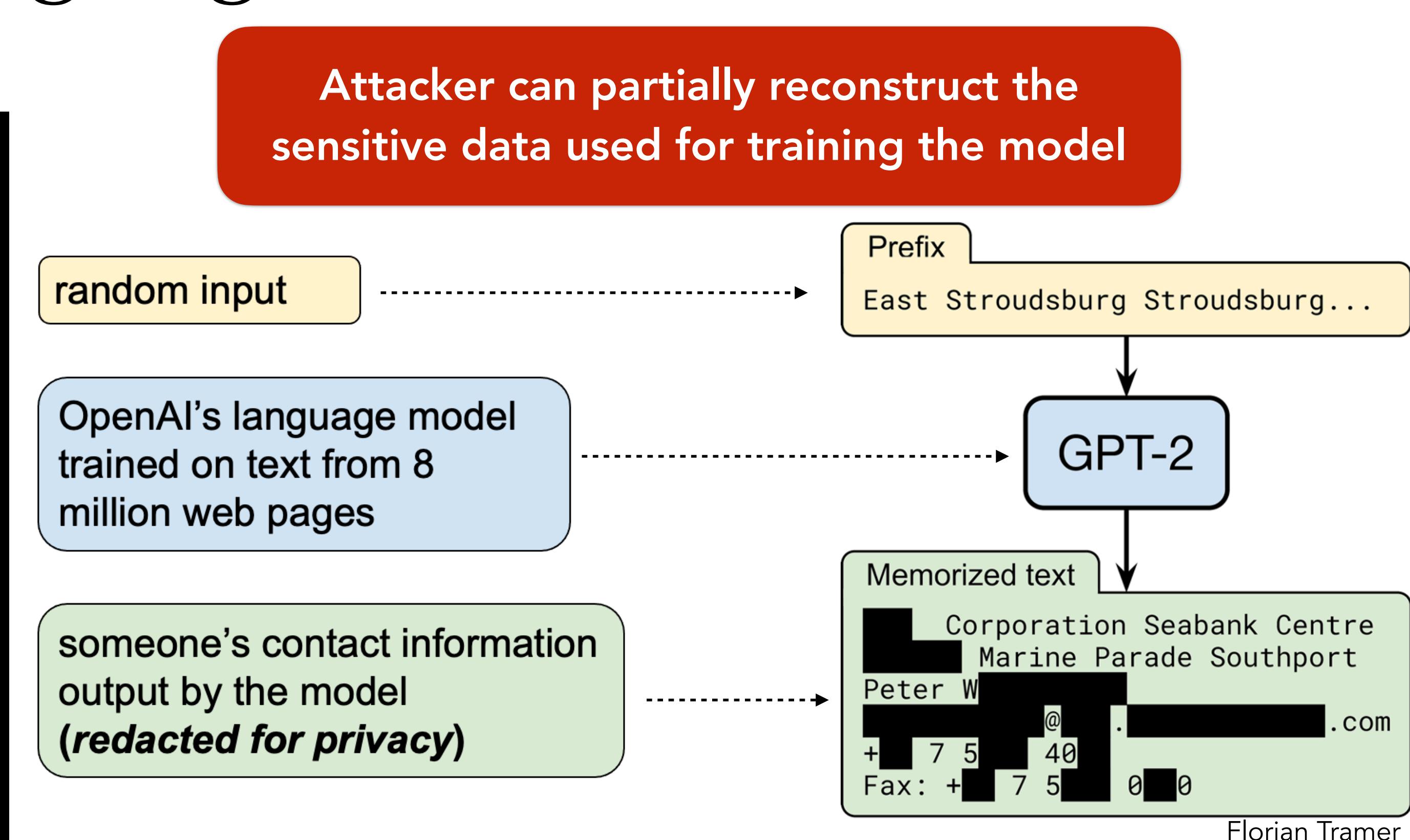
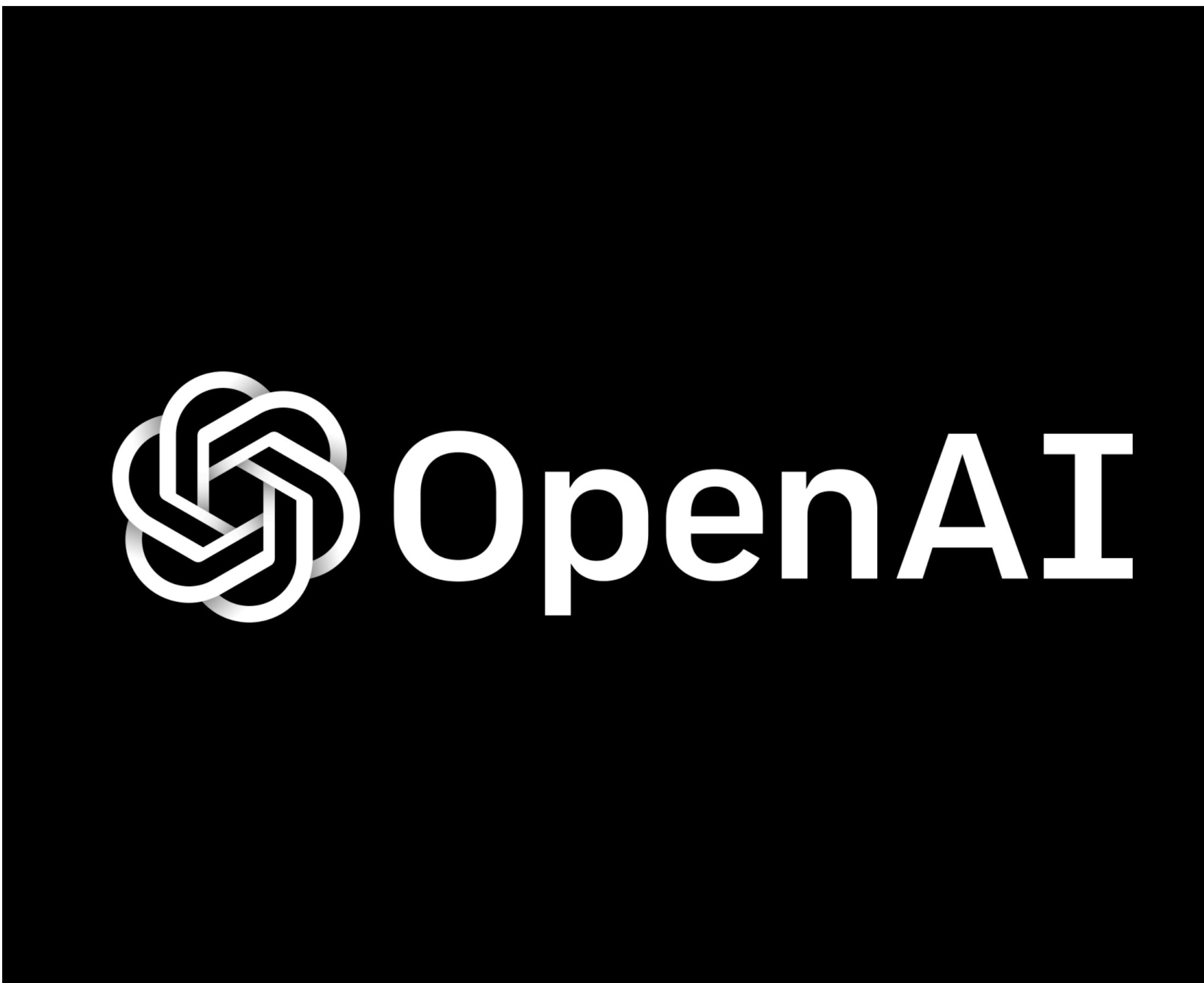
# Real World Attacks against Machine Learning as a Service Platforms



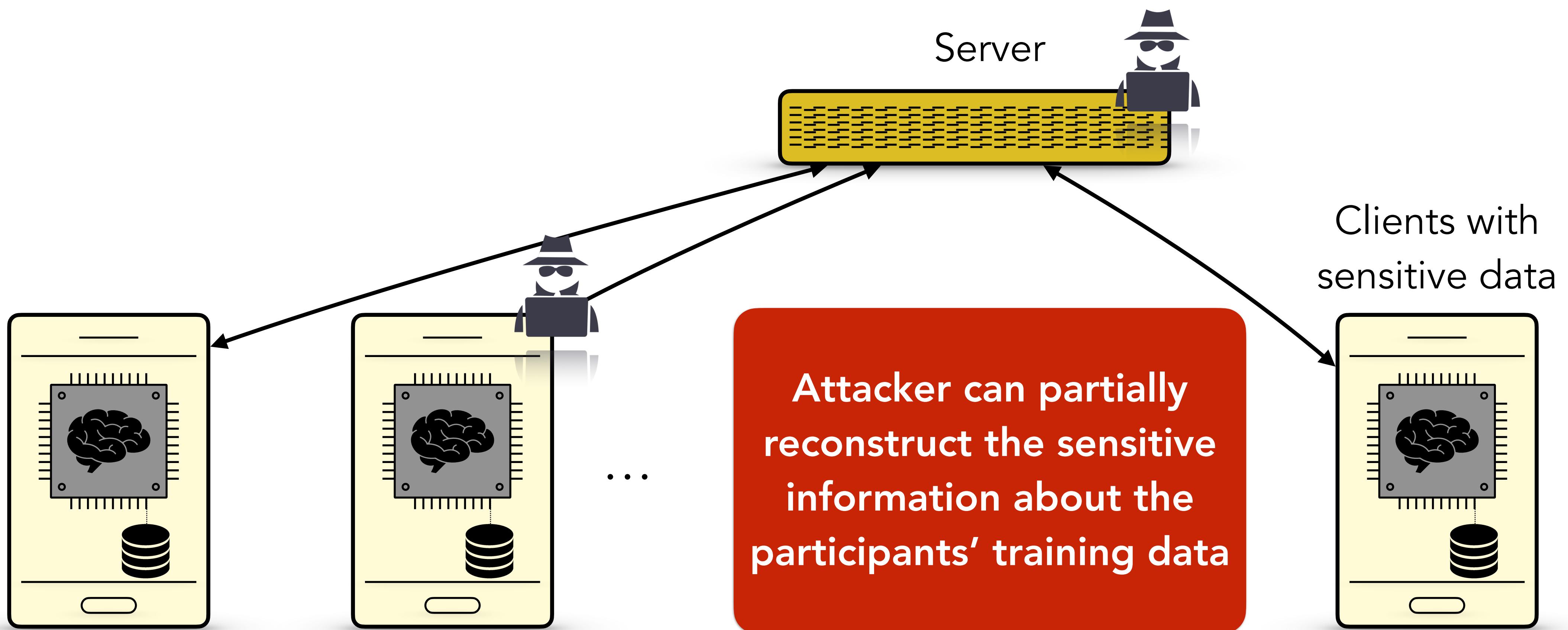
Google Cloud



# Real World Attacks against Large Language Models



# Real World Attacks against Federated Learning Algorithms



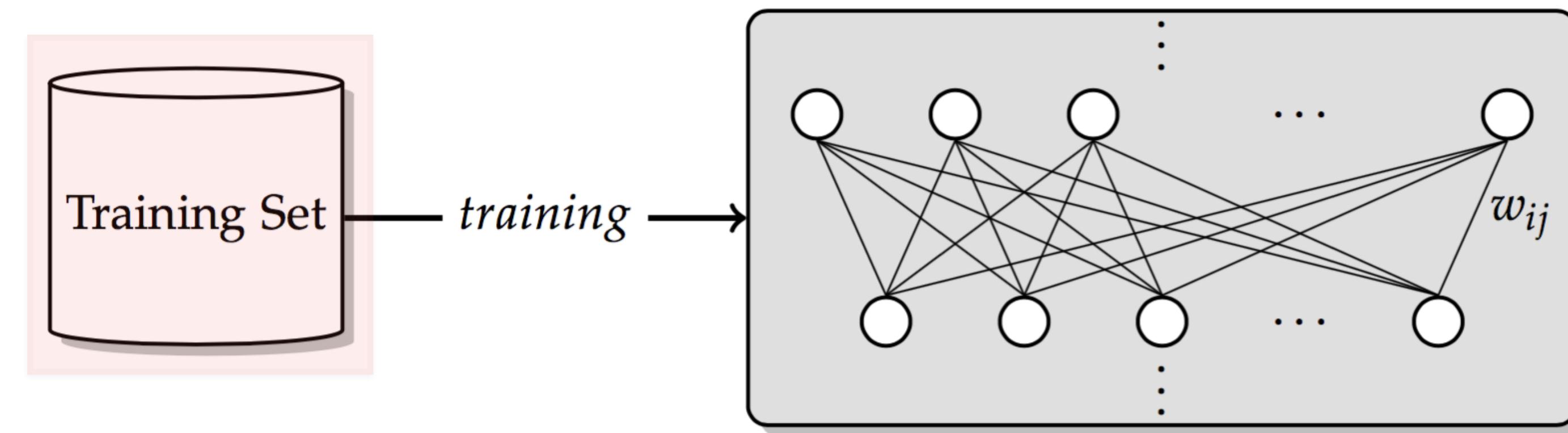
[Nasr, Shokri, Houmansadr] Comprehensive Privacy Analysis of Deep Learning: Passive and Active White-box Inference Attacks against Centralized and Federated Learning, SP'19

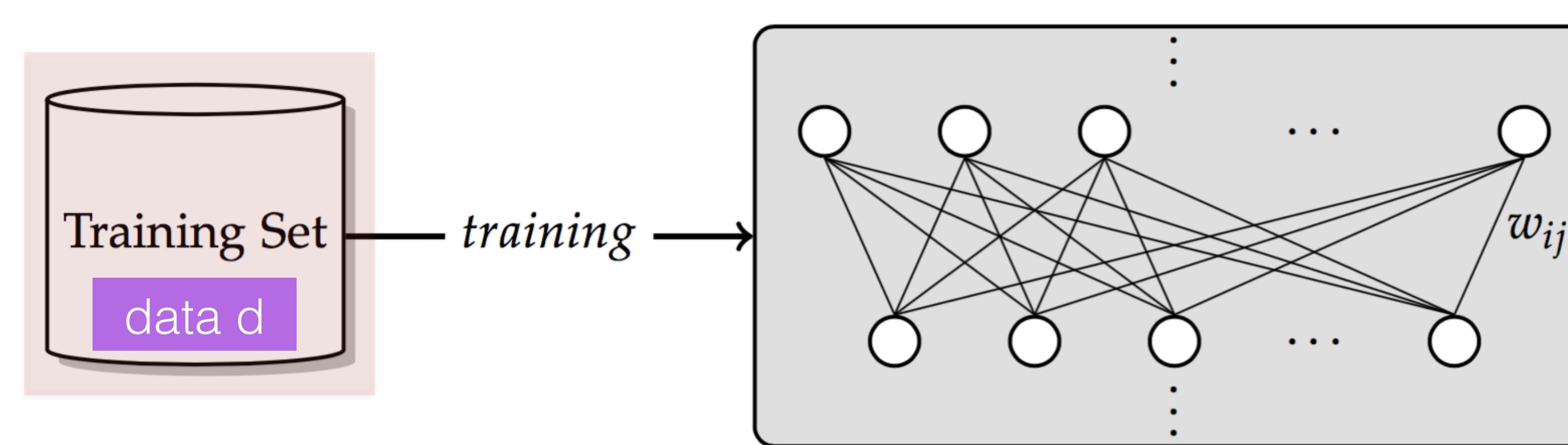
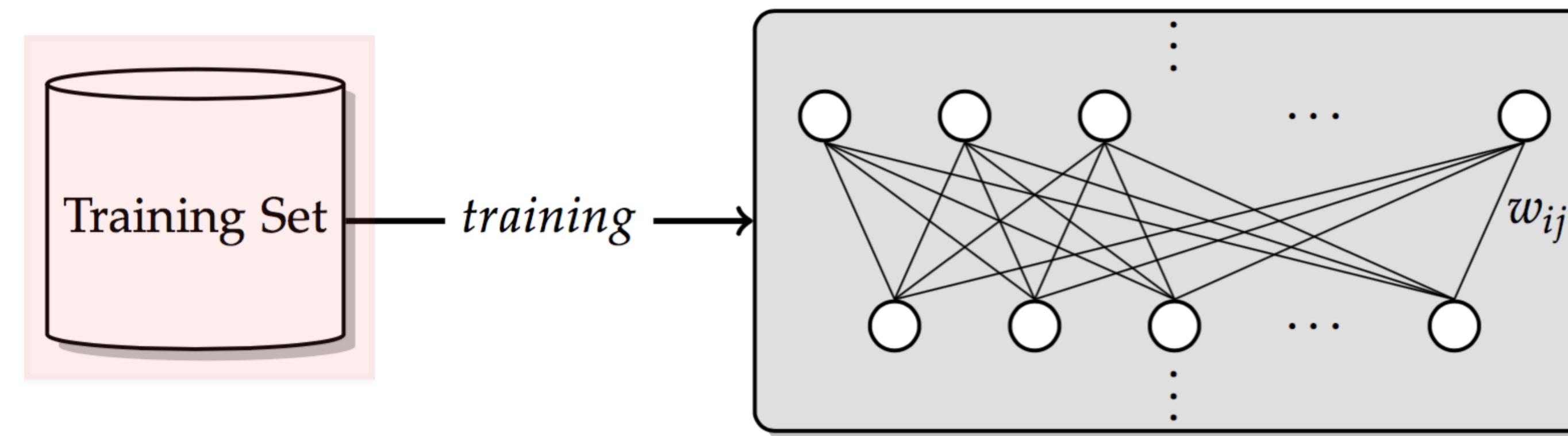
[Melis, Song, De Cristofaro, Shmatikov] Exploiting Unintended Feature Leakage in Collaborative Learning, SP'19

[Zhang, Tople, Ohrimenko] Leakage of Dataset Properties in Multi-Party Machine Learning, Usenix Security'21

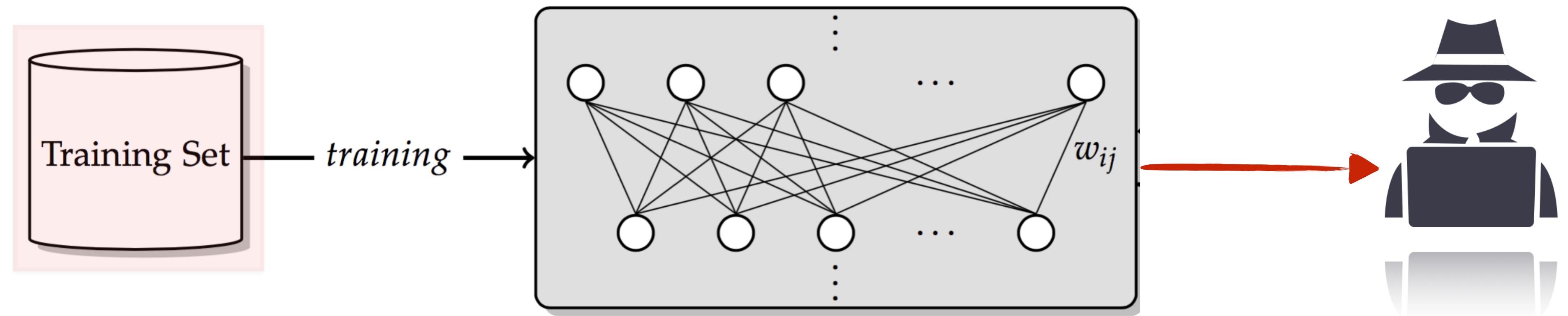
Models are **personal data**

We need a standard method for auditing data  
privacy in machine learning systems

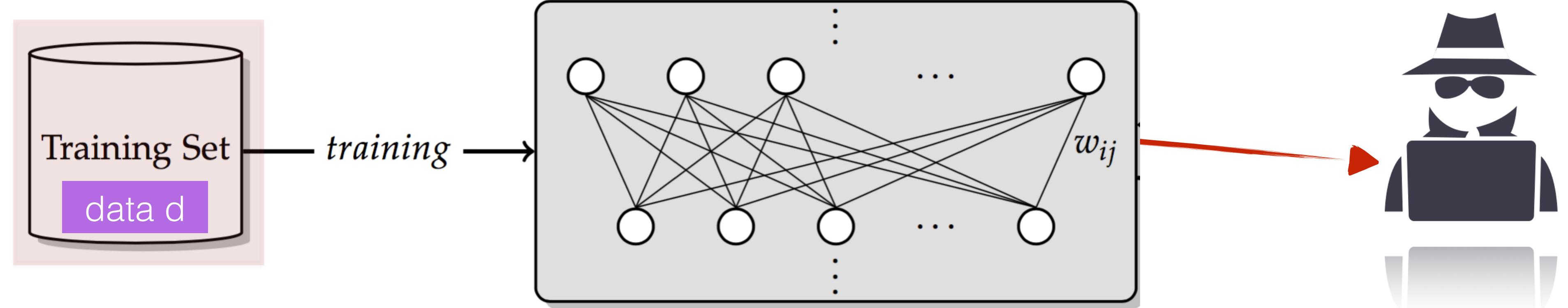




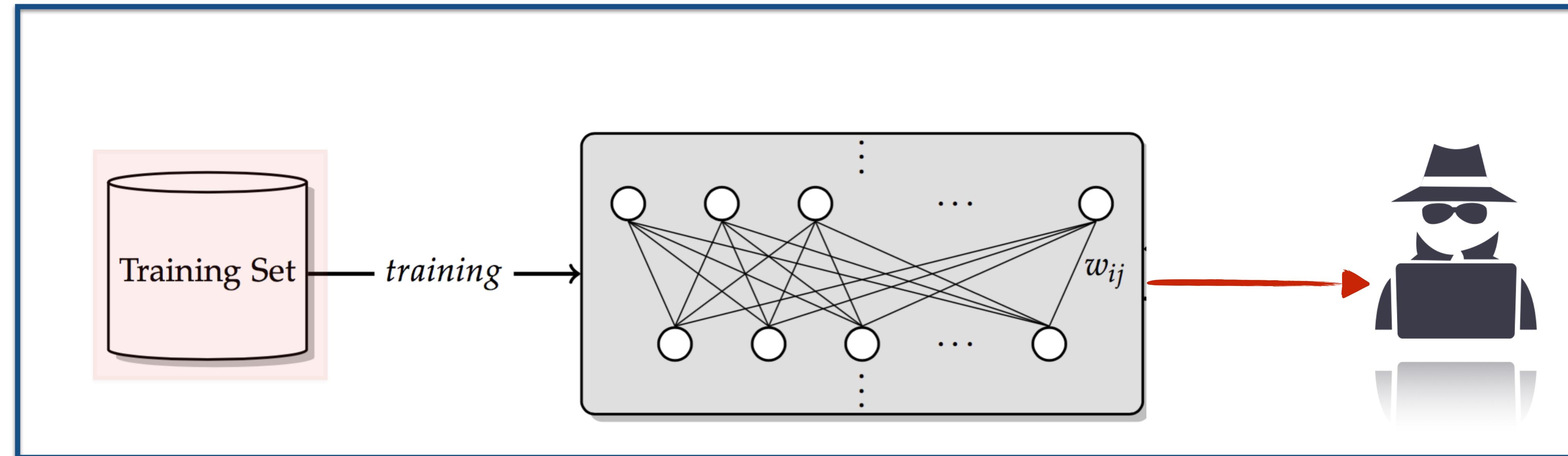
A world in which  
data  $d$  **is not** part  
of training set



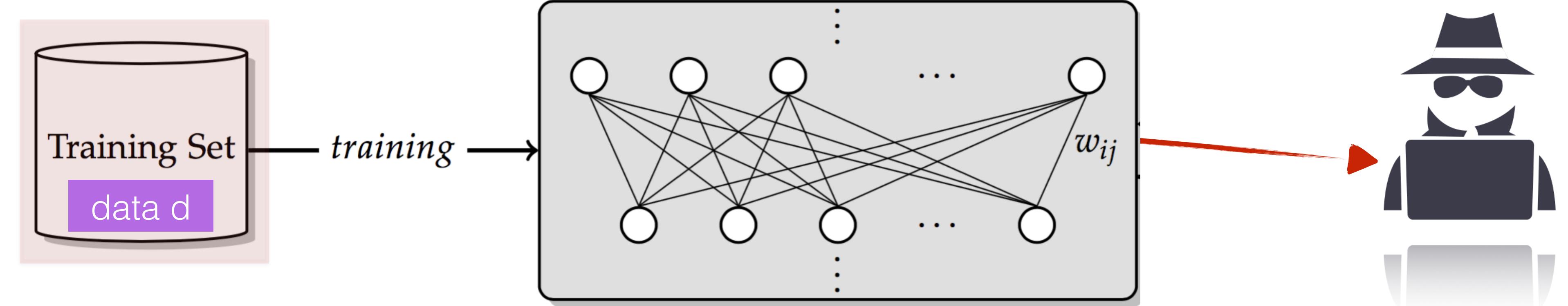
Alternative world  
where data  $d$  **is**  
part of training set



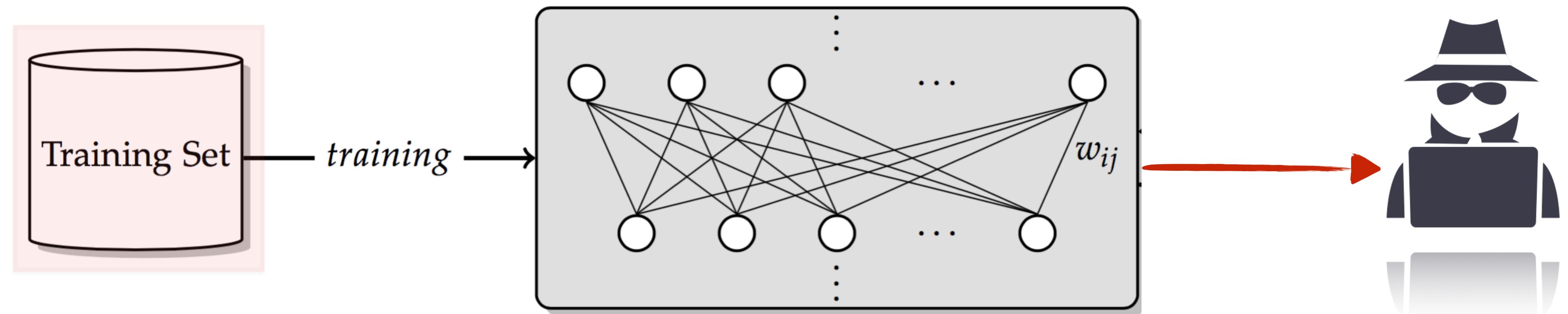
A world in which  
data  $d$  **is not** part  
of training set



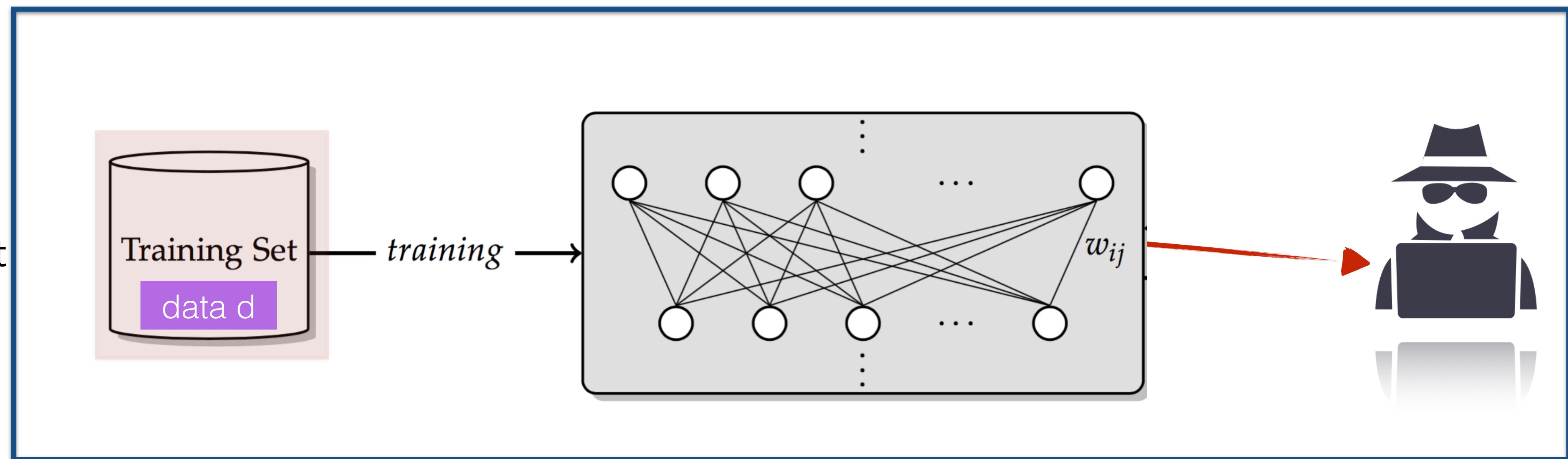
Alternative world  
where data  $d$  **is**  
part of training set



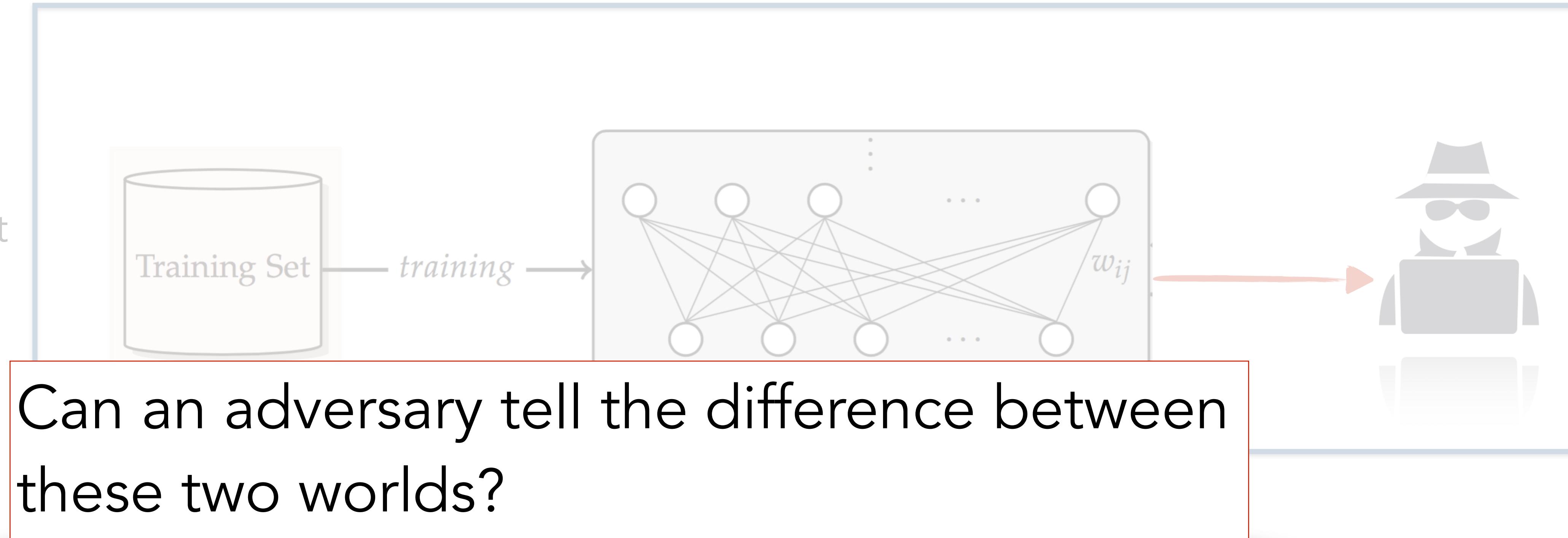
A world in which  
data  $d$  **is not** part  
of training set



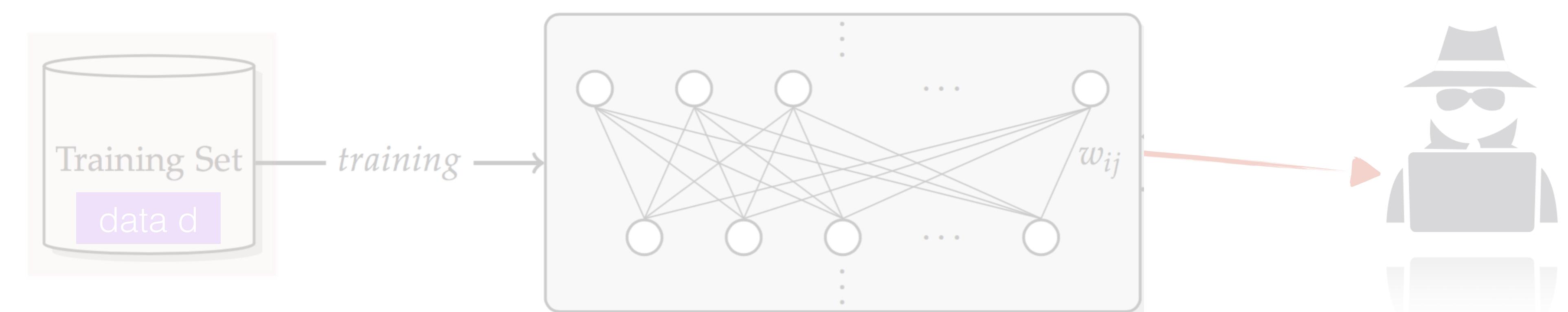
Alternative world  
where data  $d$  **is**  
part of training set



A world in which  
data  $d$  **is not** part  
of training set

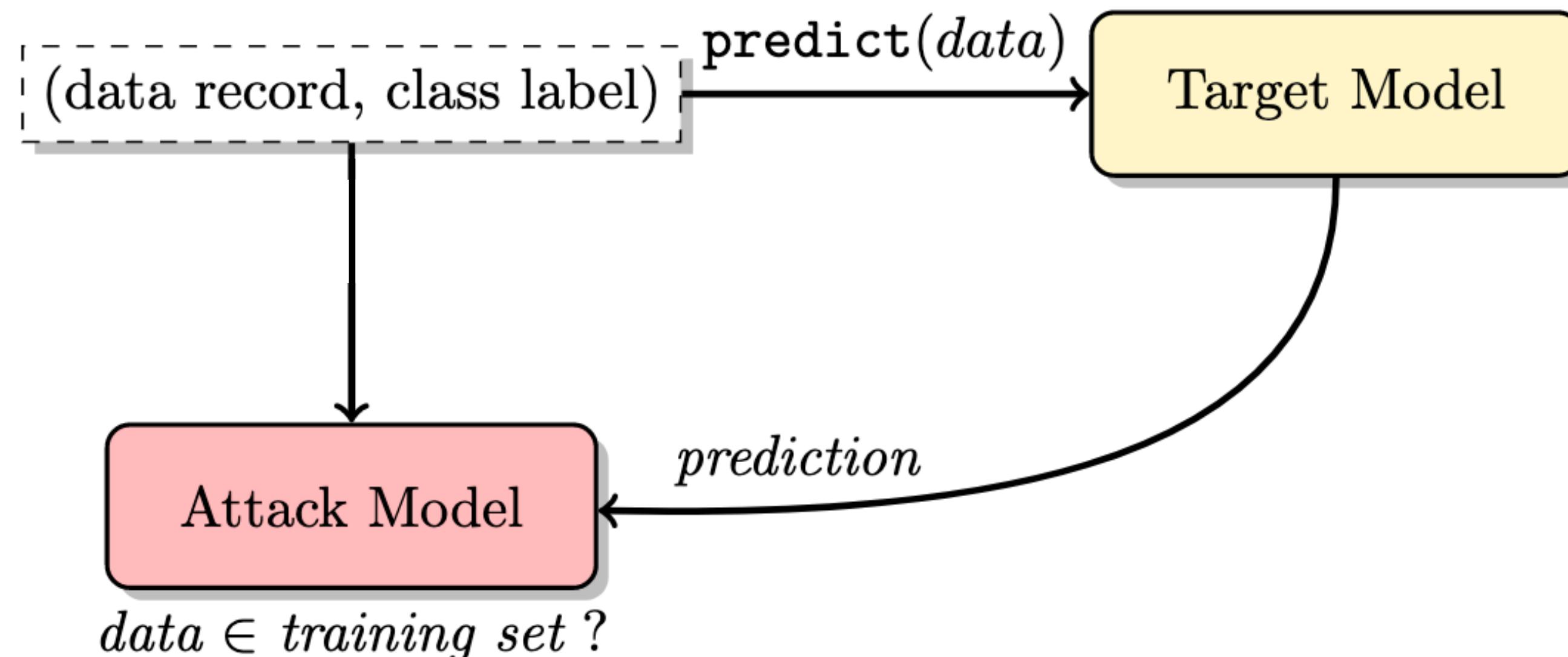


Alternative world  
where data  $d$  **is**  
part of training set

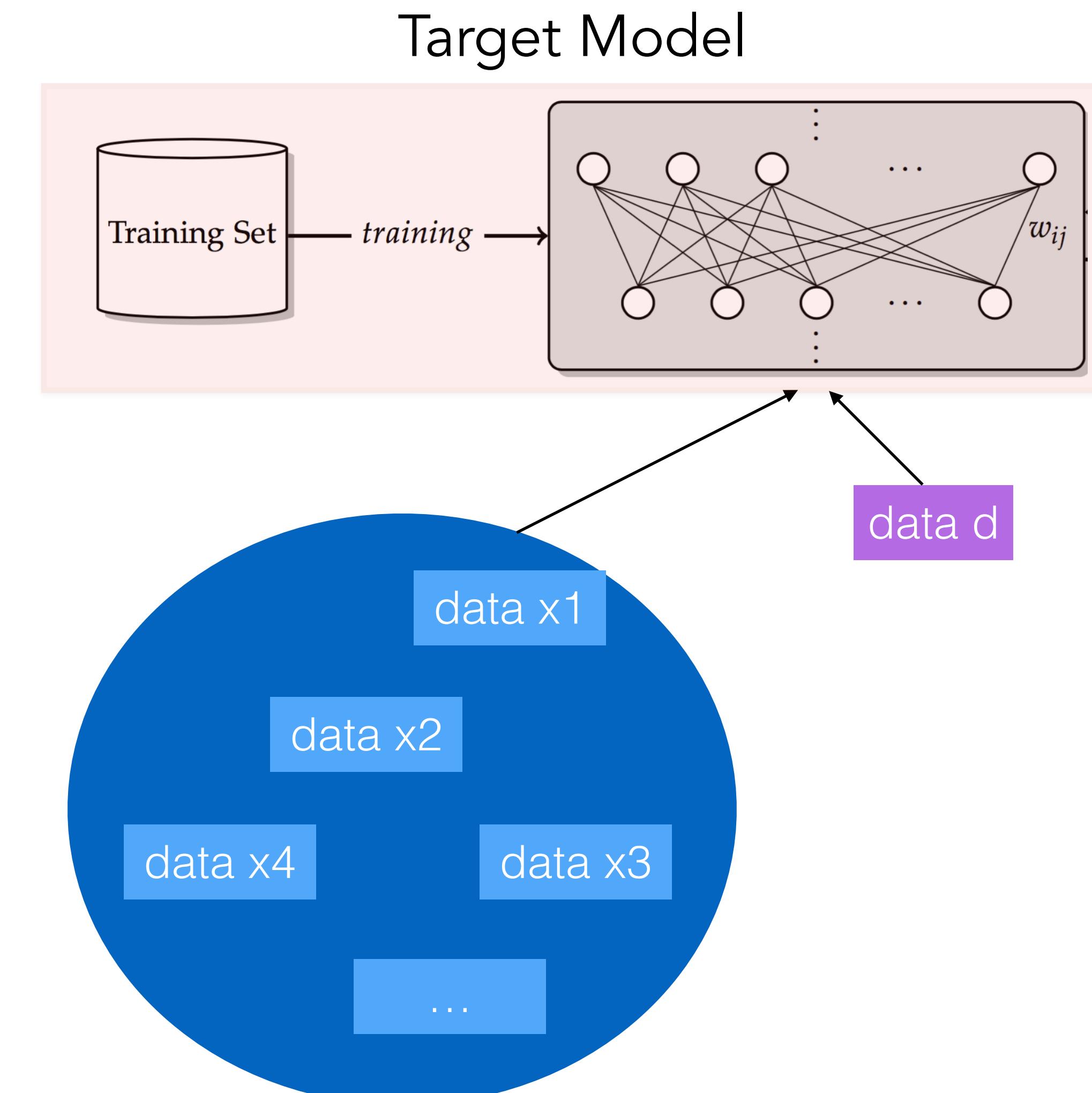


# Membership Inference

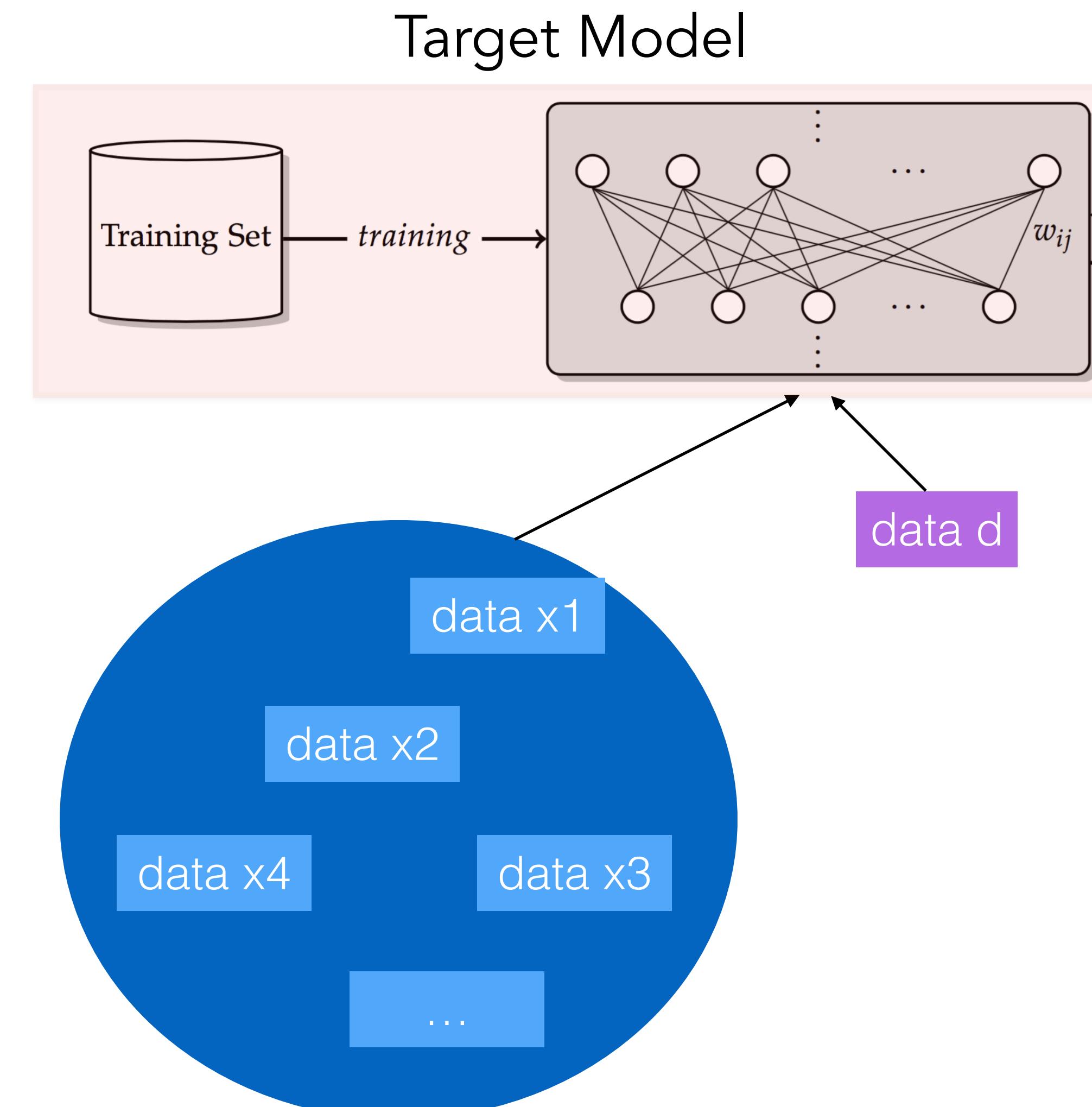
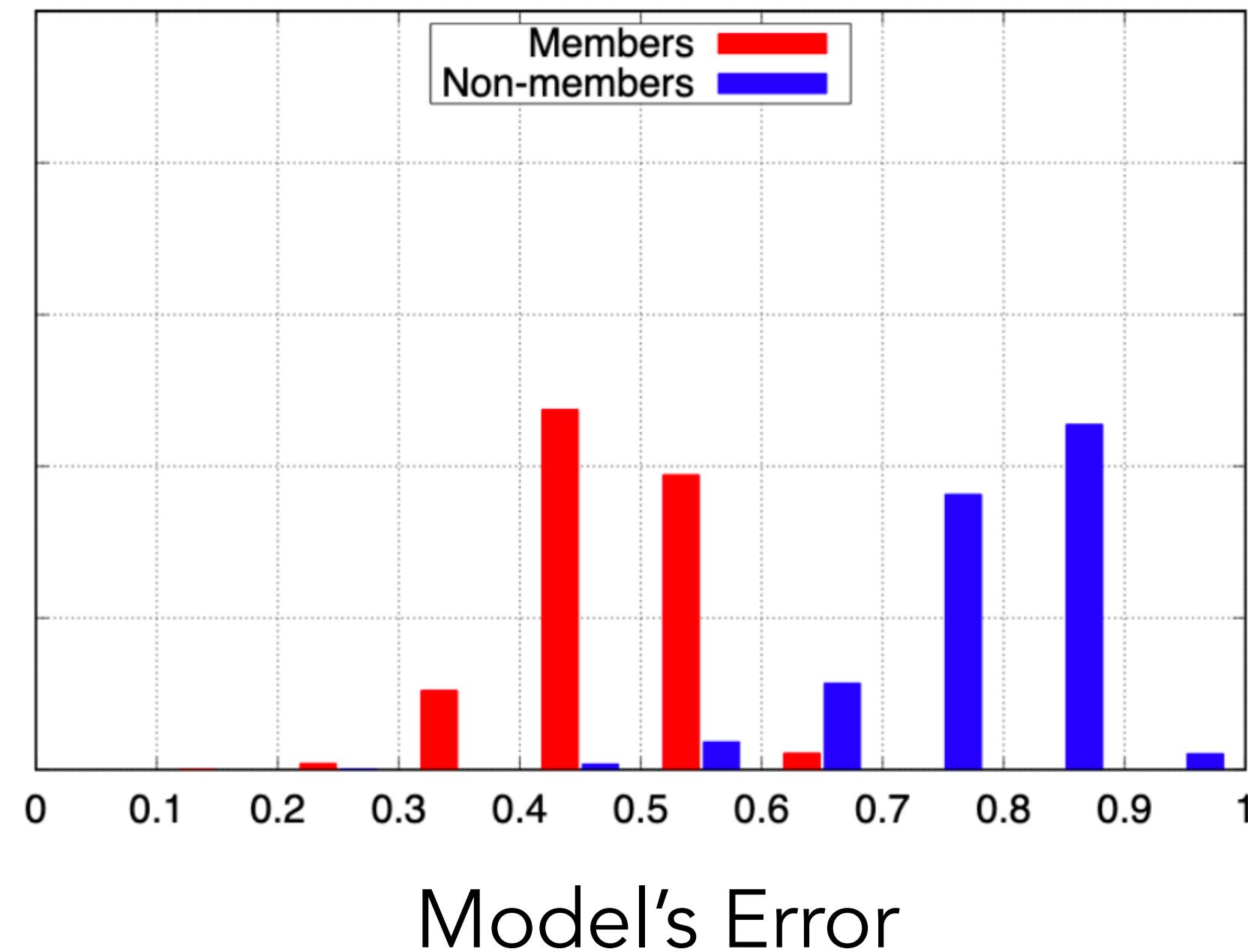
- Given a model, can an adversary infer whether a particular data point is part of its training set?
- Success of attacker is a metric for privacy loss



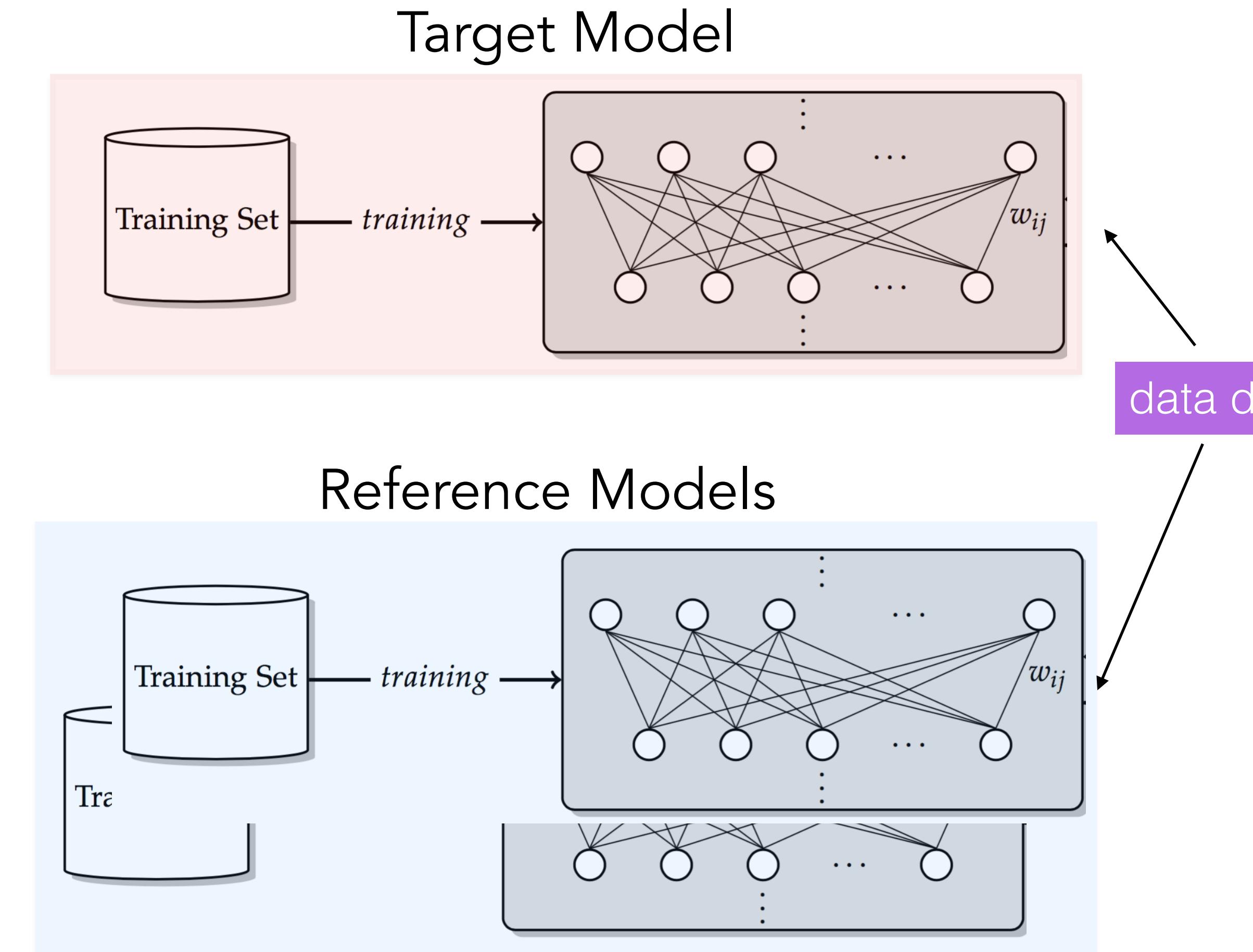
# Membership Inference Attack



# Membership Inference Attack



# Membership Inference Attack



Success rate of adversary indicates information leakage of models about their training data

# AI Regulations and Guidelines

## A Taxonomy and Terminology of Adversarial Machine Learning

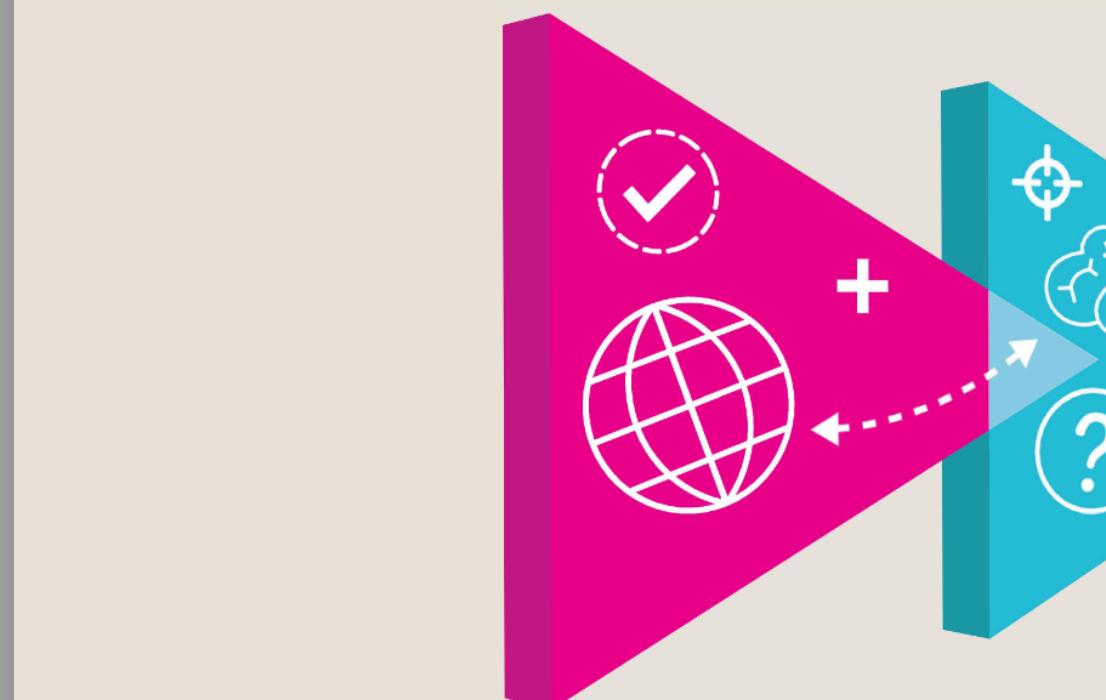
Elham Tabassi  
 Kevin J. Burns  
 Michael Hadjimichael  
 Andres D. Molina-Markham  
 Julian T. Sexton

This publication is available free of charge from:  
<https://doi.org/10.6028/NIST.IR.8269-draft>



## Guidance on the AI auditing framework

Draft guidance for consultation



 European Commission

Home > Publications > White Paper on Artificial Intelligence: a European approach to excellence and trust

WHITE PAPER

**White Paper on Artificial Intelligence: a European approach to excellence and trust**



EXECUTIVE OFFICE OF THE PRESIDENT  
 OFFICE OF MANAGEMENT AND BUDGET  
 WASHINGTON, D.C. 20503

November 17, 2020

THE DIRECTOR  
 M-21-06

MEMORANDUM FOR THE HEADS OF EXECUTIVE DEPARTMENTS AND AGENCIES

FROM: Russell T. Vought  
 Director

SUBJECT: Guidance for Regulation of Artificial Intelligence Applications



# AI Regulations and Guidelines

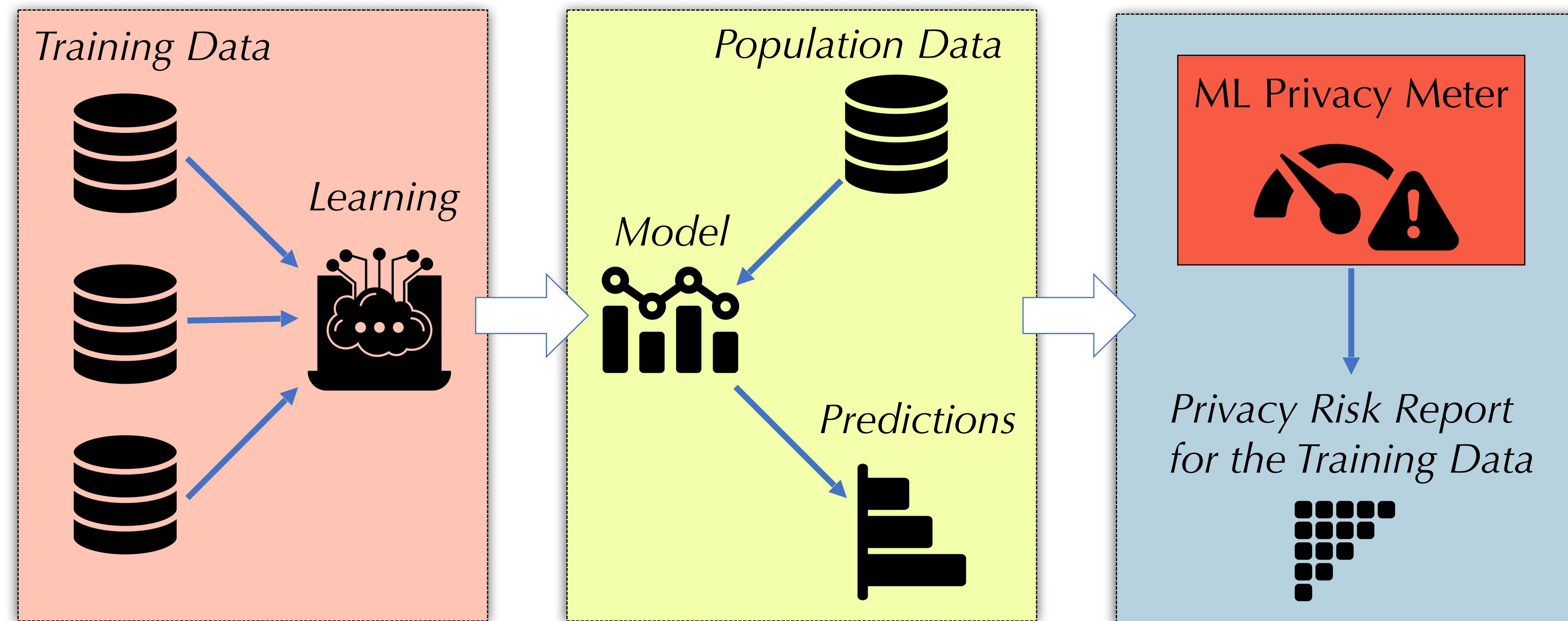
A Taxonomy and Terminology of Adversarial Machine Learning



- "... membership inferences show that AI models can inadvertently contain **personal data**"
- "Attacks that reveal confidential information about the data include membership inference ..."
- "... **should consider the risks to data throughout the design, development, and operation of an AI system**"

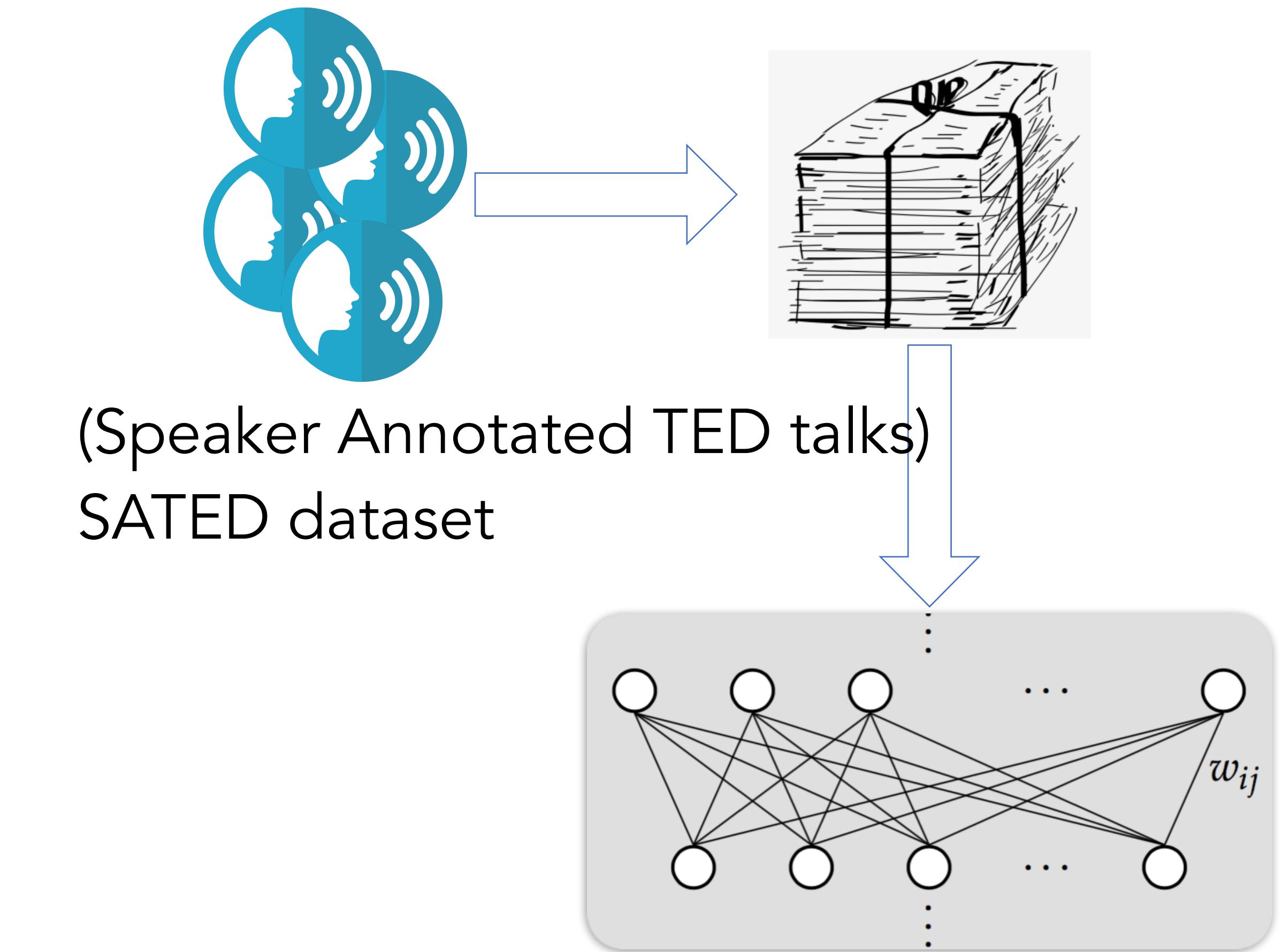
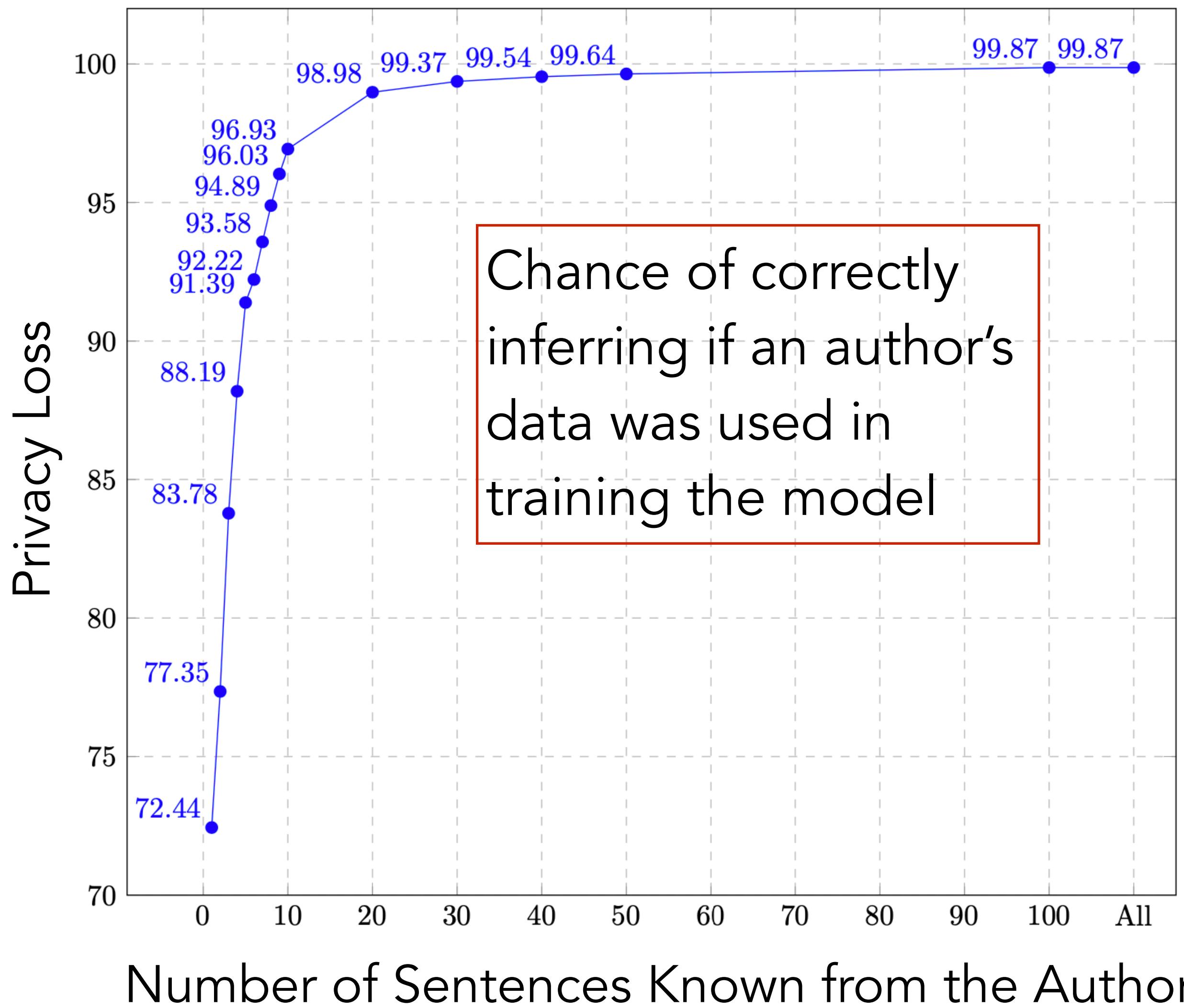
# ML Privacy Meter

[privacy-meter.com](http://privacy-meter.com)

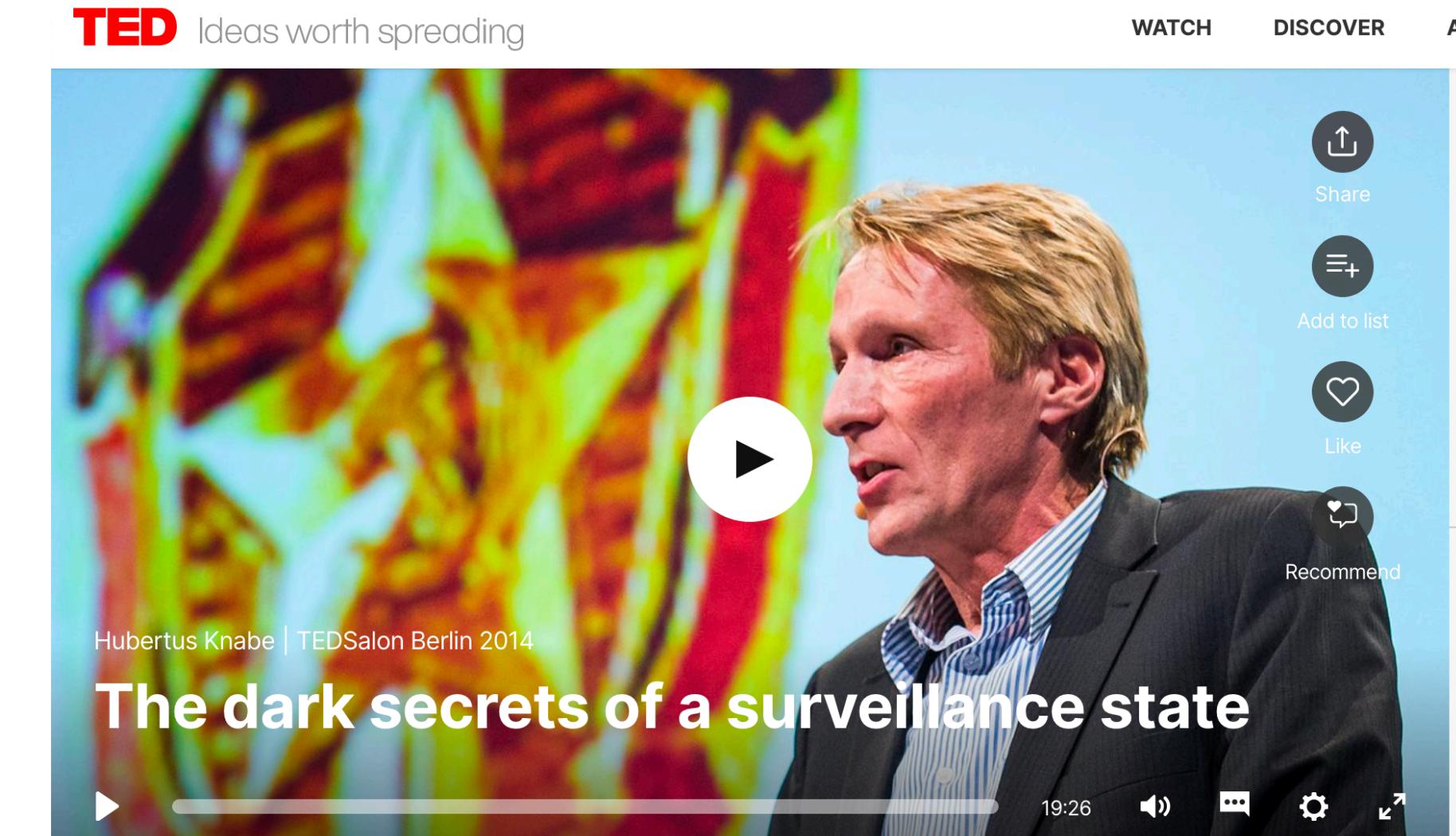
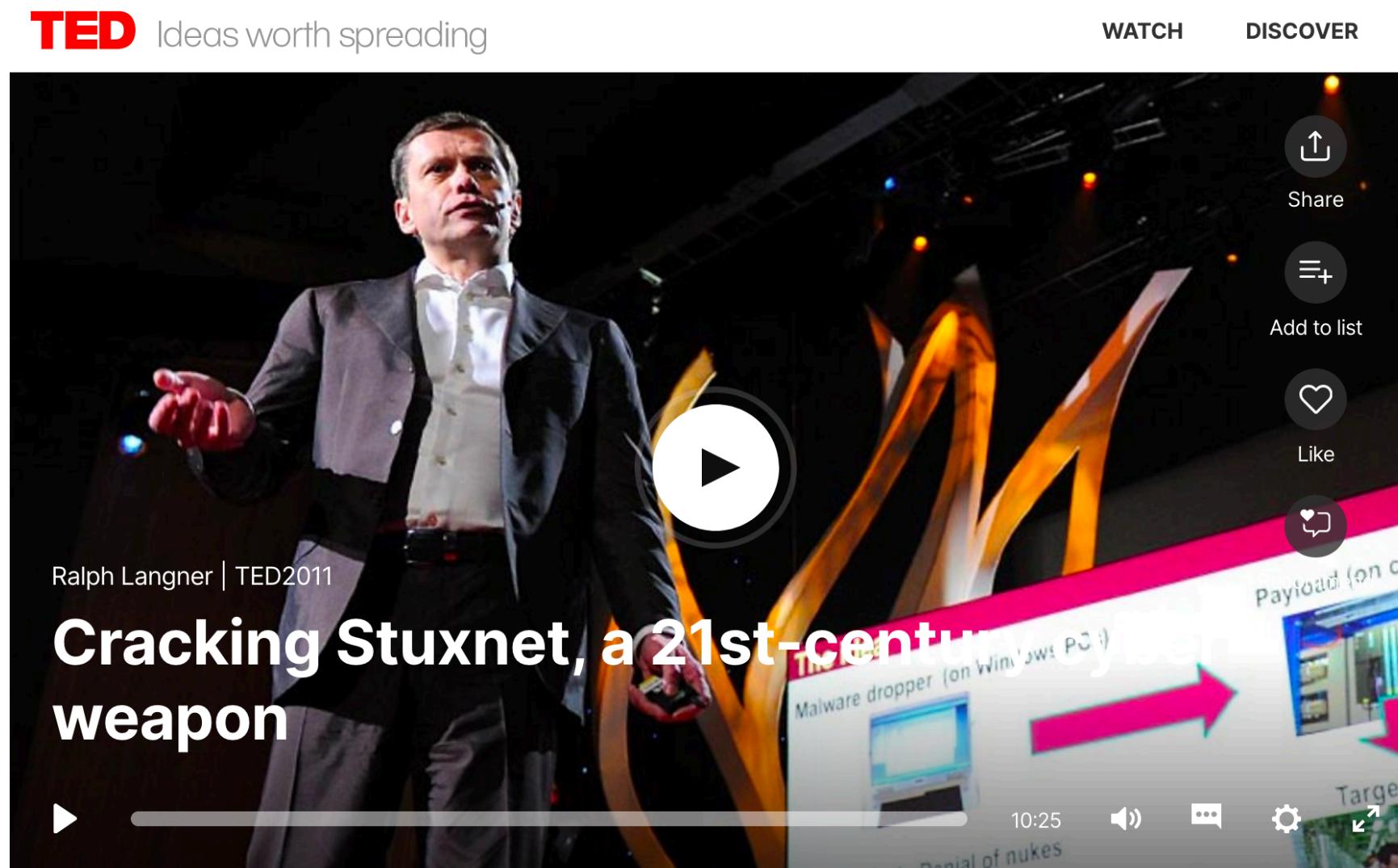


ML Privacy Meter is an open source tool that enables quantifying the privacy risks of machine learning models.

# Example: Language Generative Model



# Examples of Vulnerable Training Data



But it gets worse. And this is very important, what I'm generic. It doesn't have anything to do, in specifics, would work as well, for example, in a power plant or don't have -- as an attacker -- you don't have to deal with the case of Stuxnet. You could also use conventional

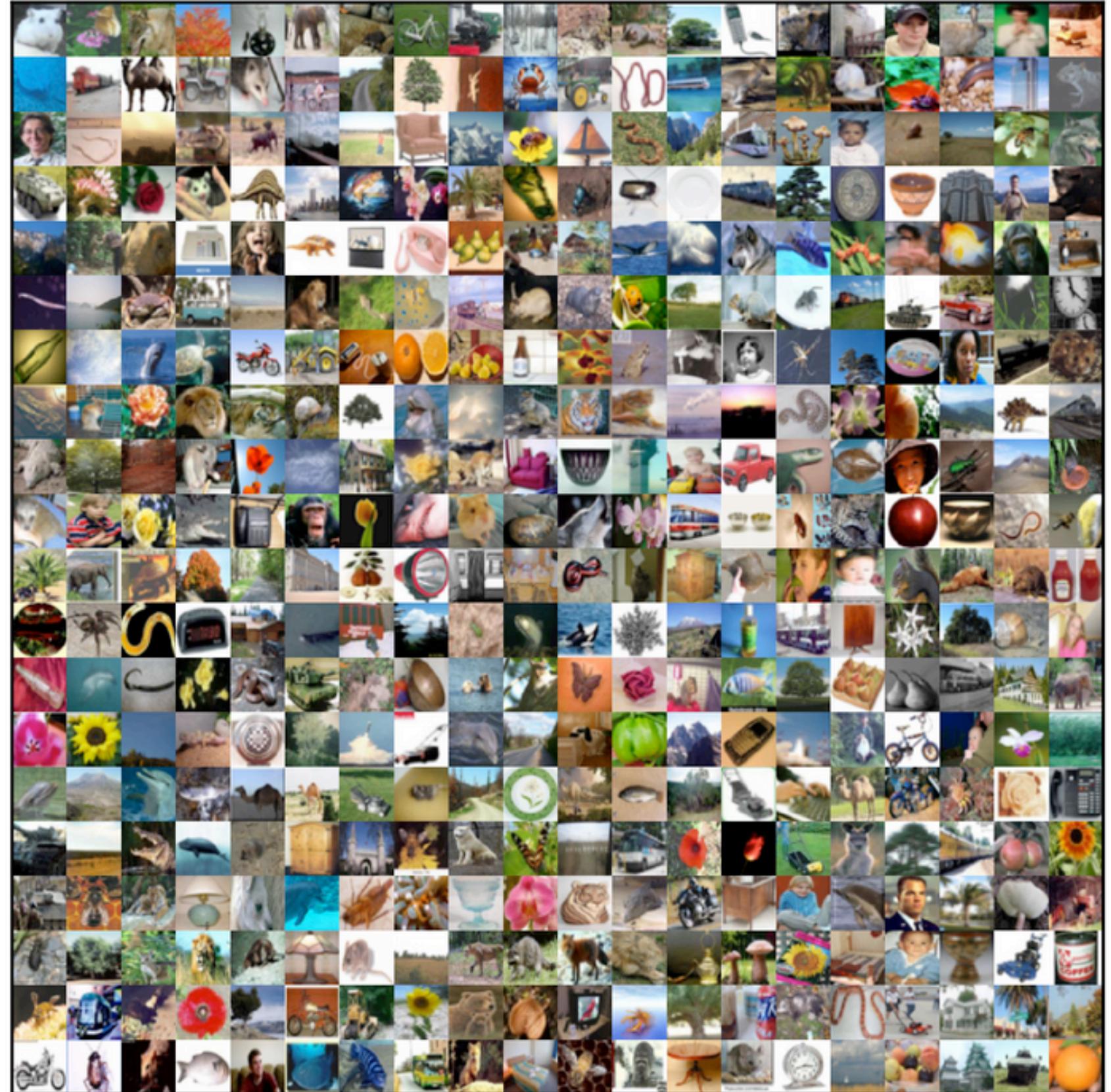
This year, Germany is celebrating the 25th anniversary of the fall of the Berlin Wall in 1989, the Communist regime was moved away, the Berlin Wall fell, and the German Democratic Republic, the GDR, in the East was joined with the West to found today's Germany. Among many changes, the secret police files of the East German secret police, known as the Stasi. Over 17 million files were opened to the public, and historians such as me have been able to learn a lot about how the GDR surveillance state functioned.

# Example: Image Classification Tasks

CIFAR100 Image classification

Model	Number of Parameters	Prediction (Test) Accuracy	Privacy Risk
AlexNet	2.47 million	44%	75.1%
ResNet	1.7 million	73%	64.3%
DenseNet	25.62 million	82%	74.3%

Large capacity      High generalizability      Low privacy



# Membership Inference Attacks

Enhanced Membership Inference Attacks  
against Machine Learning Models

Jiayuan Ye<sup>1</sup>, Aadyaa Maddi<sup>1</sup>, Sasi Kumar Murakonda<sup>2</sup>, and Reza Shokri<sup>1</sup>

# Hypothesis Testing for Membership Inference

- Given a data point “z” and black-box access to model “ $\theta$ ”,
- Tell if “z” was member of the training set of “ $\theta$ ”
- Null hypothesis: non-member
- Alternative hypothesis: member

$$H_0 : D \xleftarrow{n \text{ i.i.d. samples} \sim \pi(z)} D_{pop}, \theta \sim \mathcal{T}(D), z \xleftarrow{\text{sample} \sim \pi(z)} D_{pop}$$

$$H_1 : D \xleftarrow{n \text{ i.i.d. samples} \sim \pi(z)} D_{pop}, \theta \sim \mathcal{T}(D), z \xleftarrow{\text{sample}} D$$

# Likelihood Ratio Test

$$LR(\theta, z) = \frac{L(H_0|\theta, z)}{L(H_1|\theta, z)} \approx e^{\frac{1}{T}\ell(\theta, x_z, y_z)}$$

# Likelihood Ratio Test

$$LR(\theta, z) = \frac{L(H_0|\theta, z)}{L(H_1|\theta, z)} \approx e^{\frac{1}{T}\ell(\theta, x_z, y_z)}$$

Reject the null hypothesis:  $\{(\theta, z) : LR(\theta, z) \leq c\}$

# Likelihood Ratio Test

$$LR(\theta, z) = \frac{L(H_0|\theta, z)}{L(H_1|\theta, z)} \approx e^{\frac{1}{T}\ell(\theta, x_z, y_z)}$$

Reject the null hypothesis:  $\{(\theta, z) : LR(\theta, z) \leq c\}$

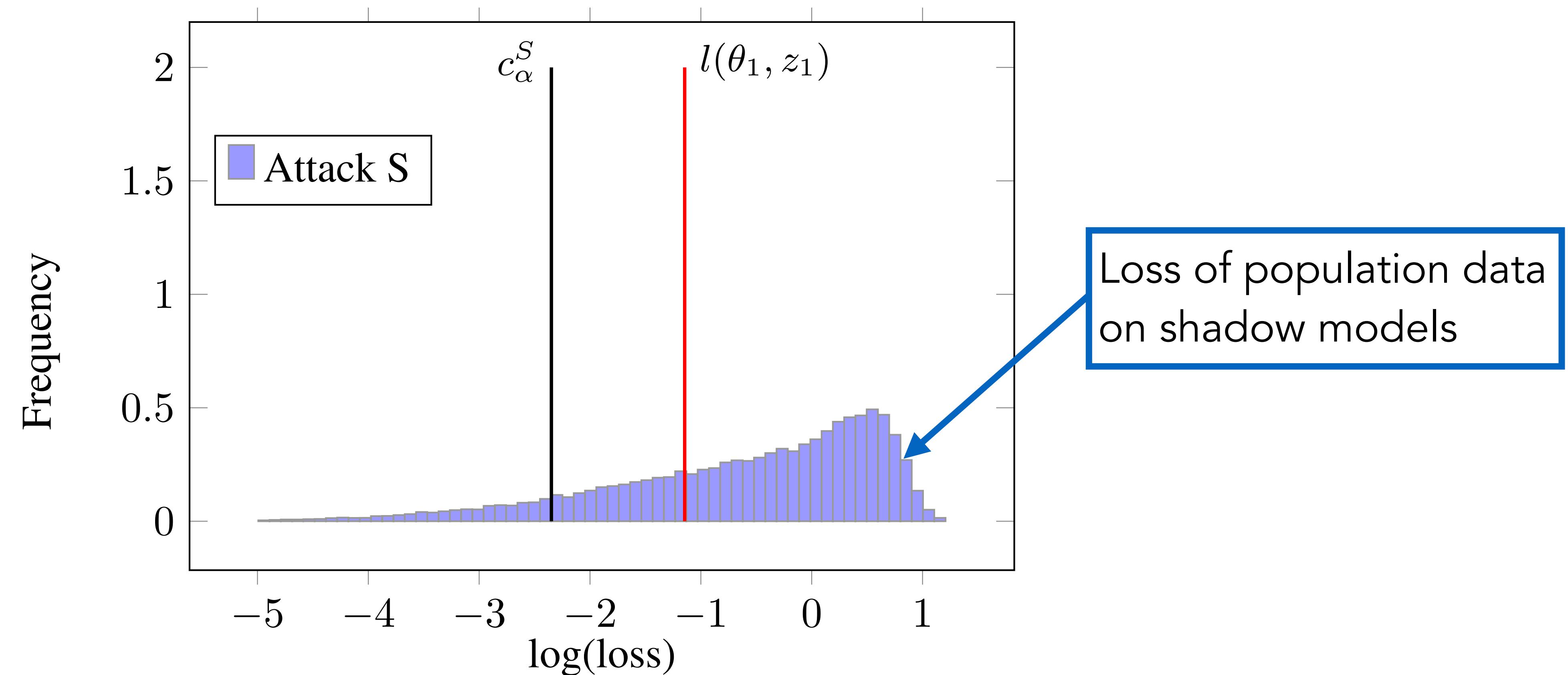
Attack: If  $\ell(\theta, x_z, y_z) \leq c_\alpha(\theta, x_z, y_z)$ , reject  $H_0$

↑  
false positive rate

# Membership Inference via Shadow Models

If  $\ell(\theta, x_z, y_z) \leq c_\alpha(y_z)$ , reject  $H_0$

- A large body of the literature is based on this technique
- Learn a threshold from the behaviour of some shadow models with their test data

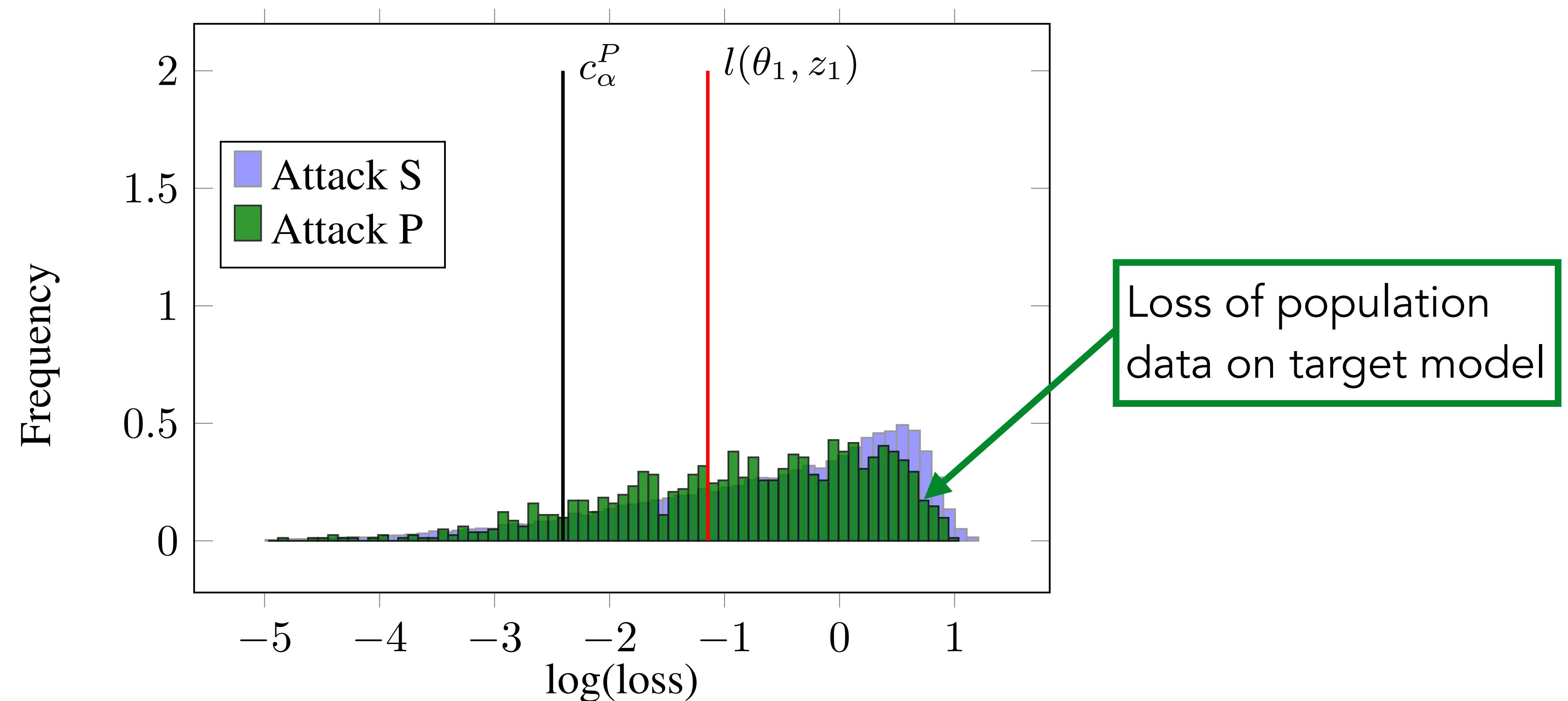


Can we do it more efficiently?

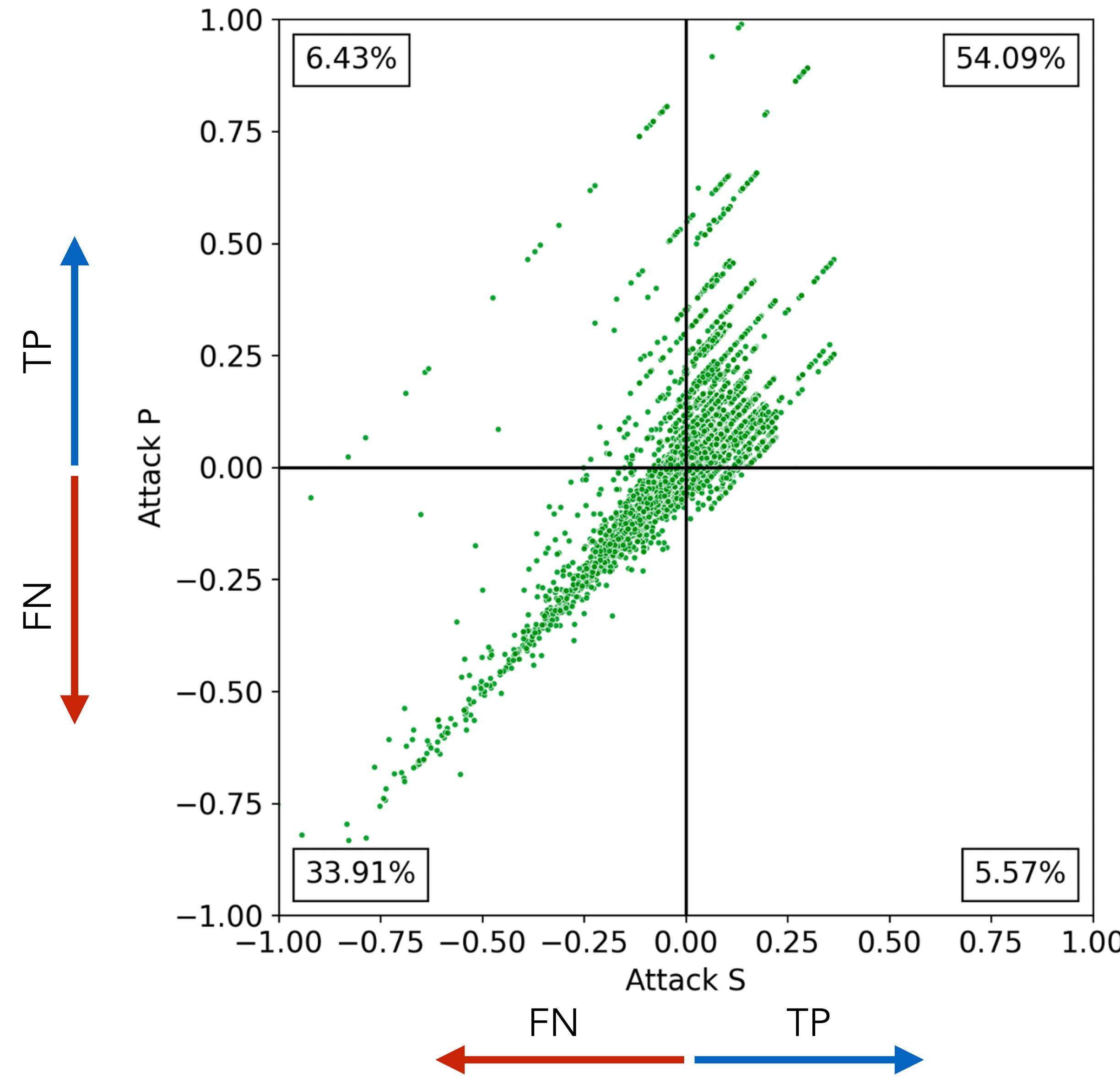
# Membership Inference via Population Data

If  $\ell(\theta, x_z, y_z) \leq c_\alpha(\theta)$ , reject  $H_0$

- Learn a threshold from the behaviour of the target model on population data



# Agreement between Attacks

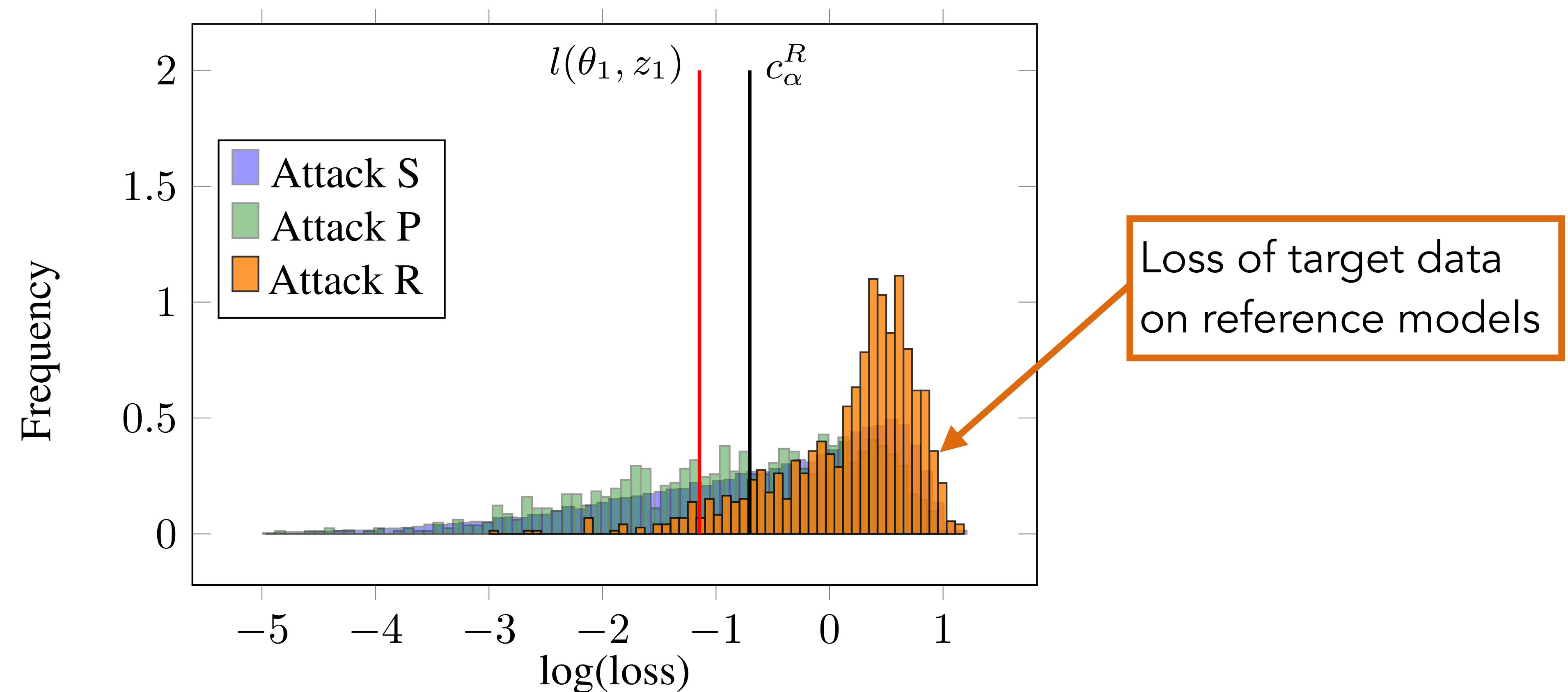


Can we do it more accurately?

# Membership Inference via Reference Models

If  $\ell(\theta, x_z, y_z) \leq c_\alpha(x_z, y_z)$ , reject  $H_0$

- Learn a threshold from the behaviour of target data on reference models



# Can we do it even more accurately?

The objective is to get as close as possible to the leave-one-out attack, where the adversary knows all “other” data in the training set

# Can we do it even more accurately?

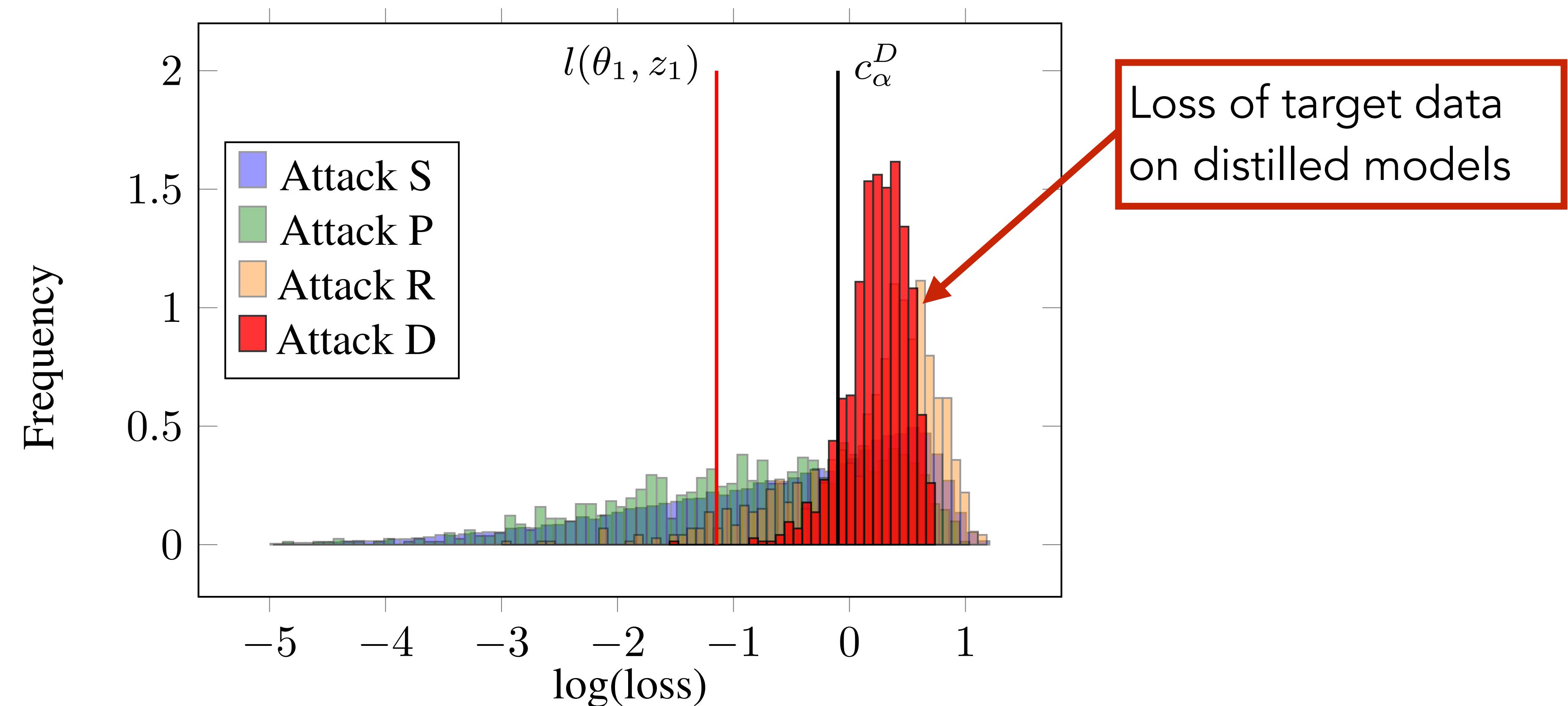
The objective is to get as close as possible to the leave-one-out attack, where the adversary knows all “other” data in the training set

- Train reference models that have a large agreement with the target model on all data, except the target data
- Idea: Model distillation — Reference models are distilled versions of the target models

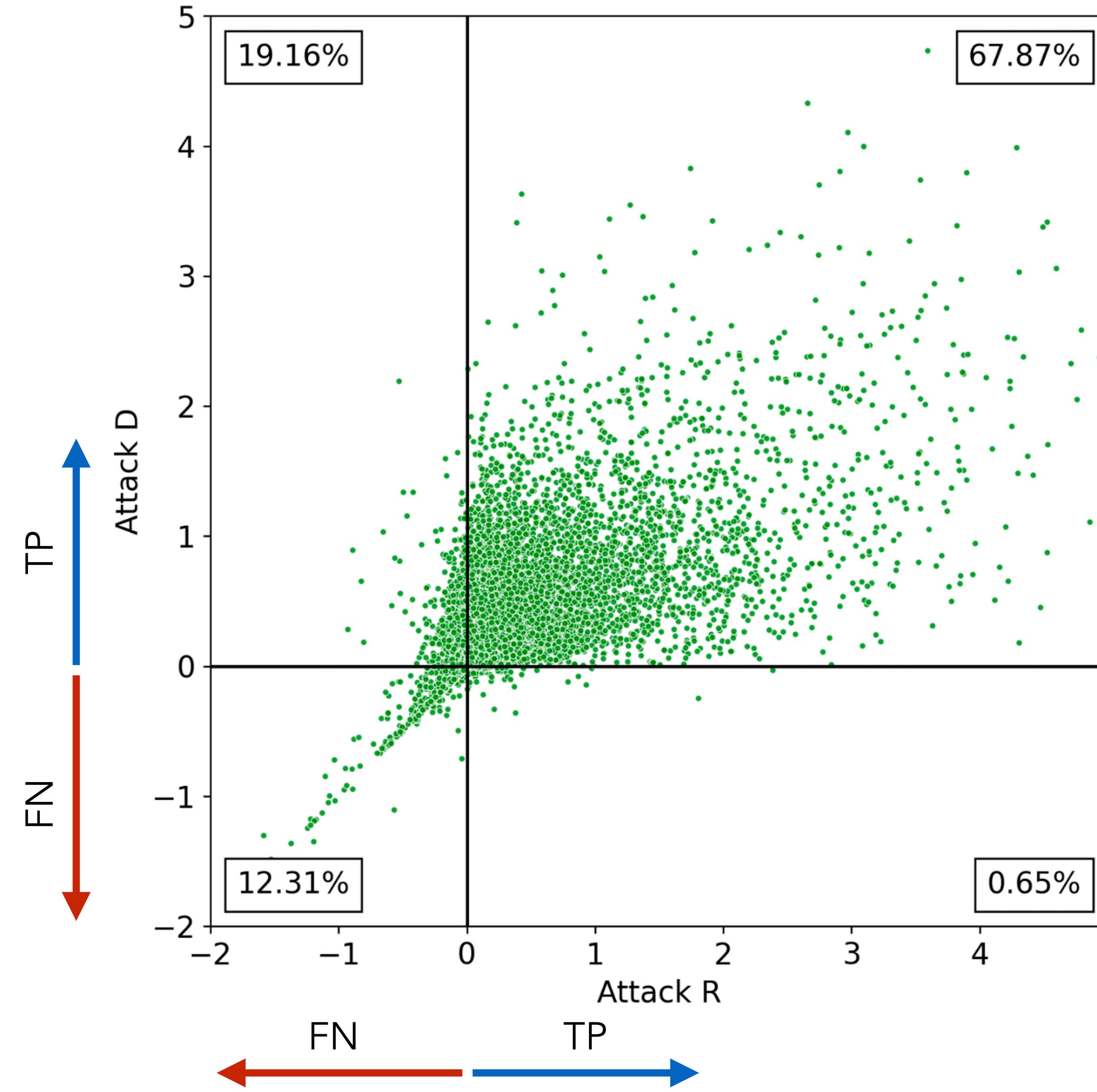
# Membership Inference via Distilled Models

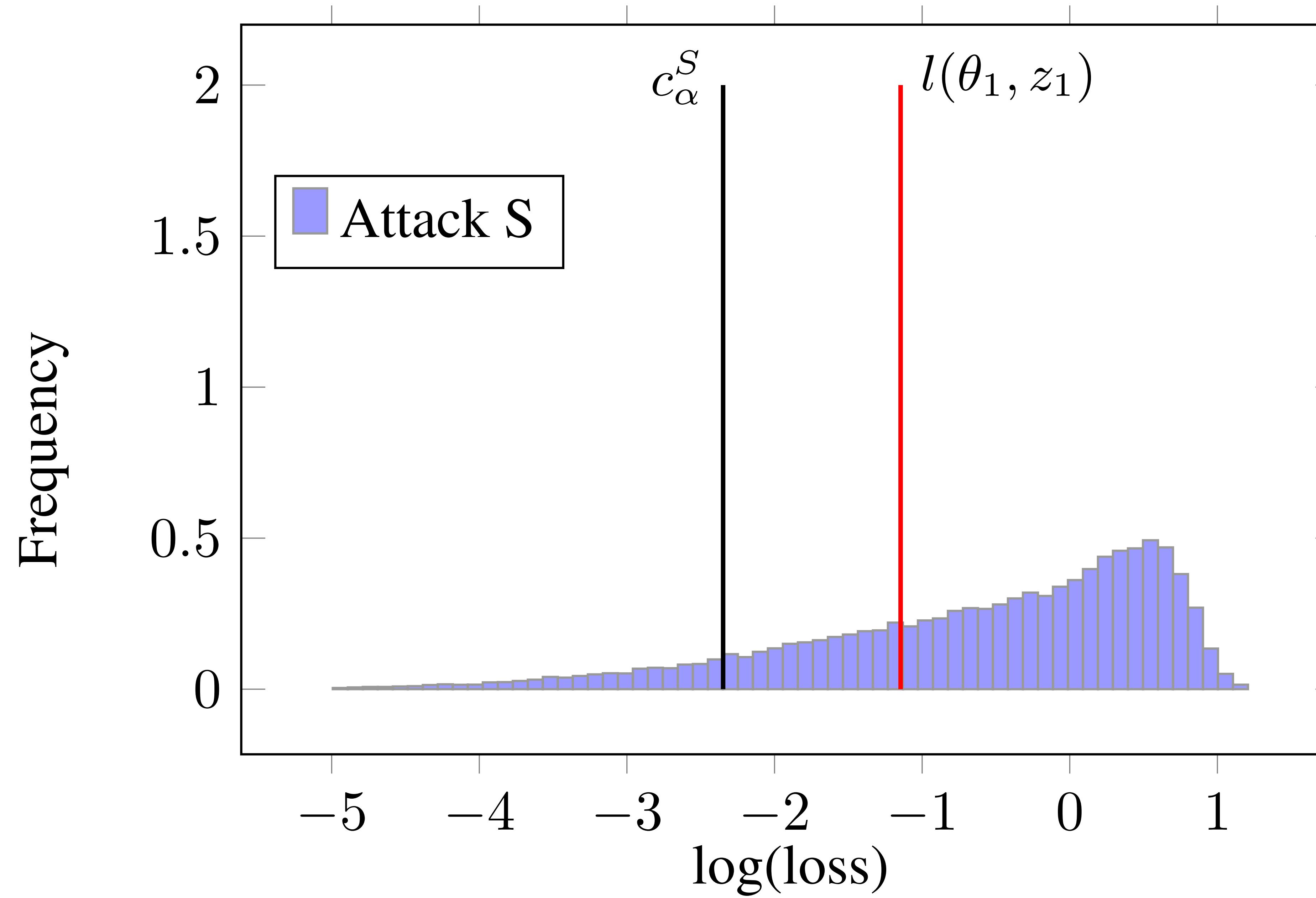
If  $\ell(\theta, x_z, y_z) \leq c_\alpha(\theta, x_z, y_z)$ , reject  $H_0$

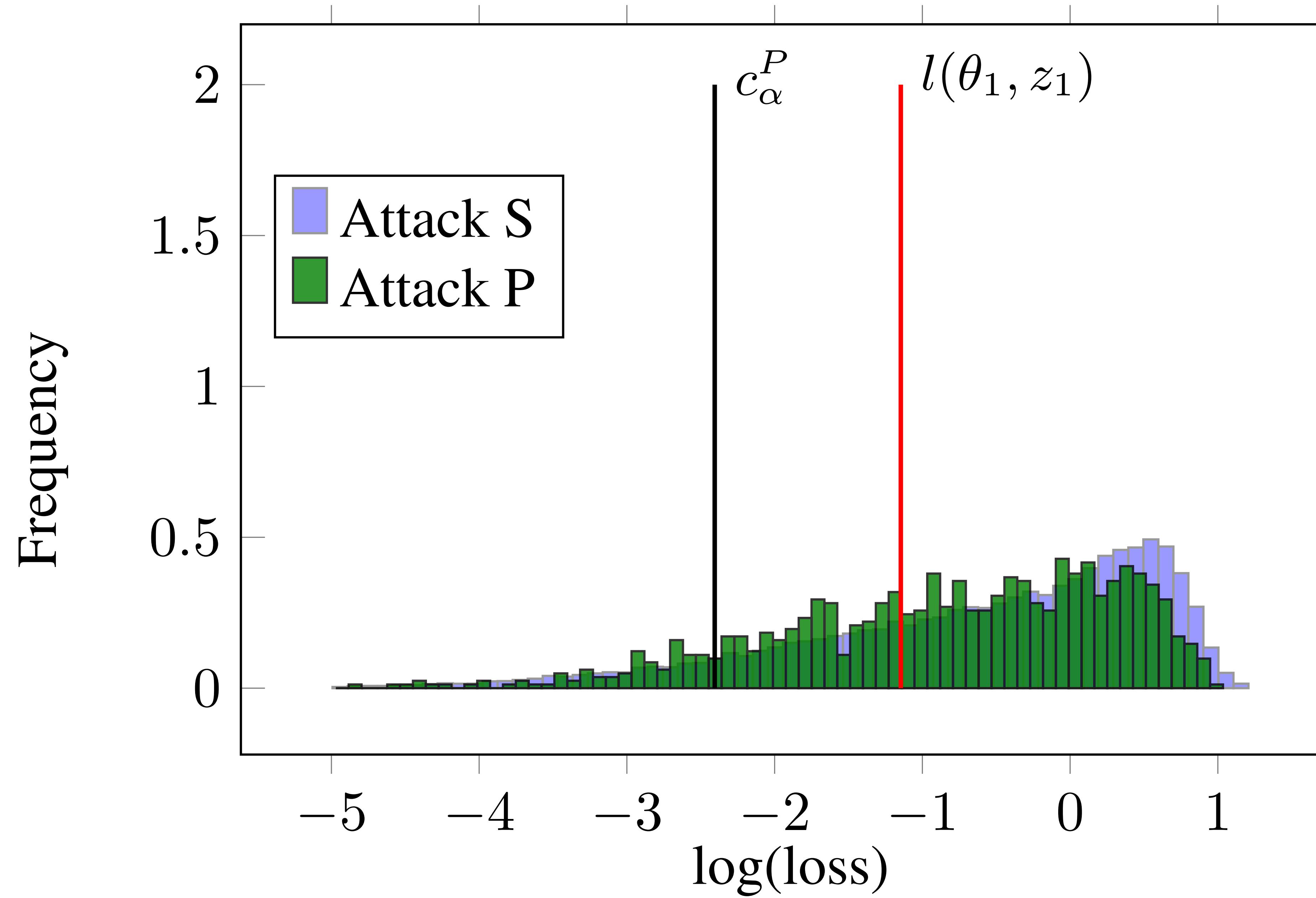
- Learn a threshold from the behaviour of target data on distilled models
- Note that the threshold depends on both target data and the target model

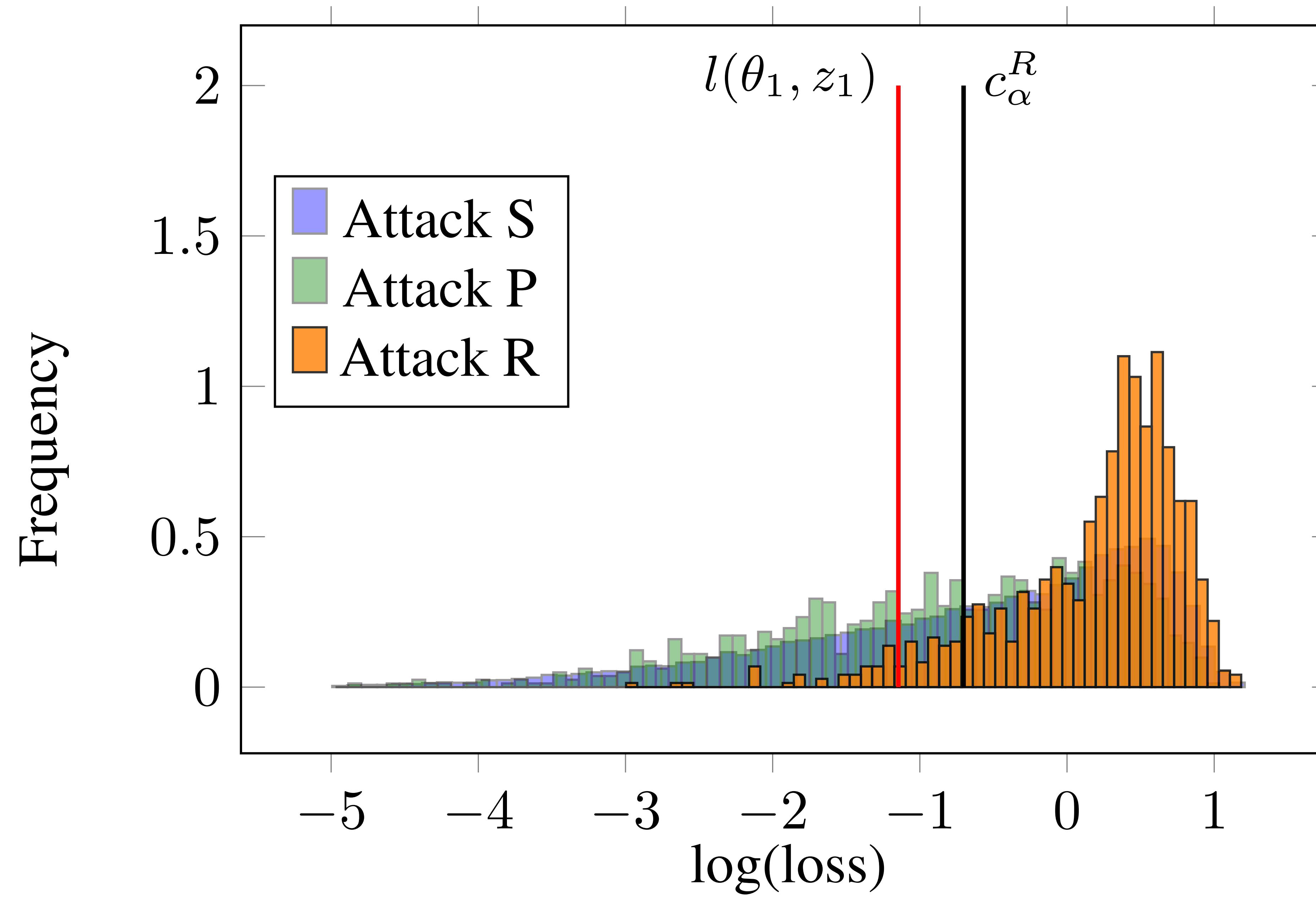


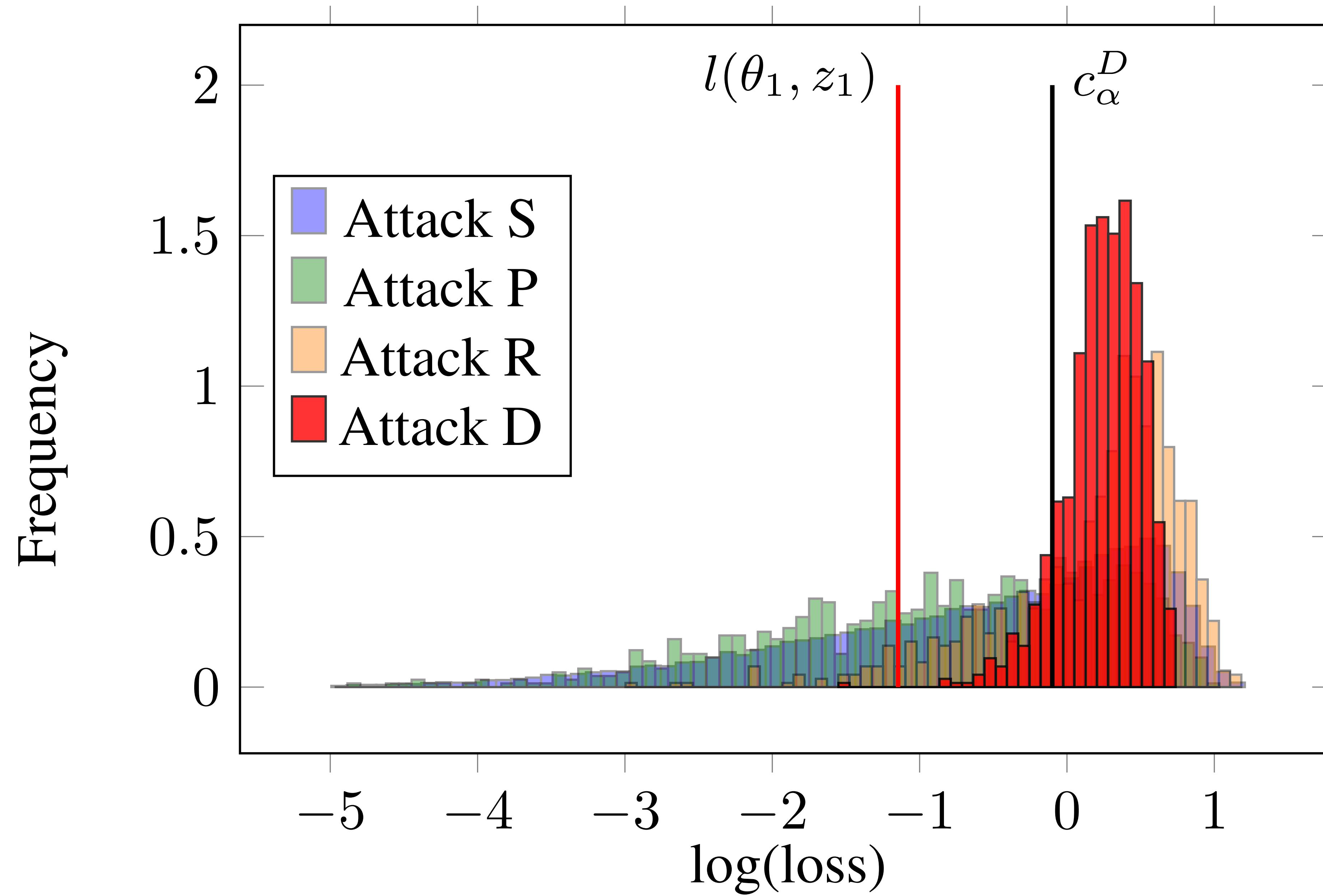
# Agreement between Attacks

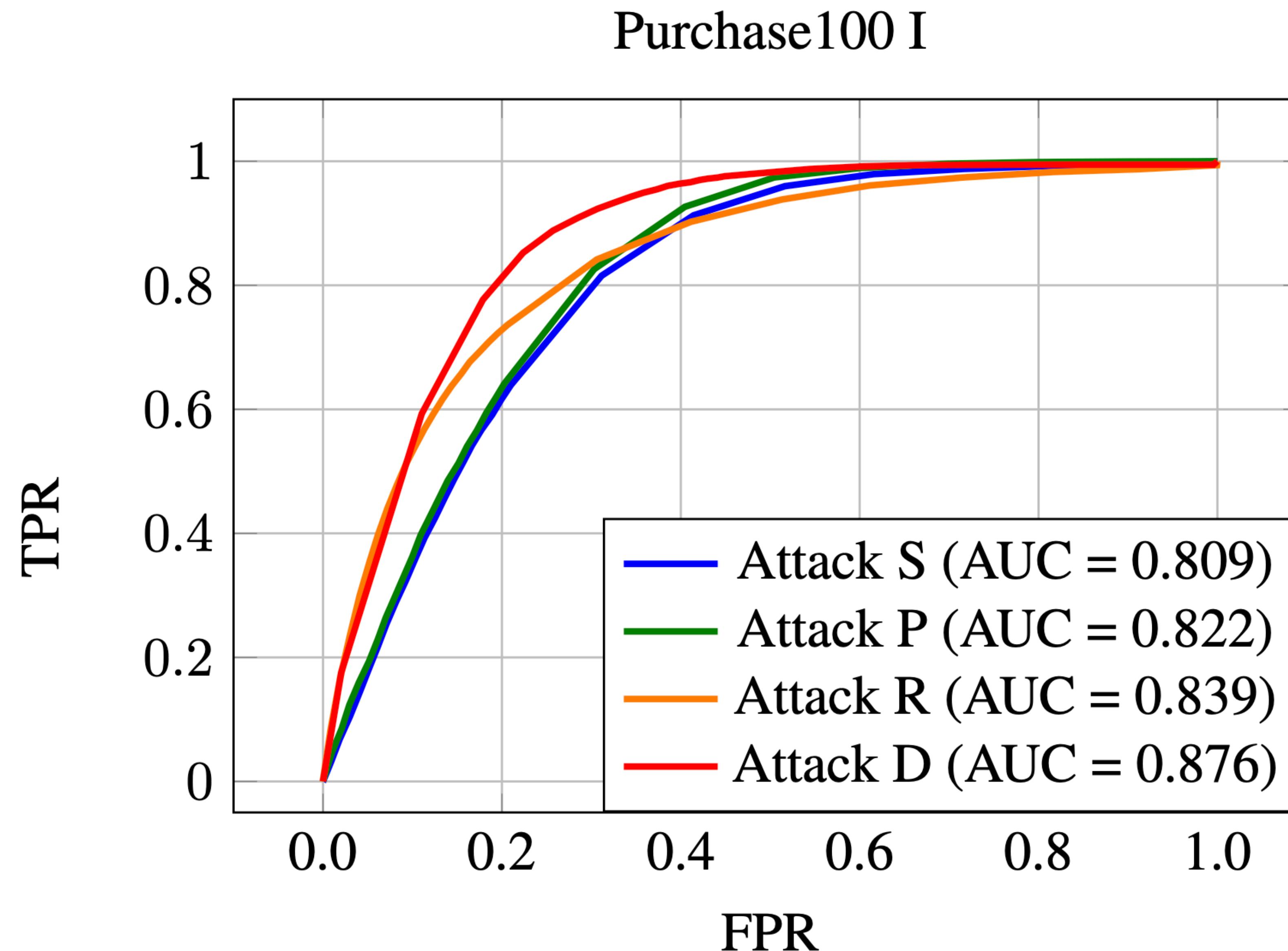












# Auditing Data Privacy for Machine Learning

- Given the privacy vulnerabilities of models, enabling access to models without auditing them (and mitigating the risks) is not much worse than allowing unauthorised access to data
- ML Privacy Meter ([privacy-meter.com](http://privacy-meter.com)) aids regulatory compliance, through a systematic method to audit data privacy for a wide range of machine learning algorithms

[Shokri, Stronati, Song, Shmatikov] Membership Inference Attacks against Machine Learning Models, SP'17

[Nasr, Shokri, Houmansadr] Comprehensive Privacy Analysis of Deep Learning: Passive and Active White-box Inference Attacks against Centralized and Federated Learning, SP'19

[Ye, Maddi, Murakonda, Shokri] Enhanced Membership Inference Attacks against Machine Learning Models, 21

