

Data Privacy in Machine Learning

Reza Shokri

Data Privacy and Trustworthy ML Research Lab
National University of Singapore

 reza@comp.nus.edu.sg  @rzshokri

Privacy Regulations

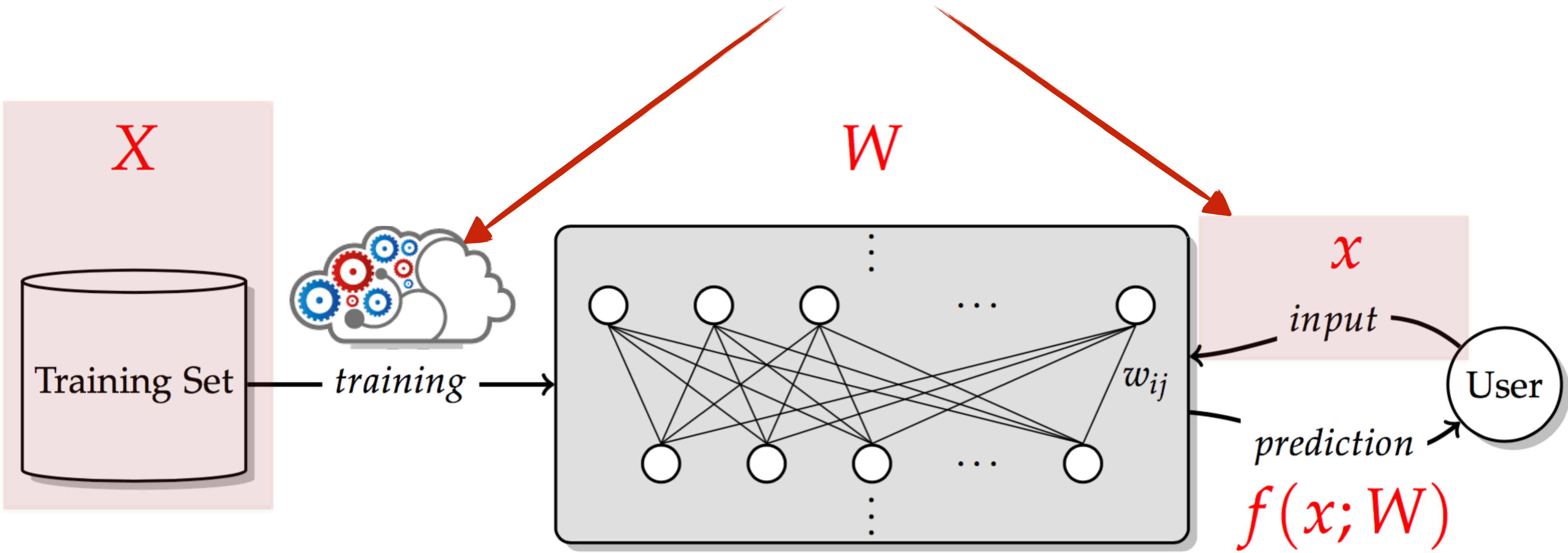


GDPR — Data Protection Impact Assessment

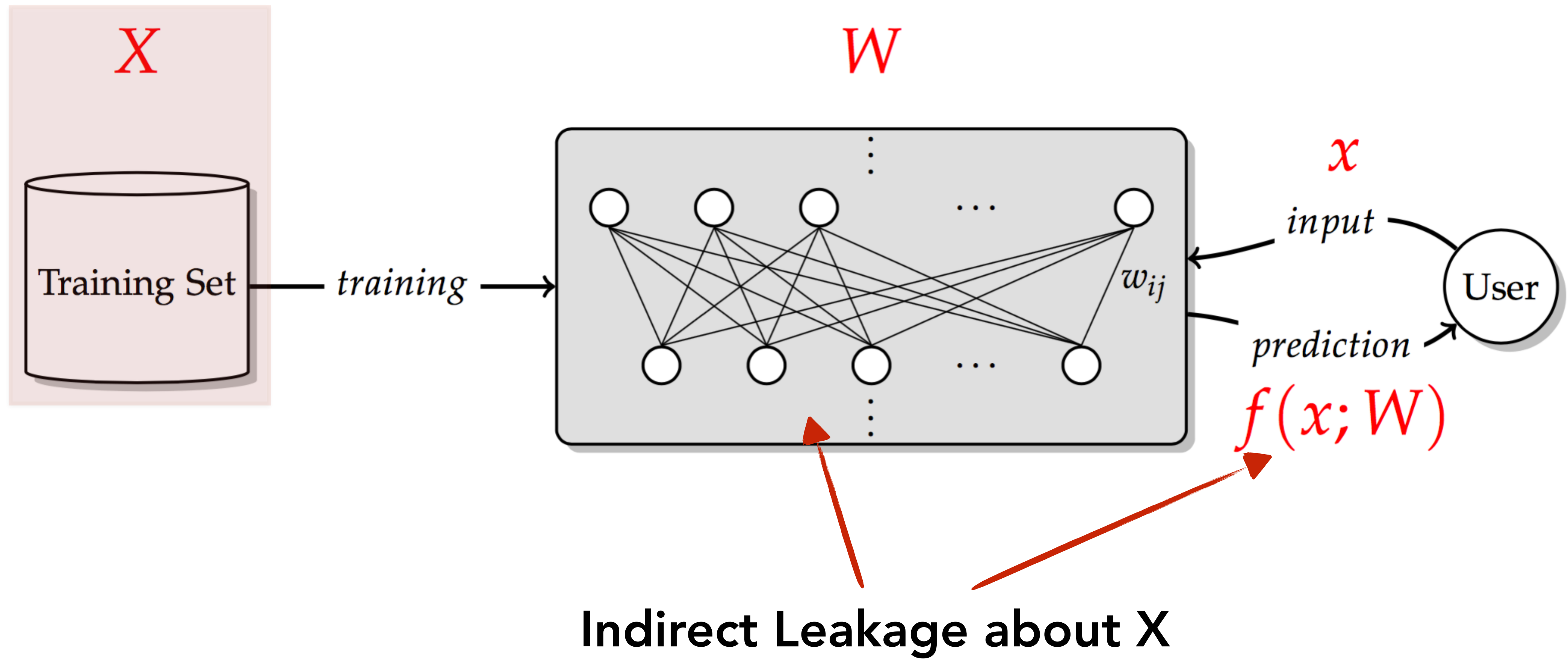
The focus is mostly on data collection, data sharing, access control, ...

Direct Privacy Risks

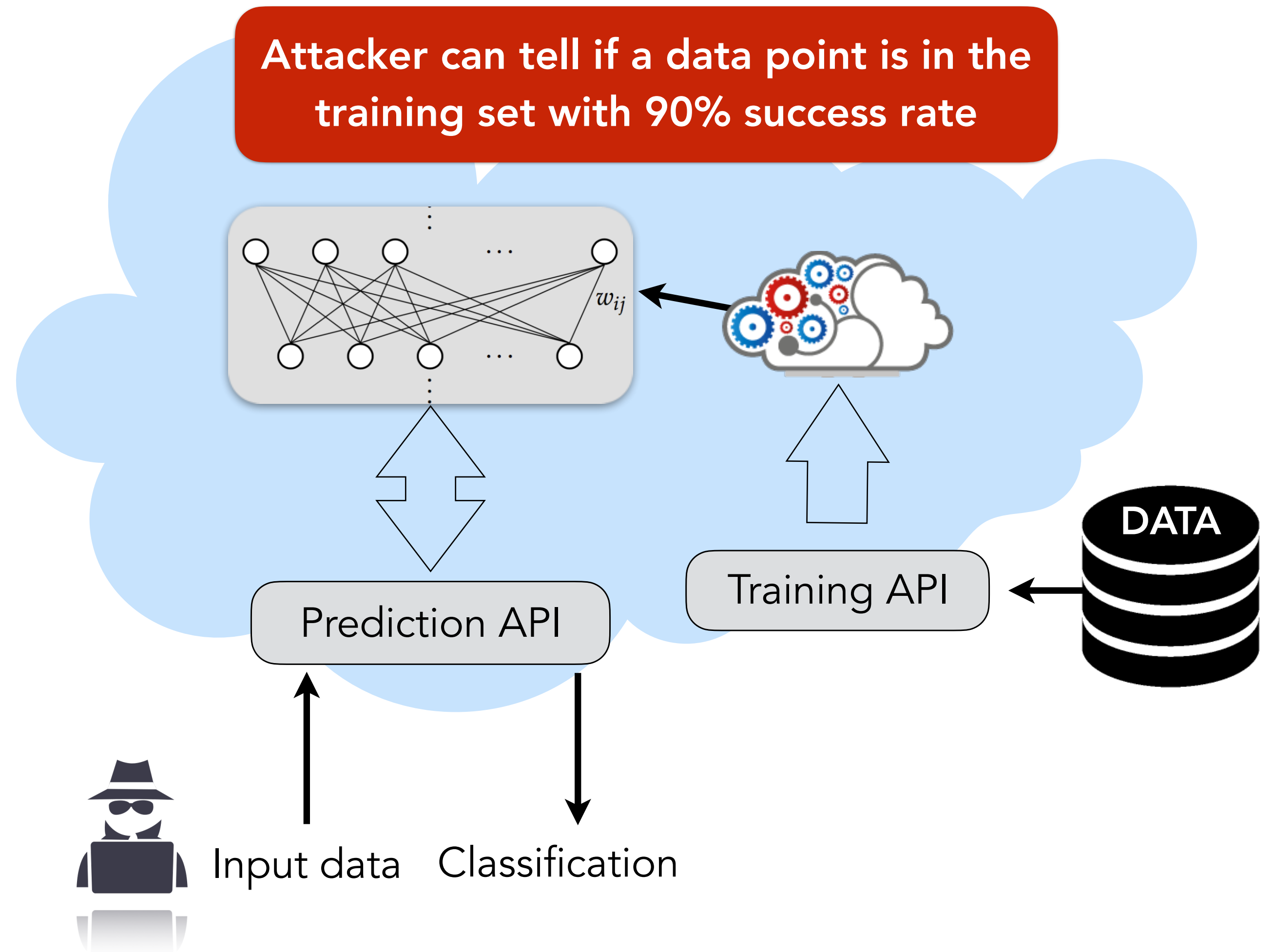
Direct Access to Sensitive Data



Indirect Privacy Risks



Real World Attacks against Machine Learning as a Service Platforms



Real World Attacks against Large Language Models



Attacker can partially reconstruct the sensitive data used for training the model

random input

Prefix
East Stroudsburg Stroudsburg...

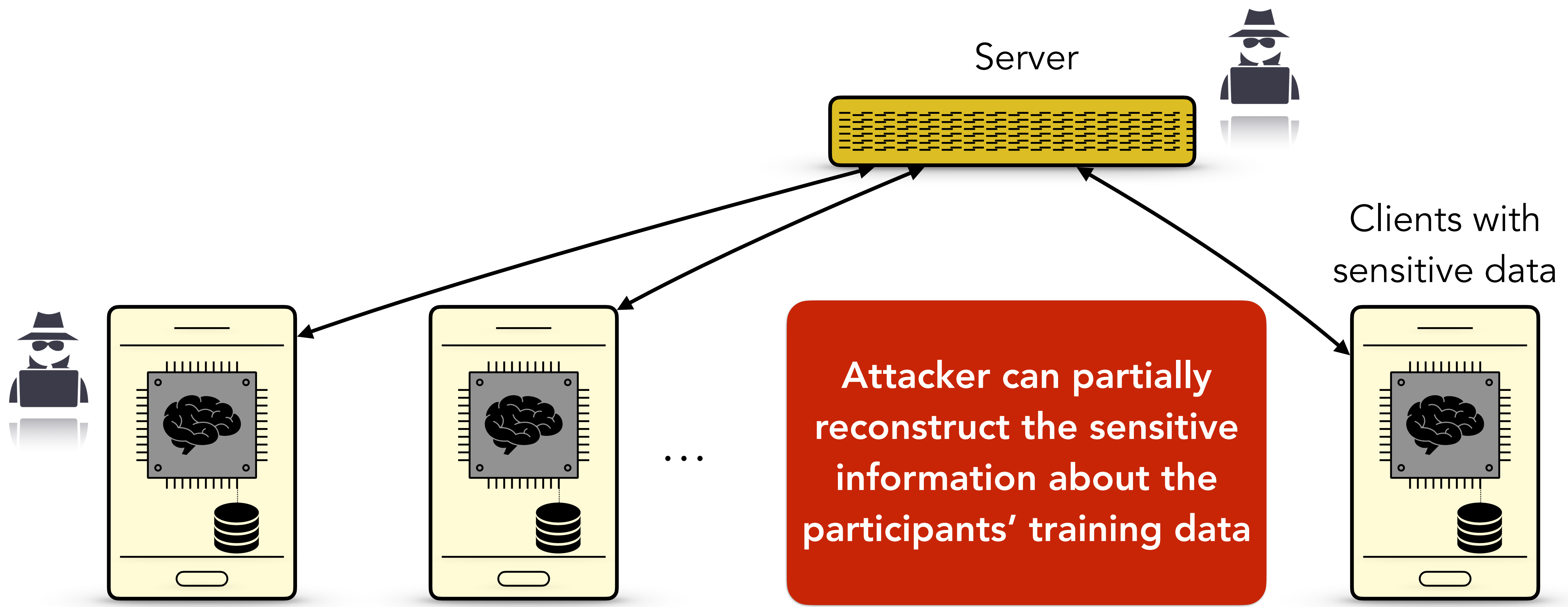
OpenAI's language model trained on text from 8 million web pages

GPT-2

someone's contact information output by the model
(redacted for privacy)

Memorized text
[redacted] Corporation Seabank Centre
[redacted] Marine Parade Southport
Peter W [redacted]
[redacted]@ [redacted].com
+ [redacted] 7 5 [redacted] 40 [redacted]
Fax: + [redacted] 7 5 [redacted] 0 [redacted]

Real World Attacks against Federated Learning Algorithms



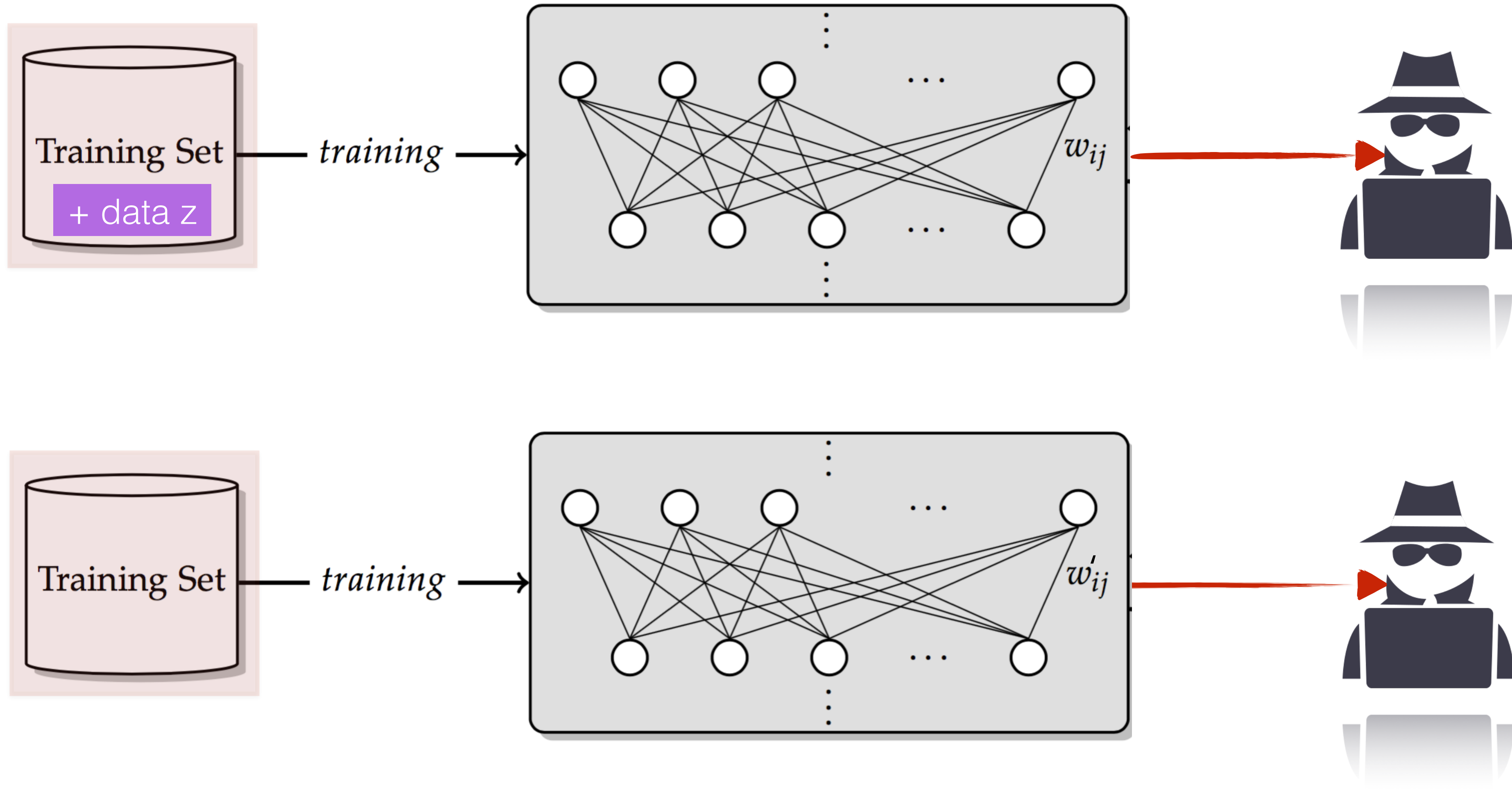
[Nasr, Shokri, Houmansadr] Comprehensive Privacy Analysis of Deep Learning: Passive and Active White-box Inference Attacks against Centralized and Federated Learning, SP'19

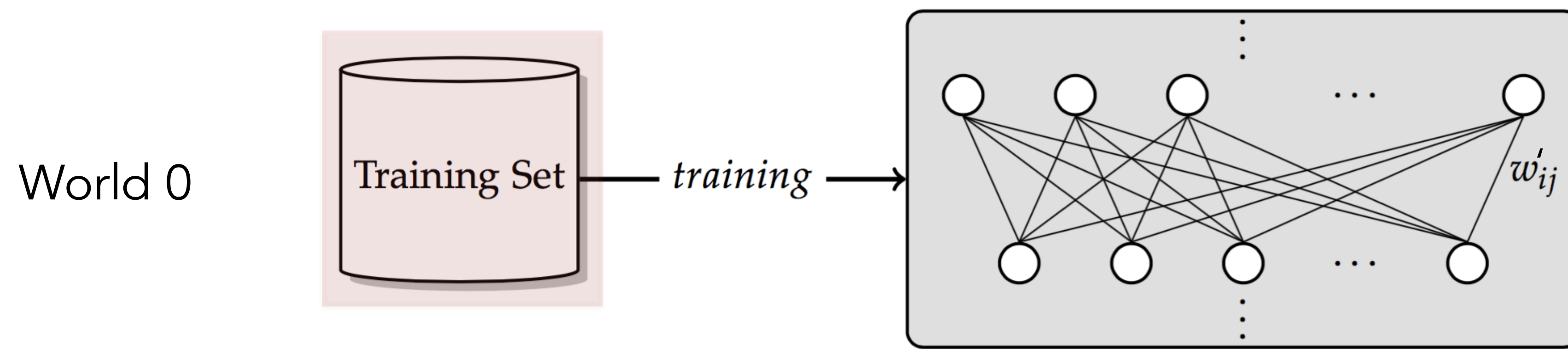
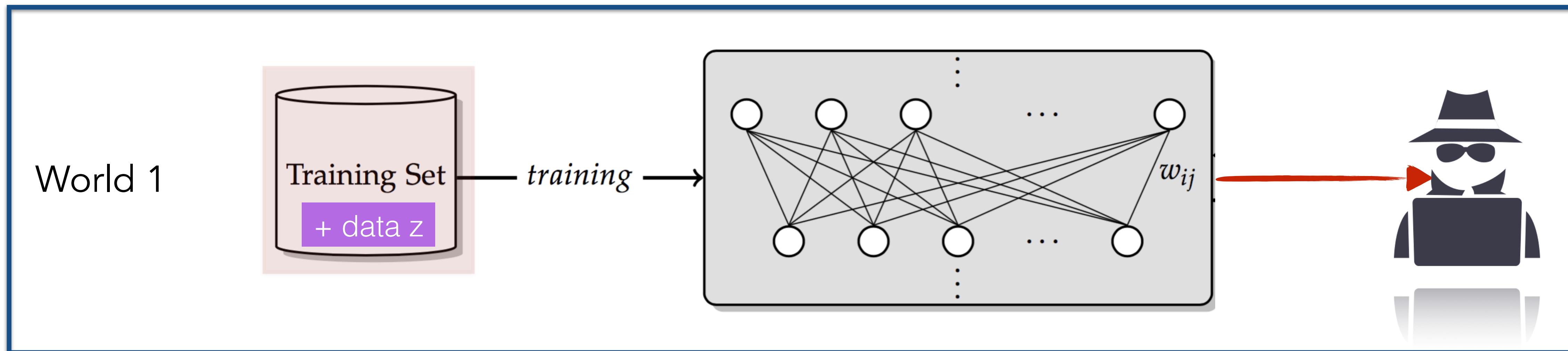
[Melis, Song, De Cristofaro, Shmatikov] Exploiting Unintended Feature Leakage in Collaborative Learning, SP'19

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

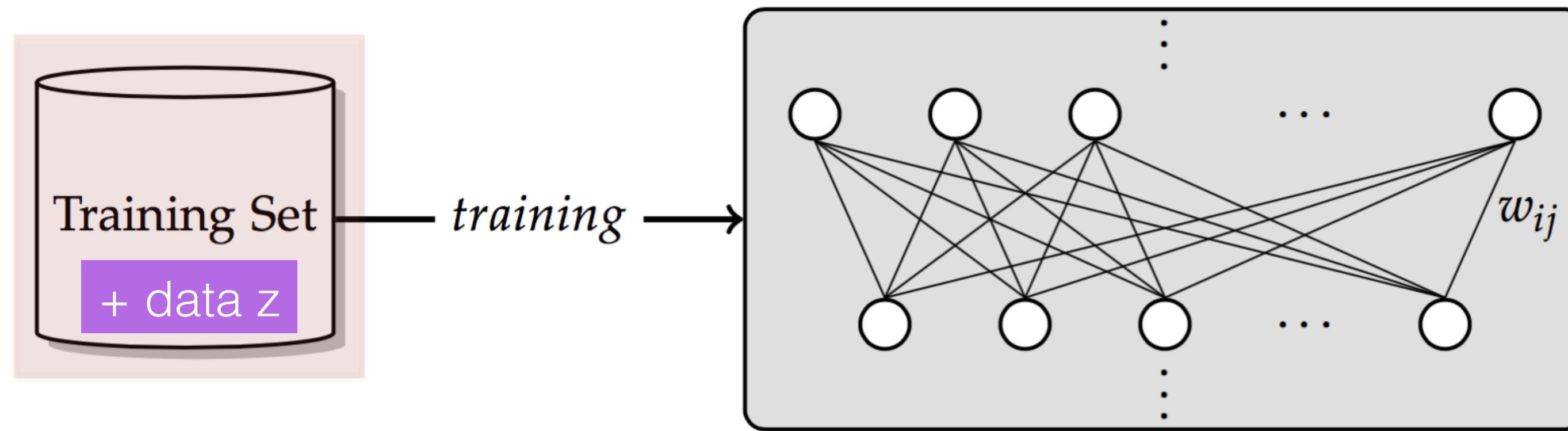
Models are **personal data**

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

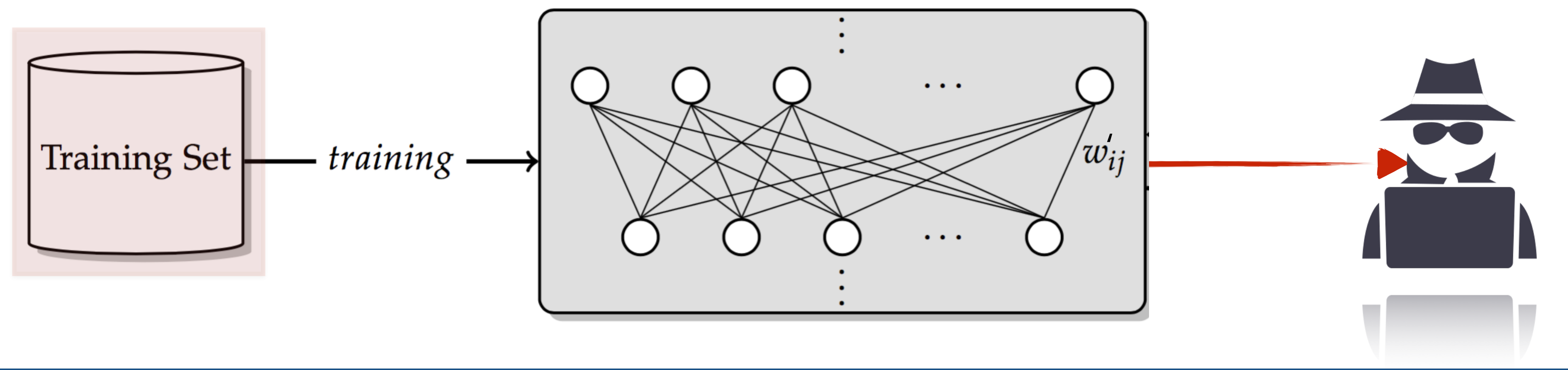




World 1

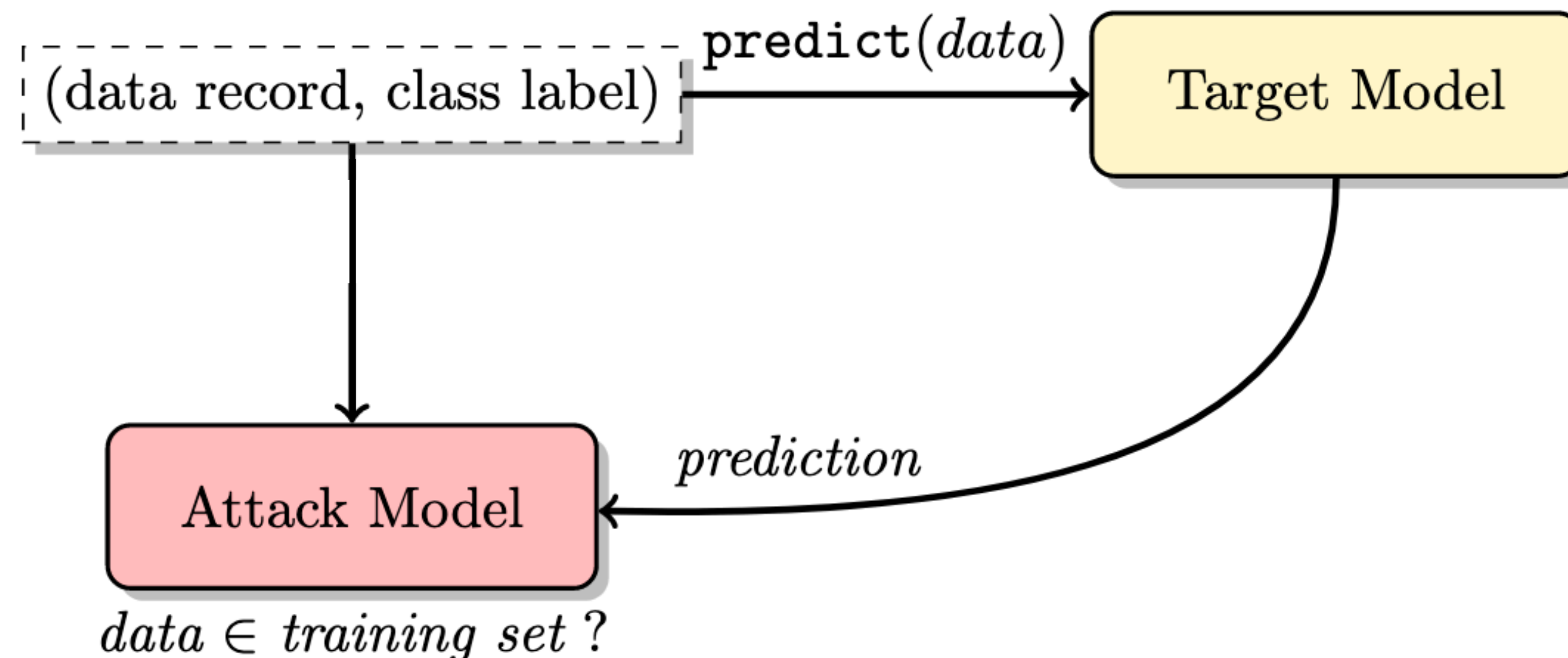


World 0



Membership Inference Attacks

- 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



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**"

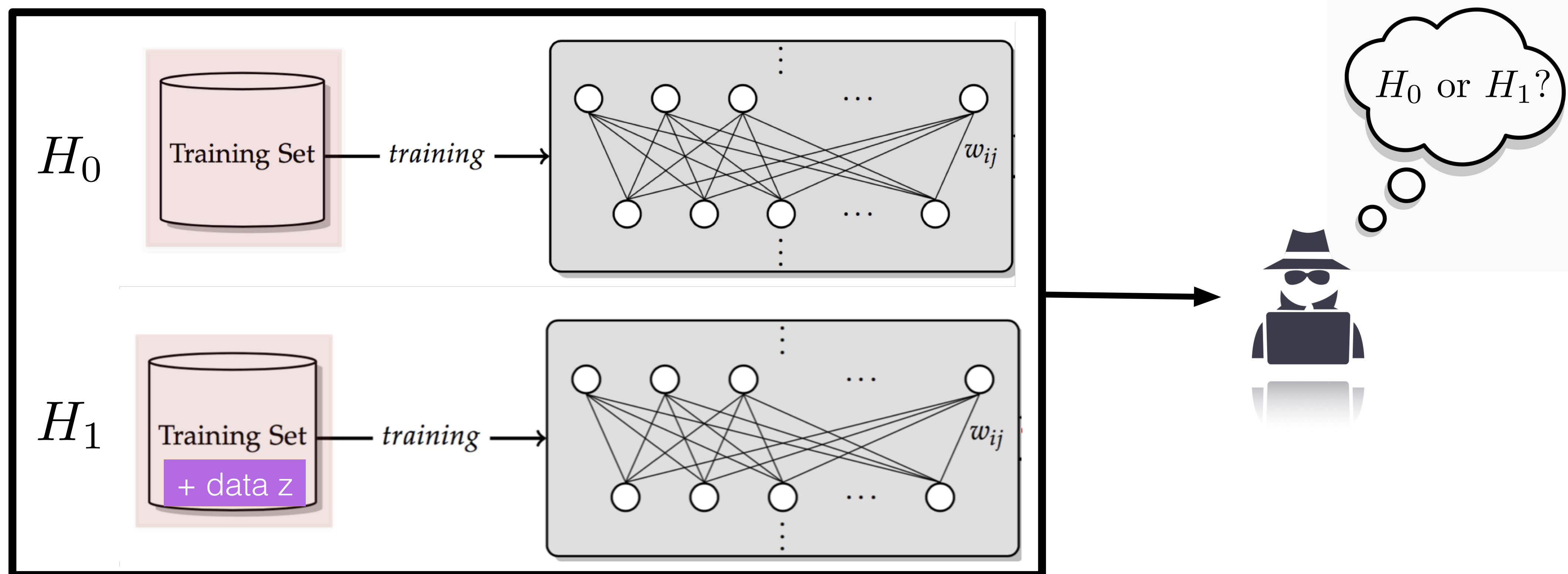
NIST
National Institute of
Standards and Technology
U.S. Department of Commerce

ico.
Information Commissioner's Office

FROM: Russell T. Vought
Director

SUBJECT: Guidance for Regulation of Artificial Intelligence Applications

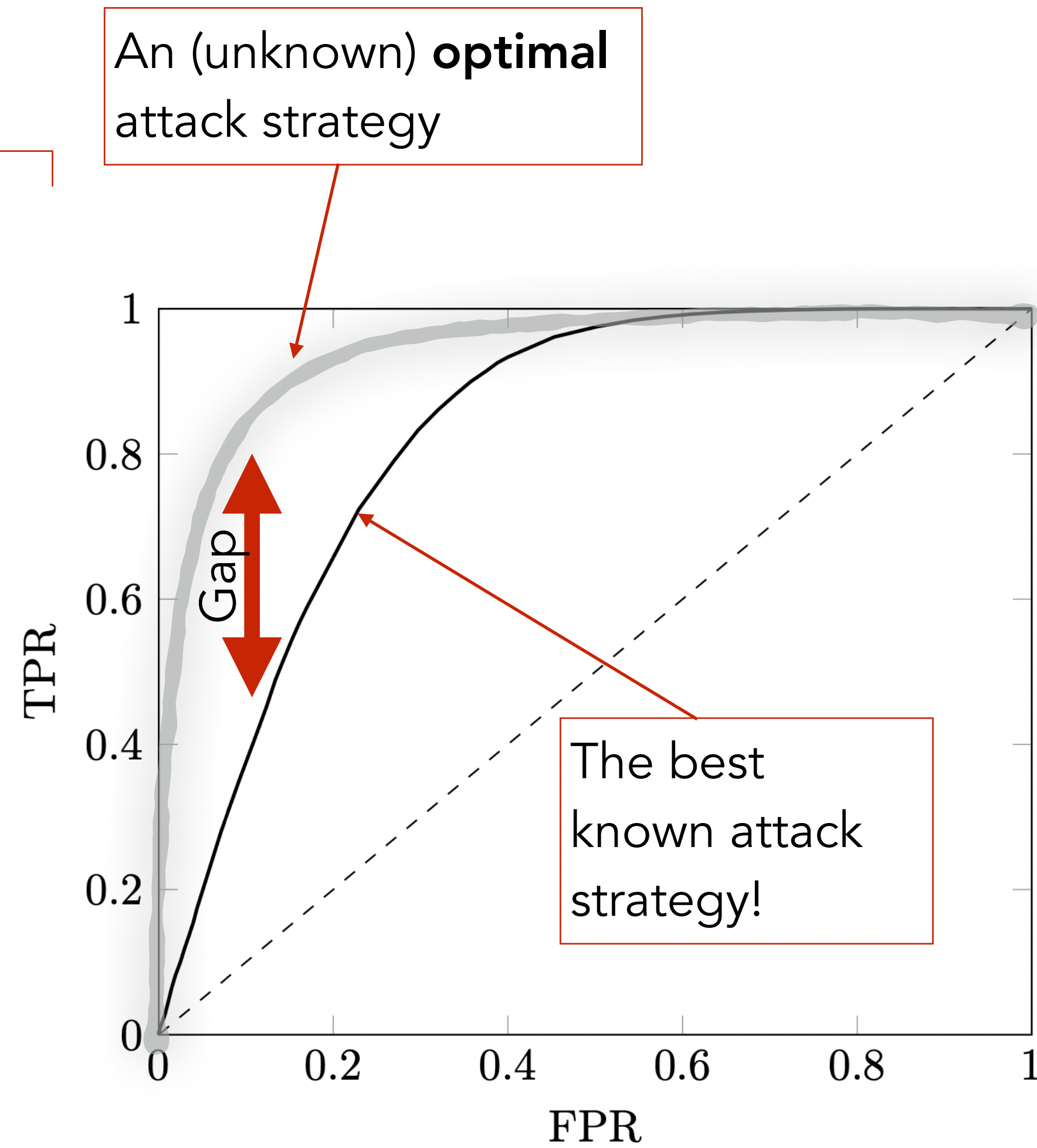
Membership Inference Attack (MIA) Game



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

power: members are correctly predicted as member

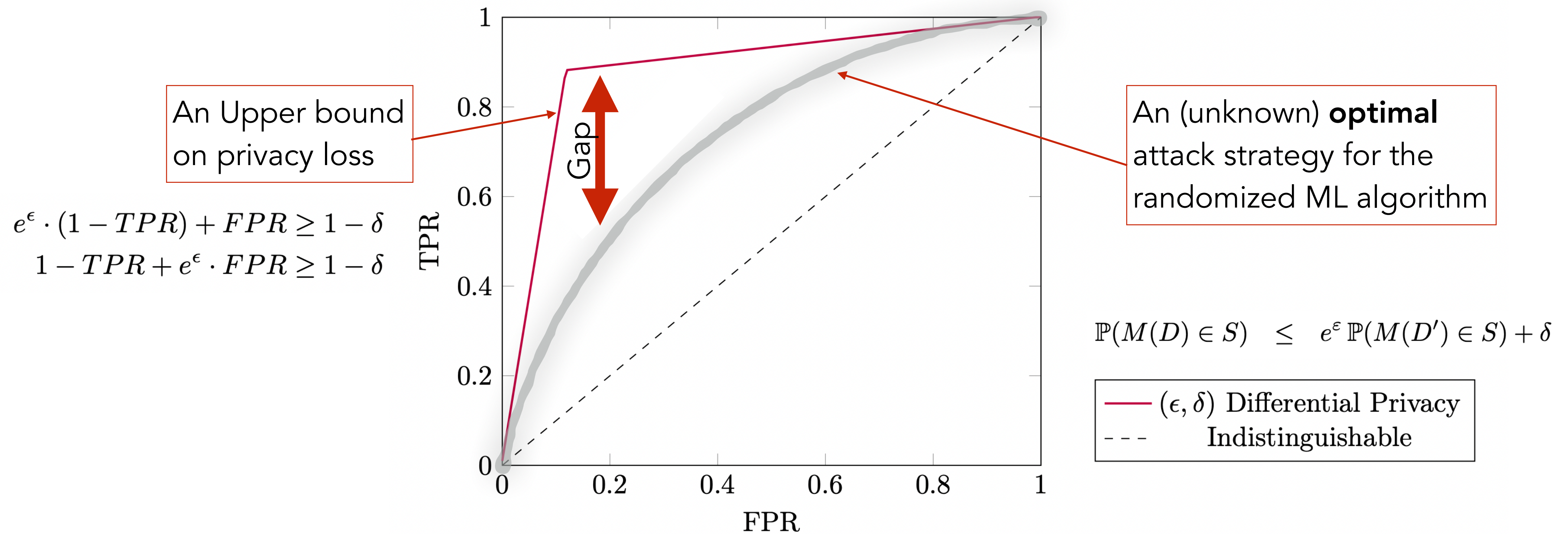
Ev
be
att
str



error: non-members are wrongly predicted as member

An attack strategy gives a **lower-bound** on the privacy risk of the target ML algorithm

This is very useful to rule out vulnerable algorithms, ...
but, lack of a known powerful attack is not a guarantee for privacy!



Prove an **upper-bound** for the privacy risk of a randomized algorithm...

A differential privacy guarantee is an **upper-bound** on the privacy risk of a randomized ML algorithm

If the bound is loose, we are over-estimating the risk, thus we unnecessarily over-randomize the algorithm, ... which could result in a high utility drop (e.g., prediction error) due to the algorithm.

How to Design Powerful Auditing Algorithms?

[Ye, Maddi, Murakonda, Bindschaedler, Shokri]

Enhanced Membership Inference Attacks against Machine Learning Models

ACM CCS'22

Hypothesis Testing for Membership Inference

- Given a data point " z " and black-box access to a model " θ ",
 - Determine if " z " was a member of the training set of " θ "

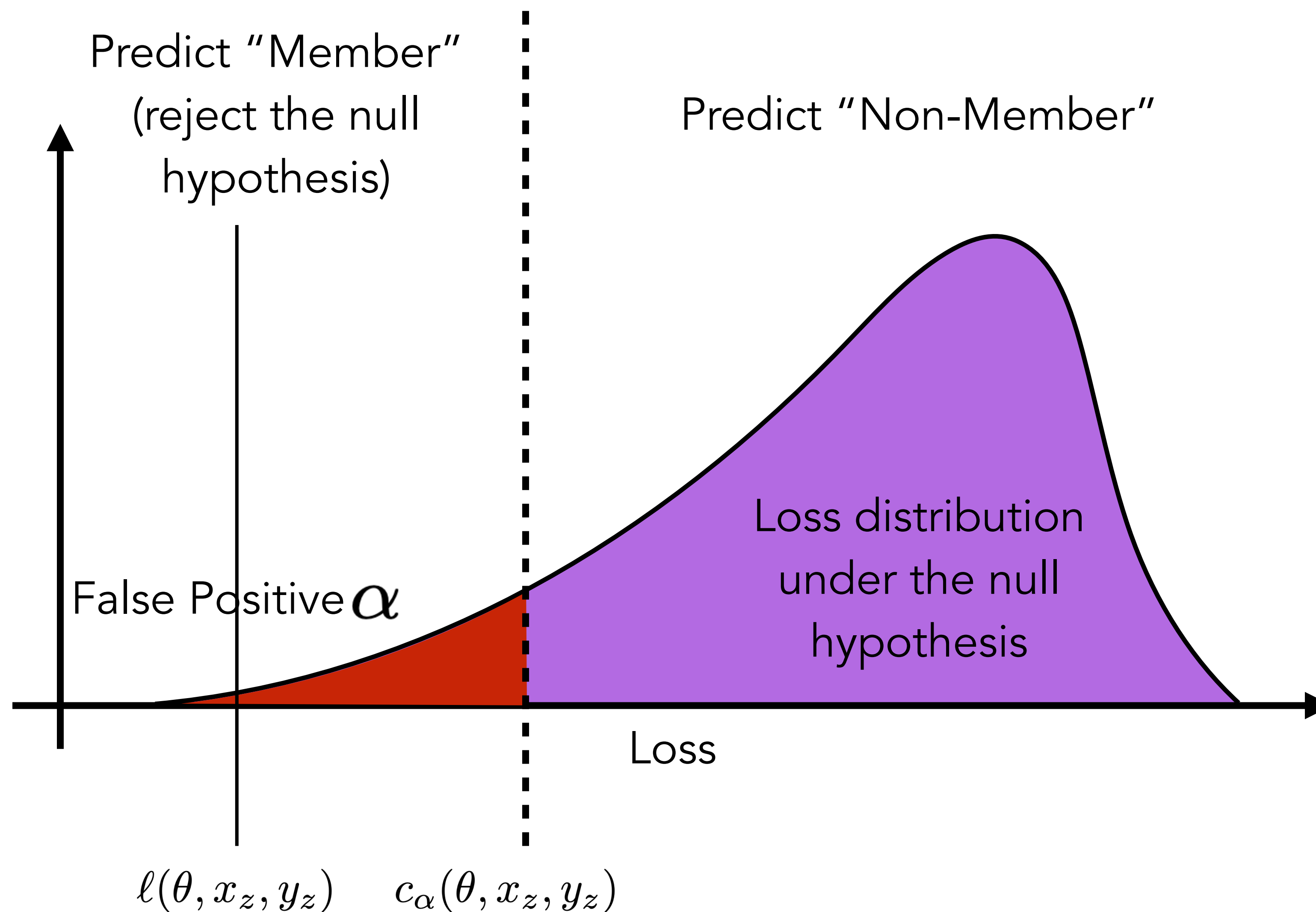
Likelihood Ratio Test

$$LR(\theta, z) = \frac{L(H_0|\theta, z)}{L(H_1|\theta, z)}$$

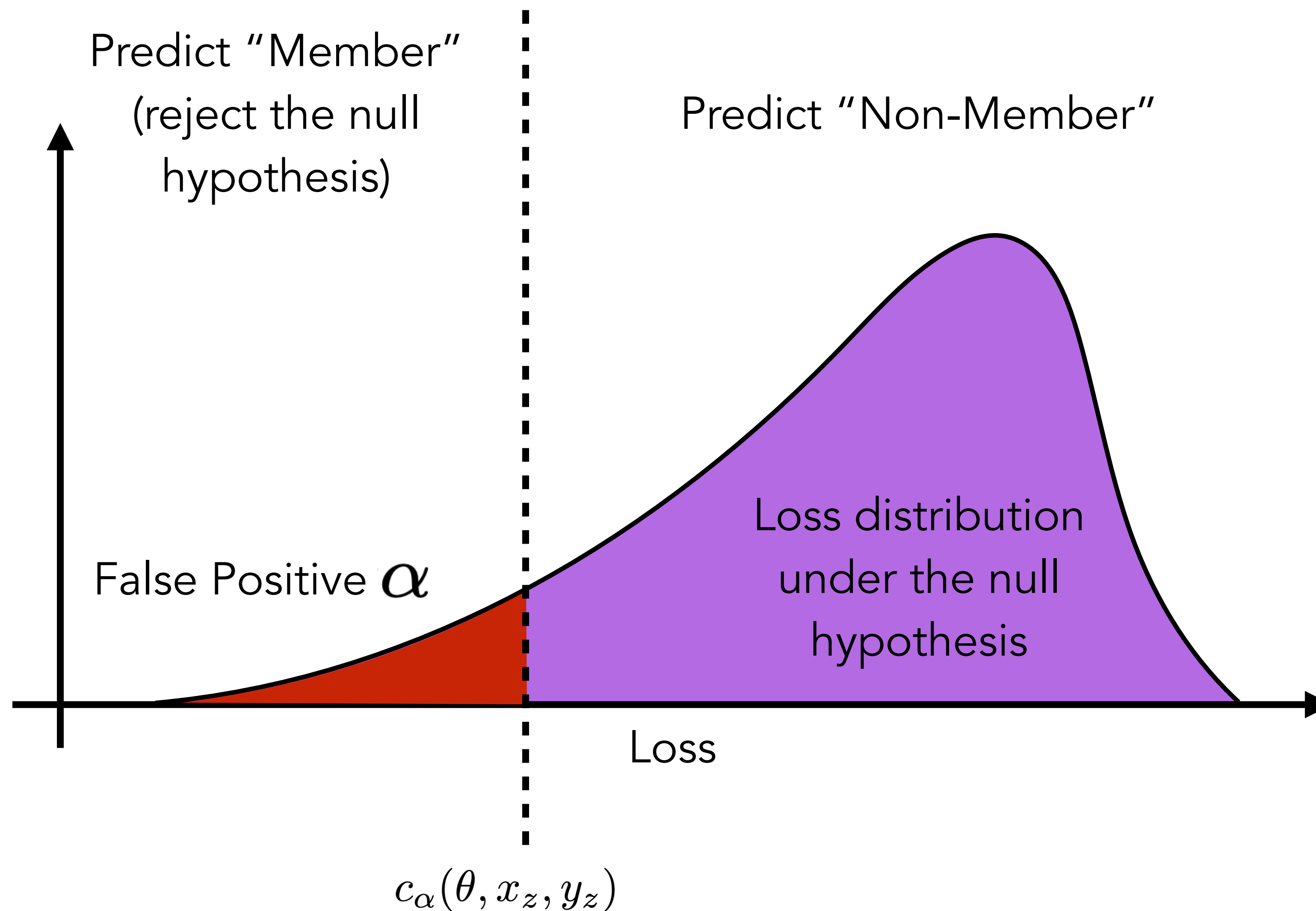
Attack: **If $\ell(\theta, x_z, y_z) \leq c_\alpha(\theta, x_z, y_z)$, reject H_0**

↑
false positive rate

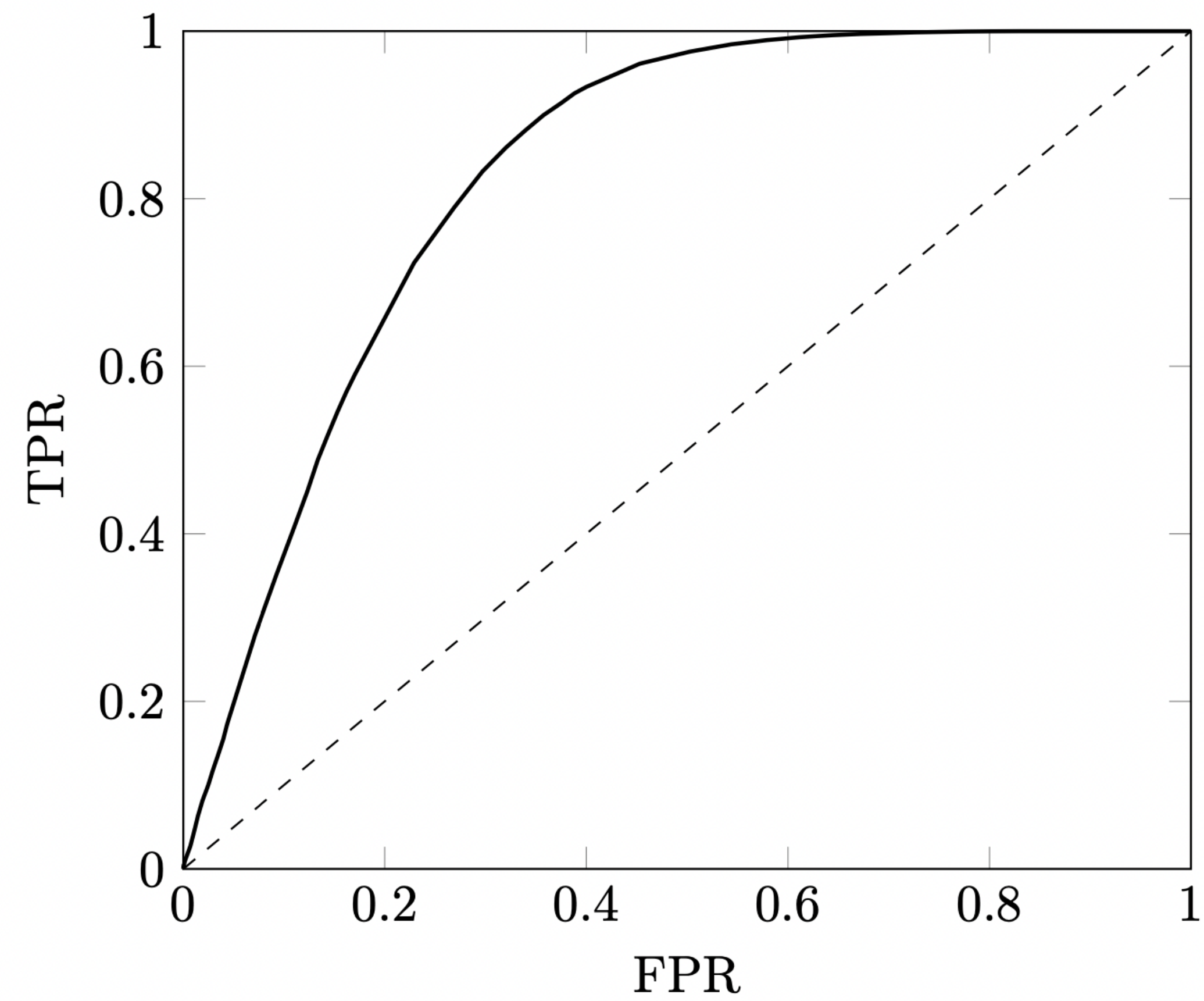
How to Interpret the Test?



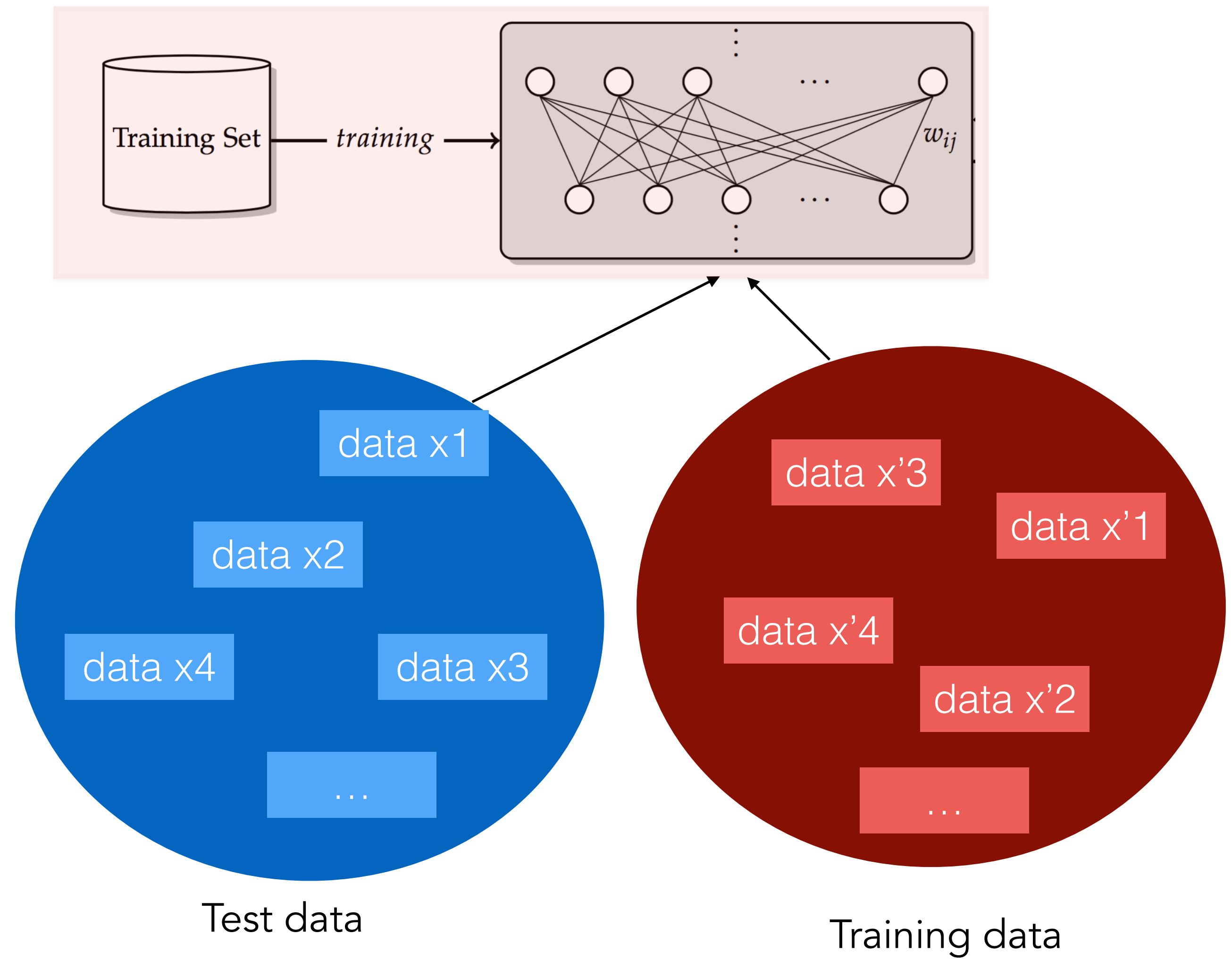
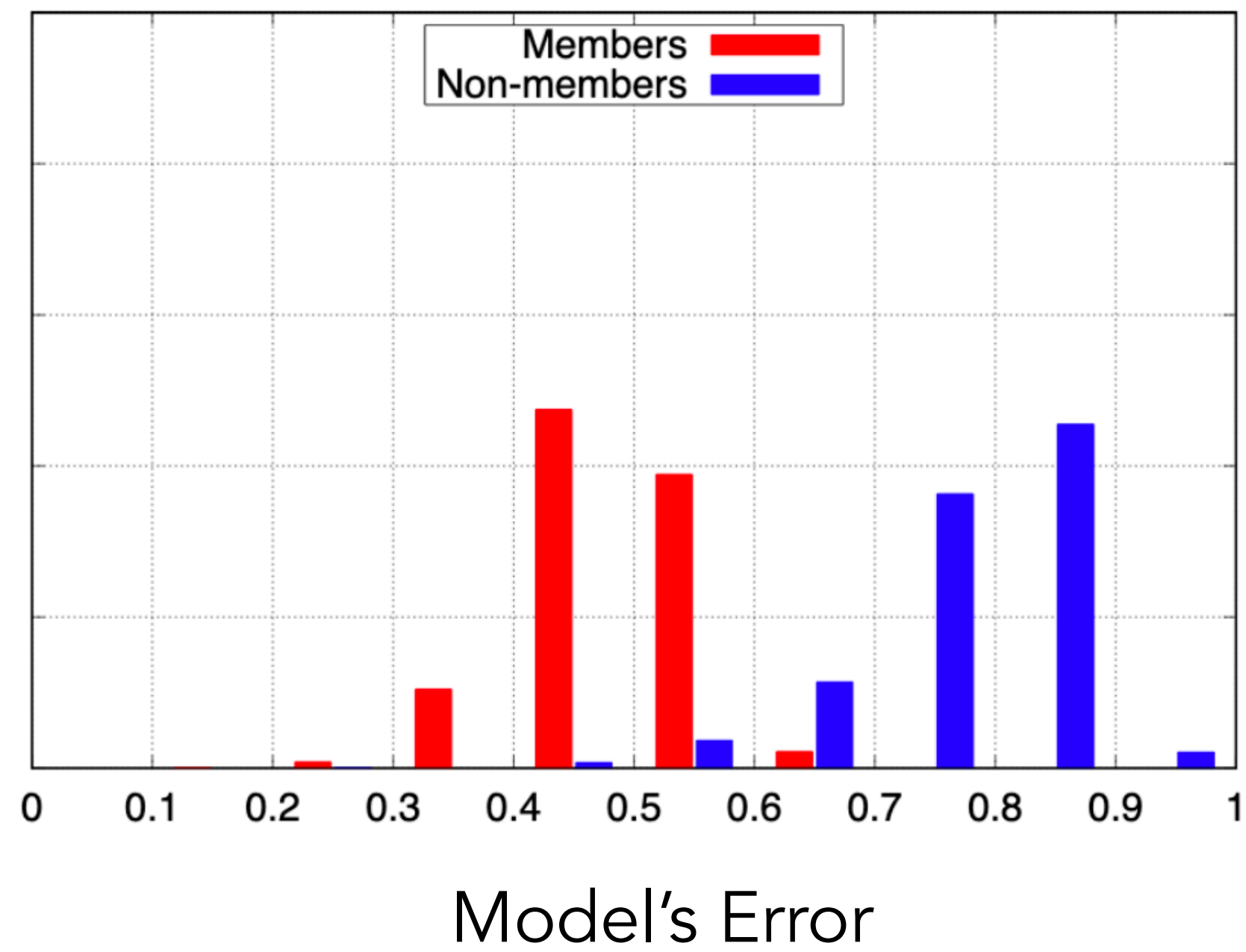
How to Construct the Test?



Power vs Error of the Test



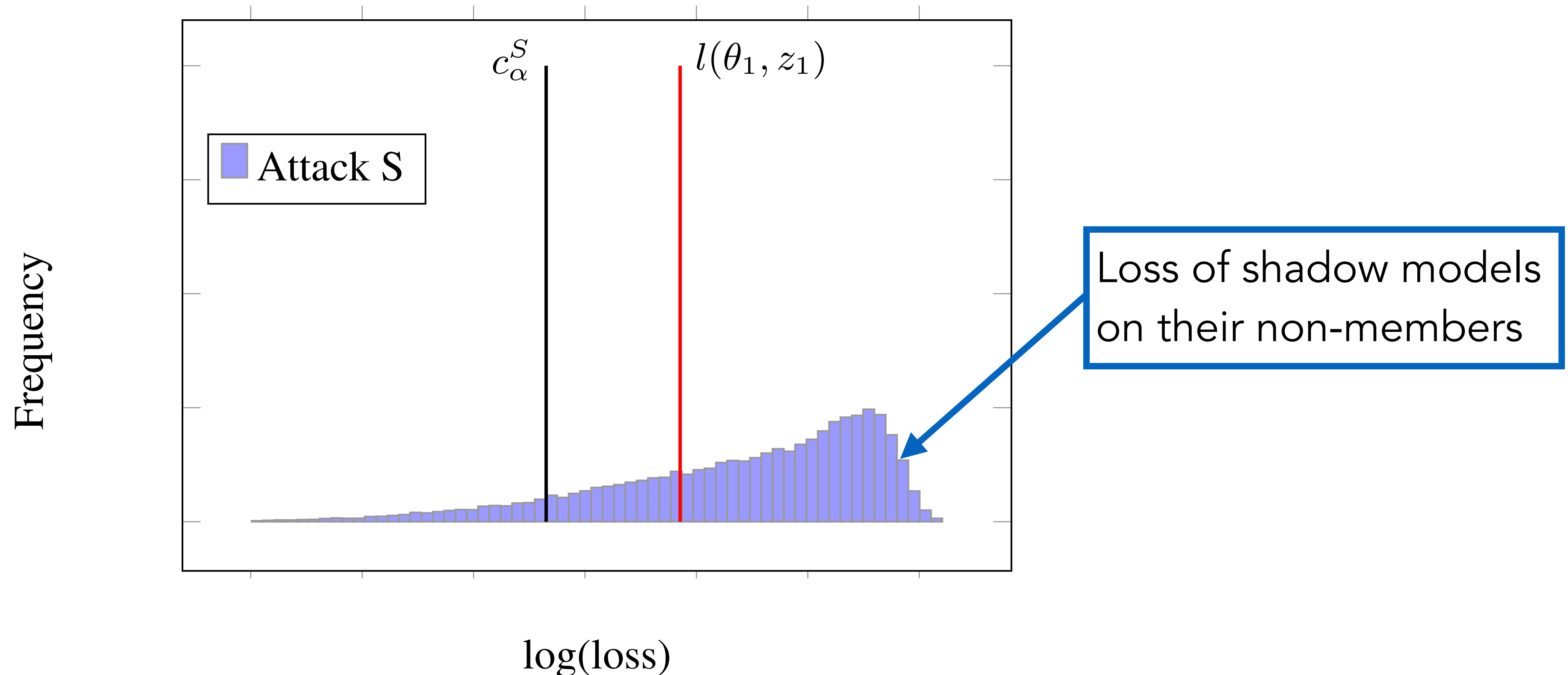
Constructing the Test ...

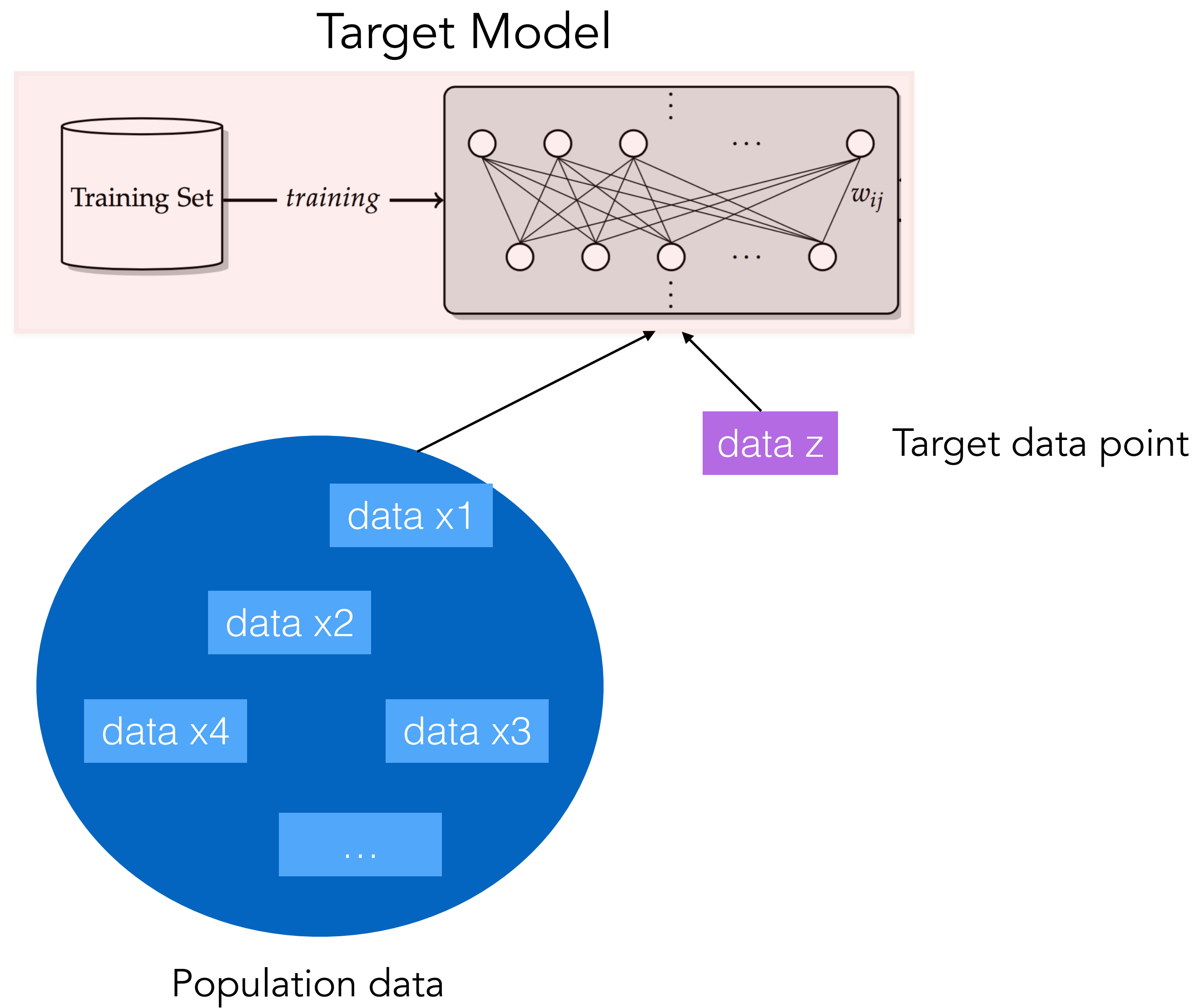


Membership Inference via Shadow Models

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 [SSSS2017]
- Learn a threshold from the behavior of some shadow models on their test data

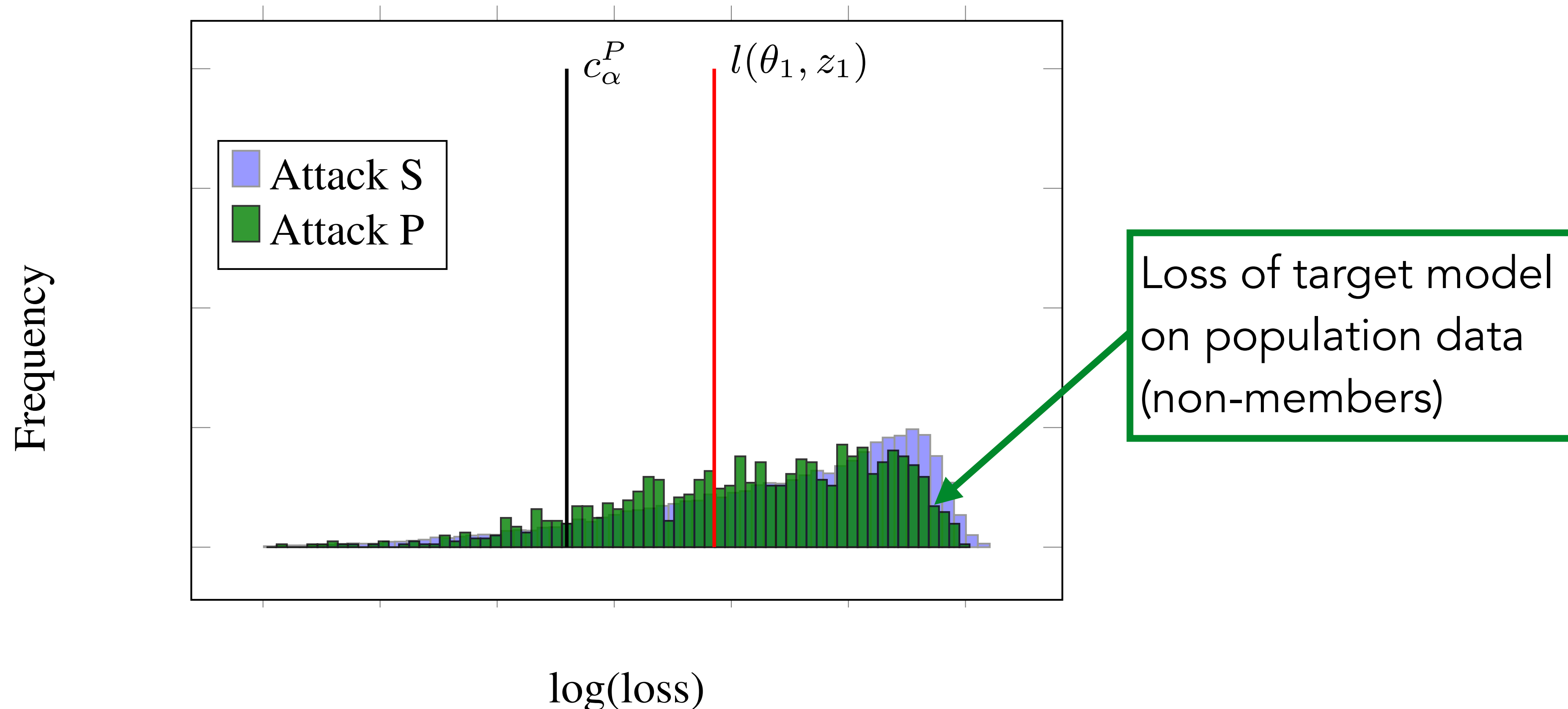




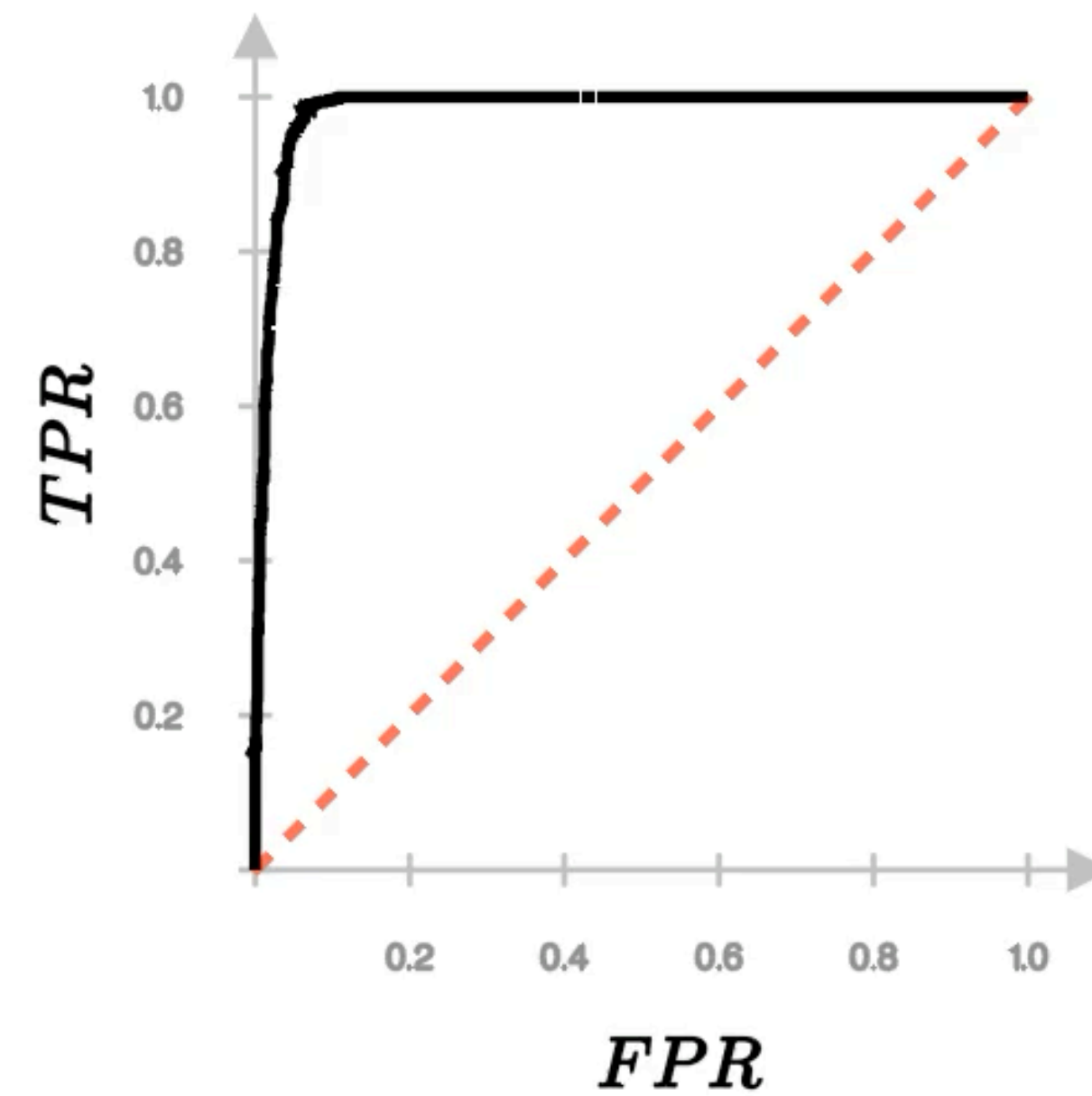
Membership Inference via Population Data

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

- Directly learn a threshold from the loss distribution of the target model on population data

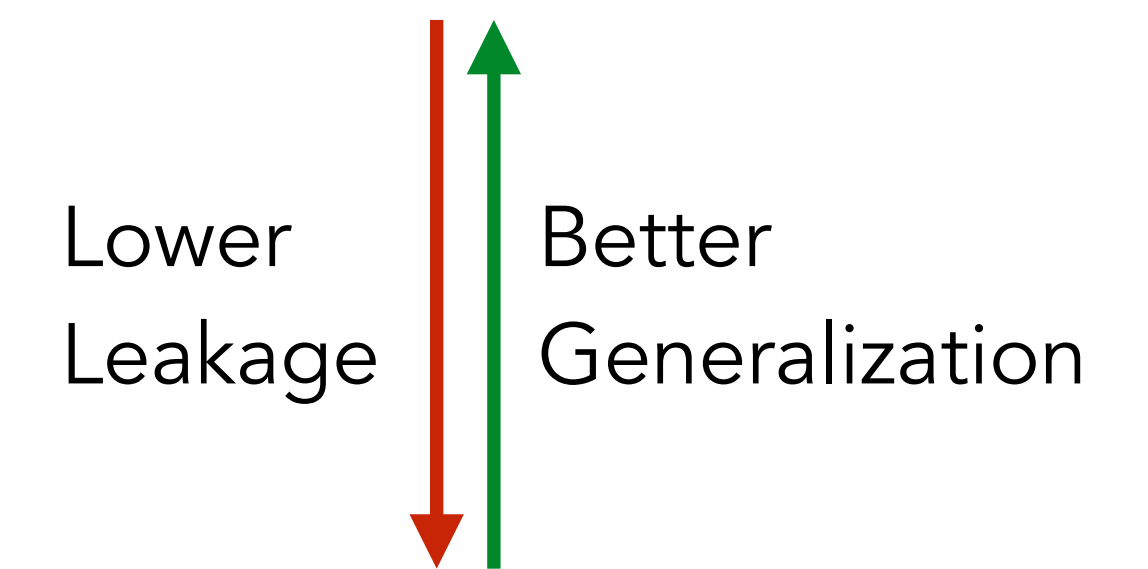


Reason for Leakage?

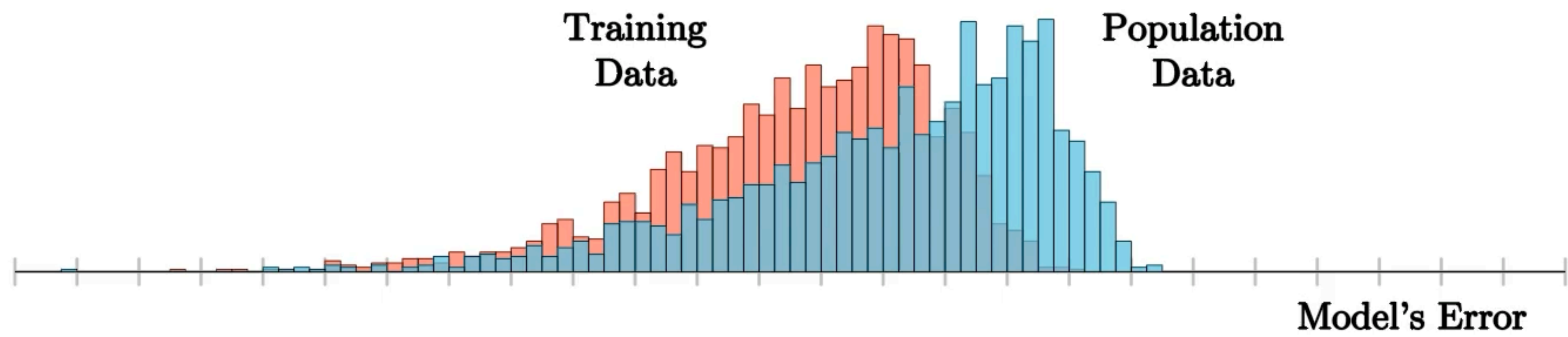
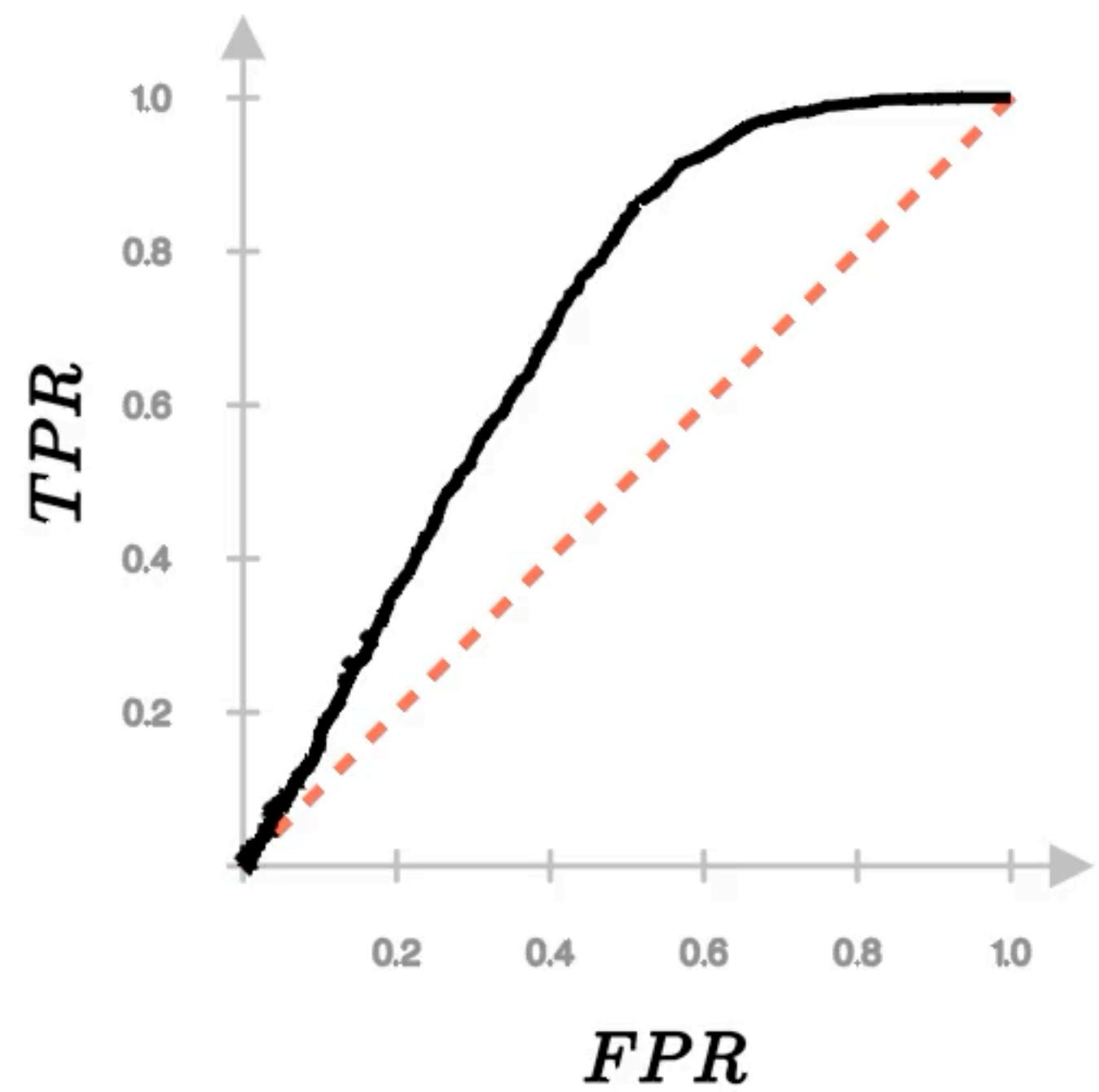


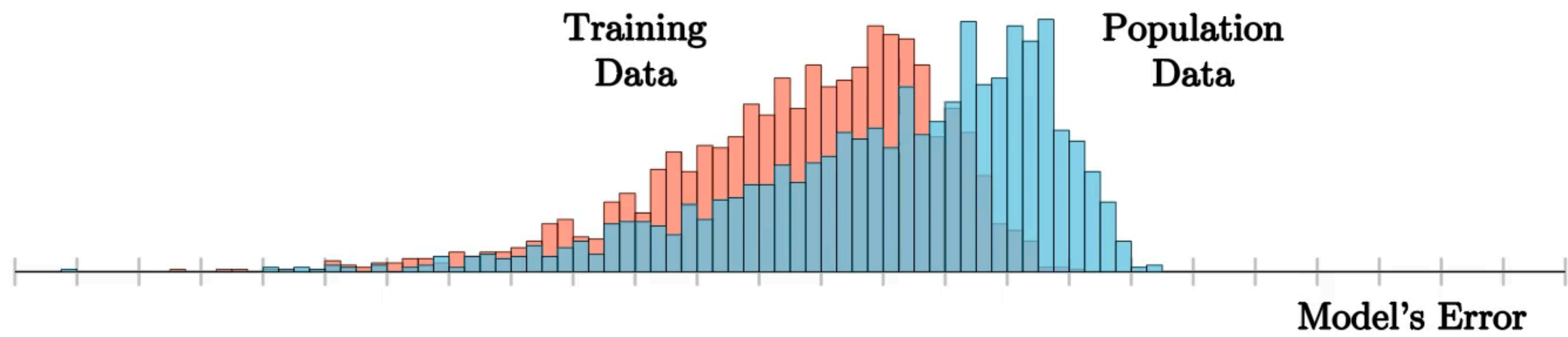
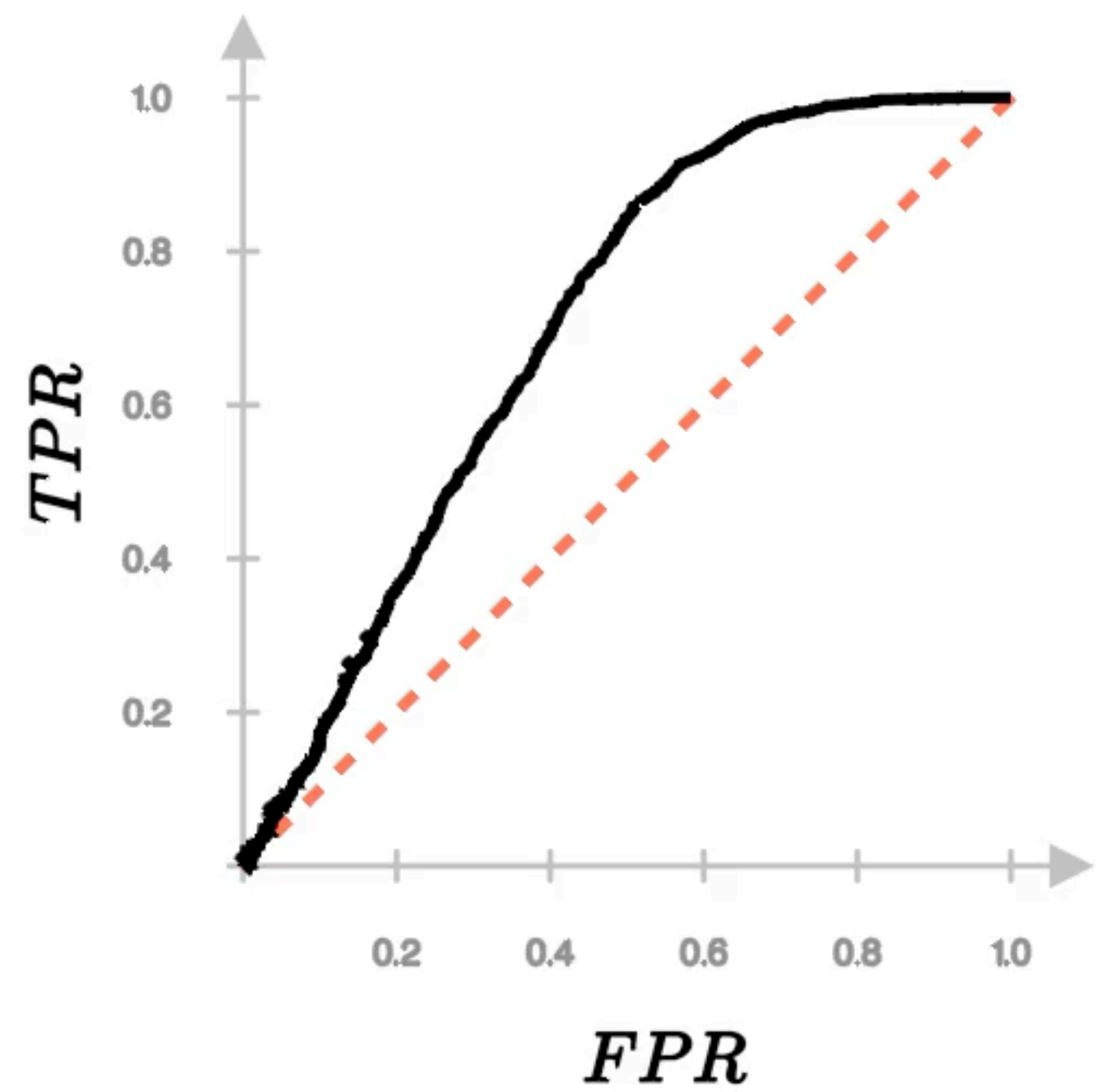
Overfitting

The behavior of the model on data distributions

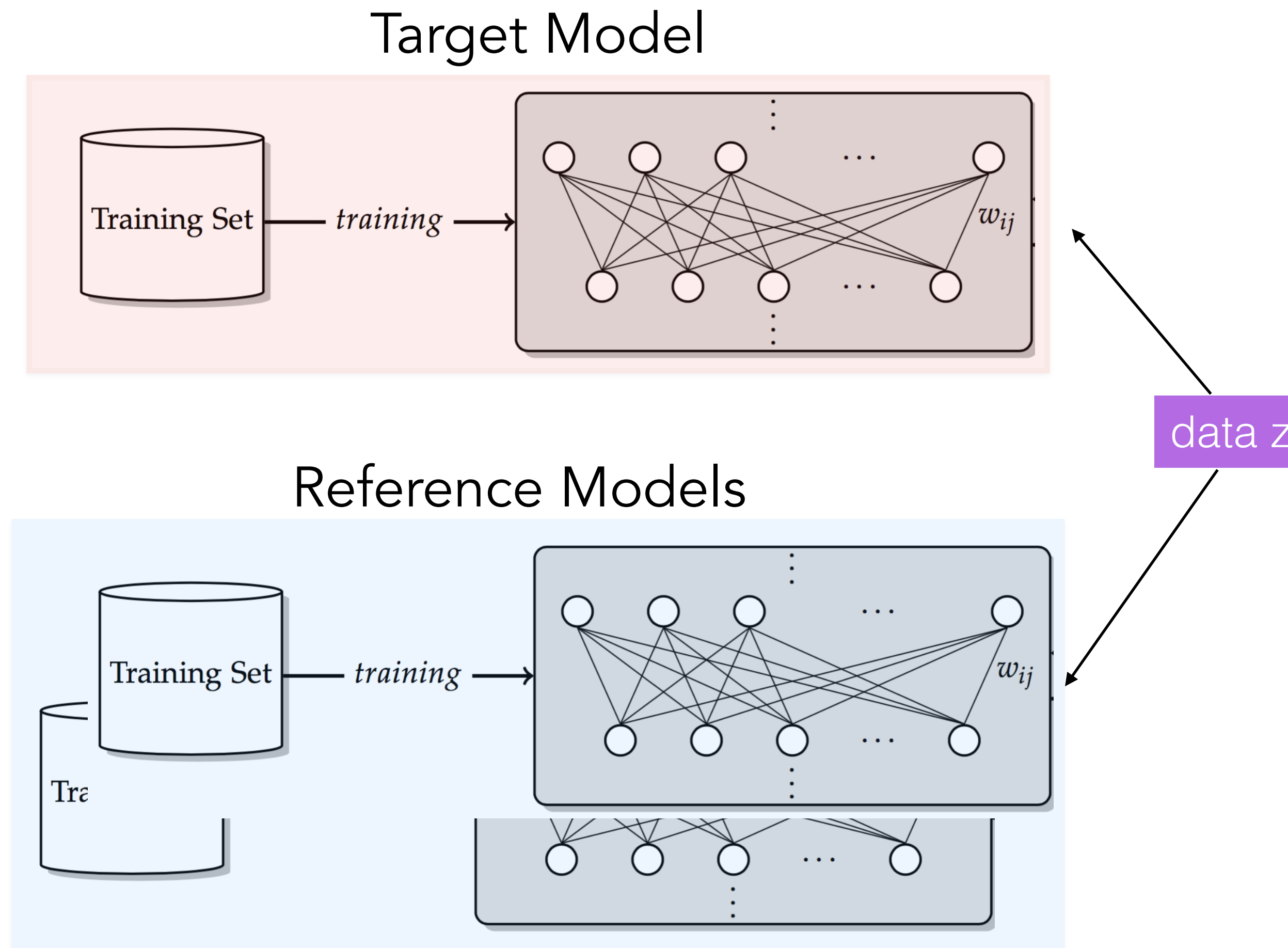


Where does this attack make errors?





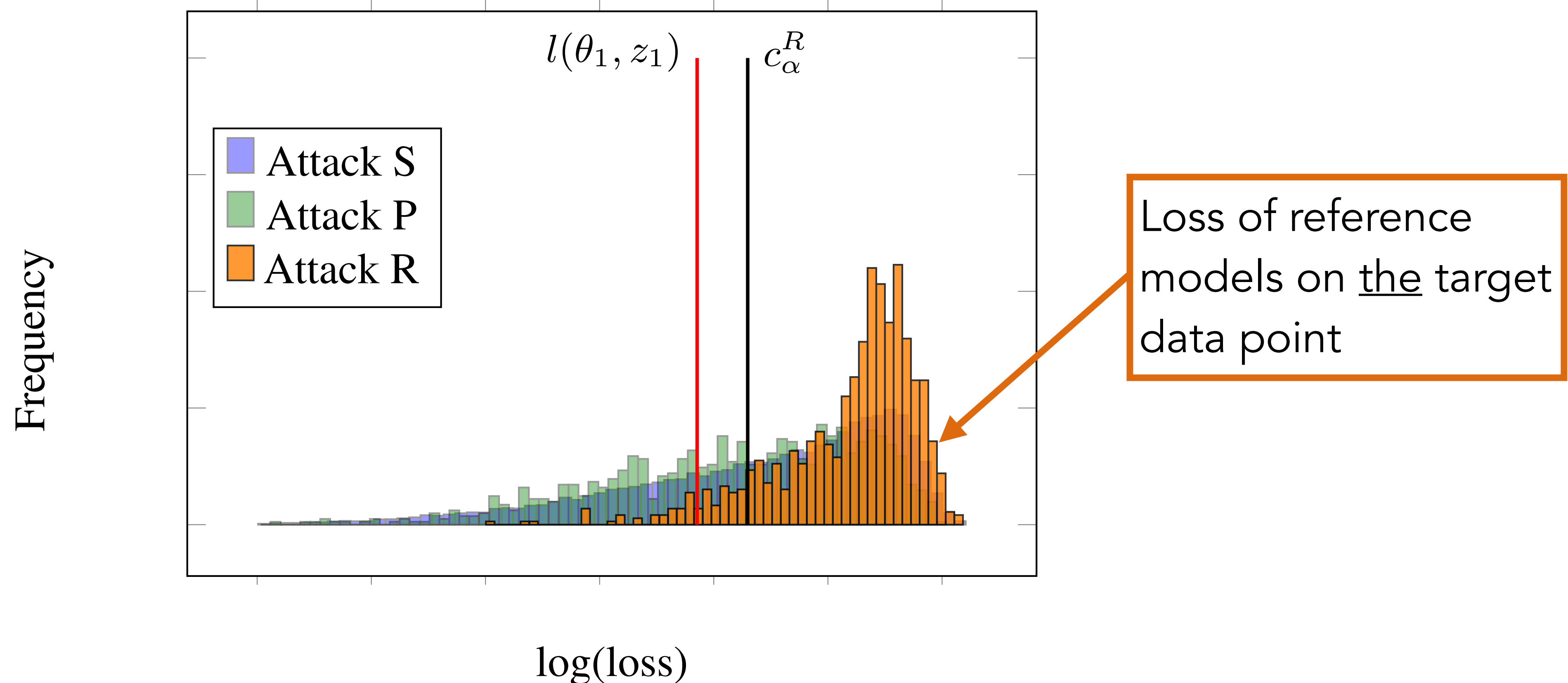
How to perform a more accurate analysis?



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 loss distribution of target data on reference models





Reason for Leakage?

Average Memorization The behavior of models on a data point, averaged over the remaining training data having been sampled from a distribution



Model's Error on x

Can we perform an even more accurate privacy analysis?

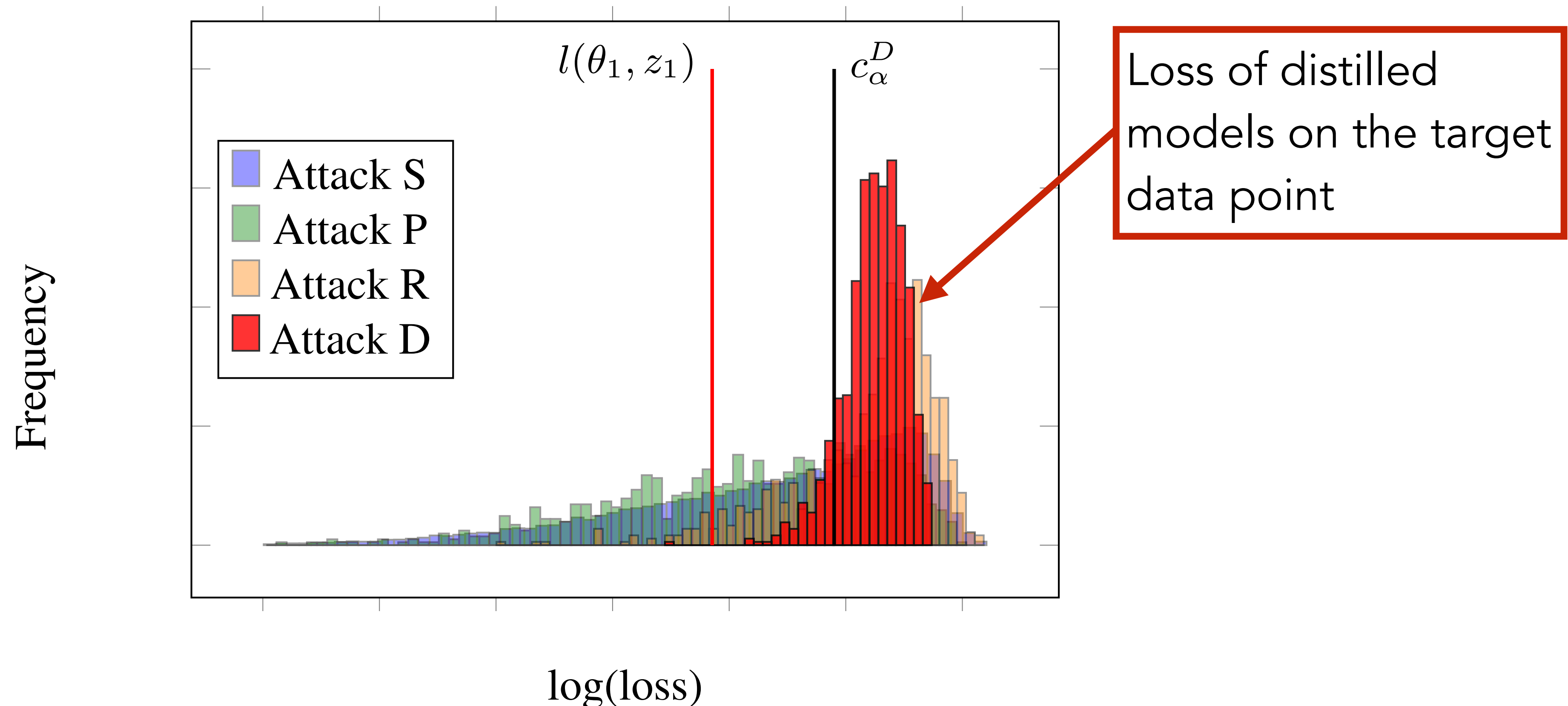
The objective is to get as close as possible to the leave-one-out attack, where the adversary knows all “other” data in the training set

- Train reference models that have a large agreement with the target model on all the training data, except the target data
- How? Use 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

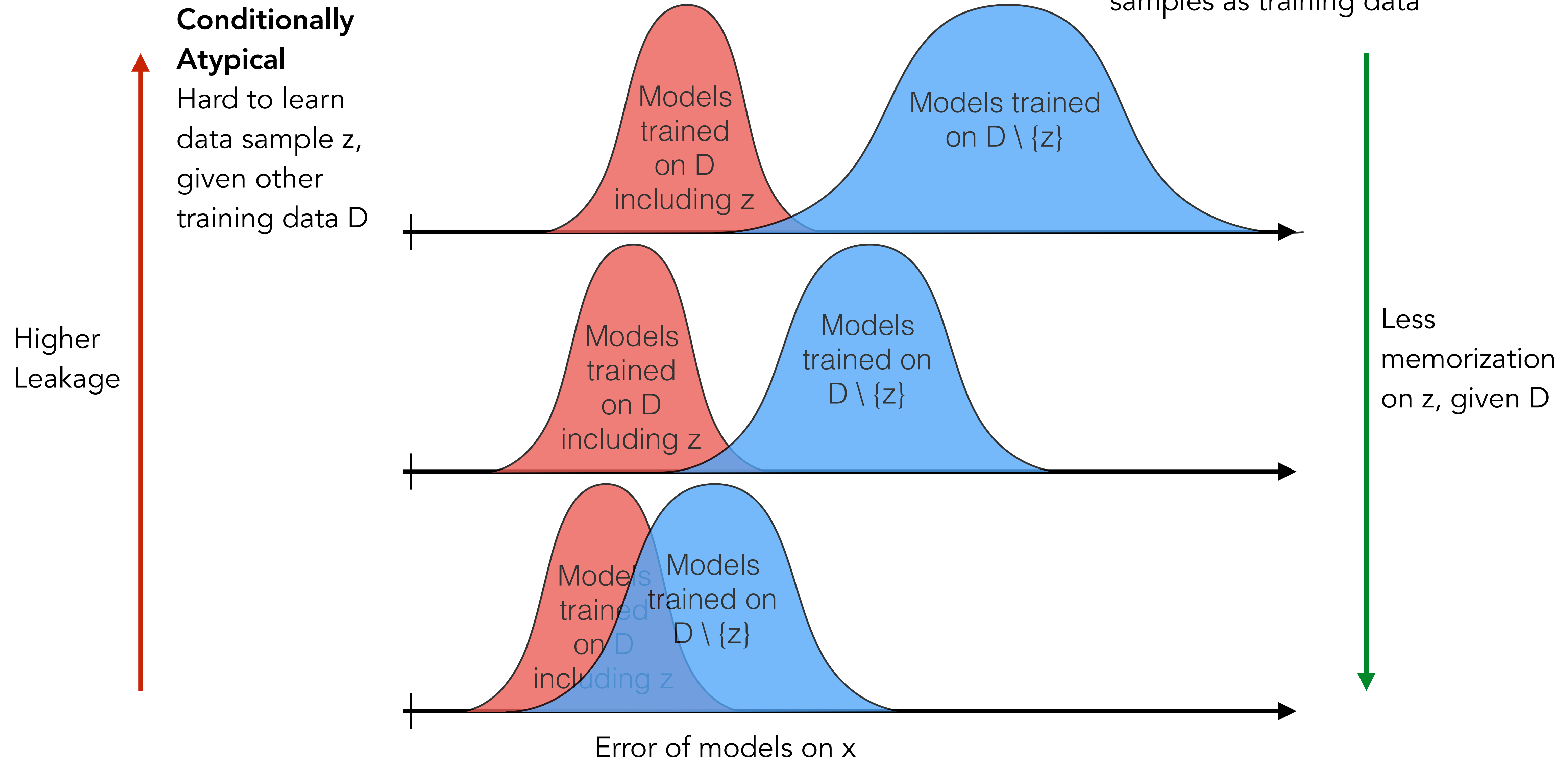
- Learn a threshold from the loss distribution of target data on distilled models
- Note that the threshold depends on both target data and the target model

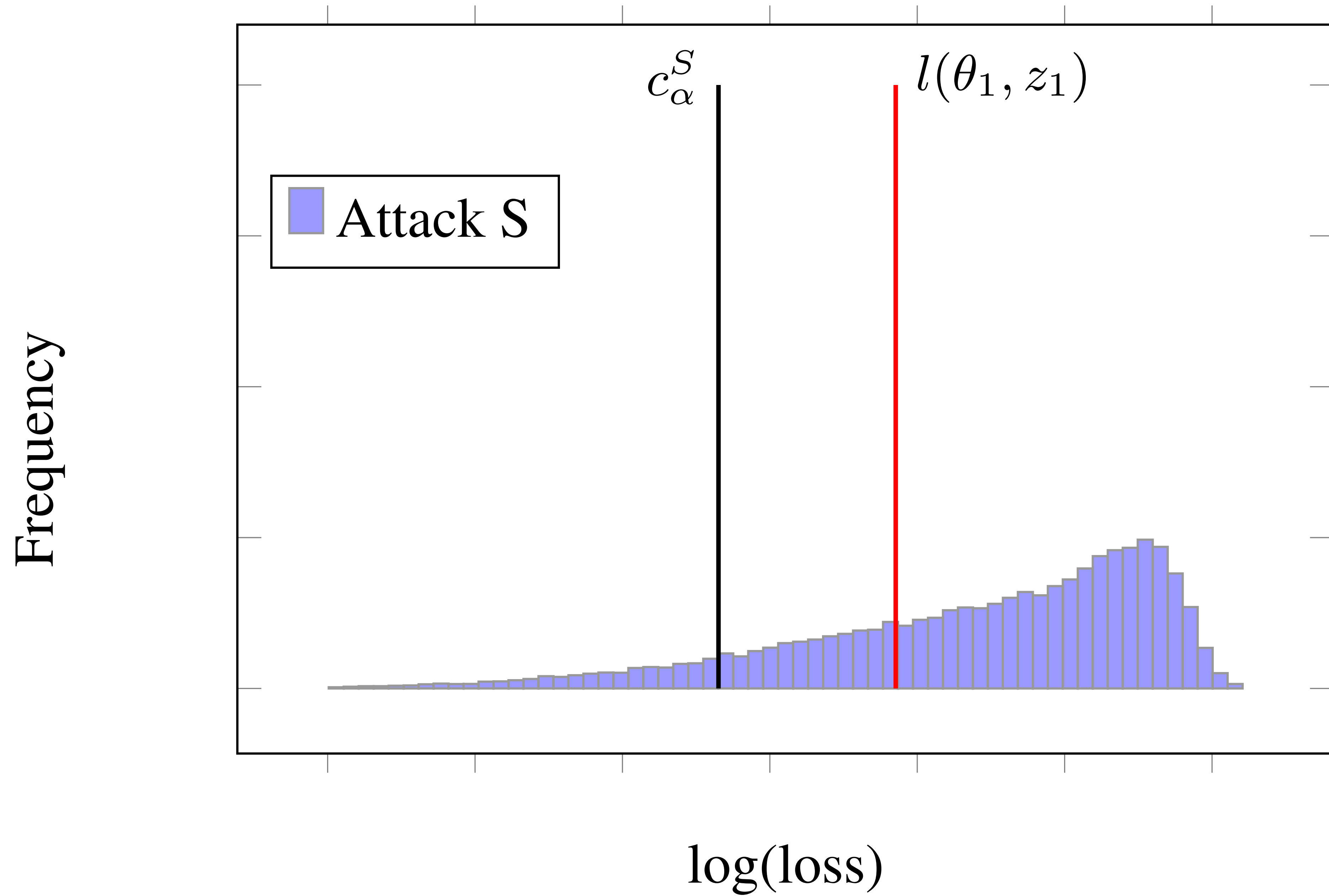


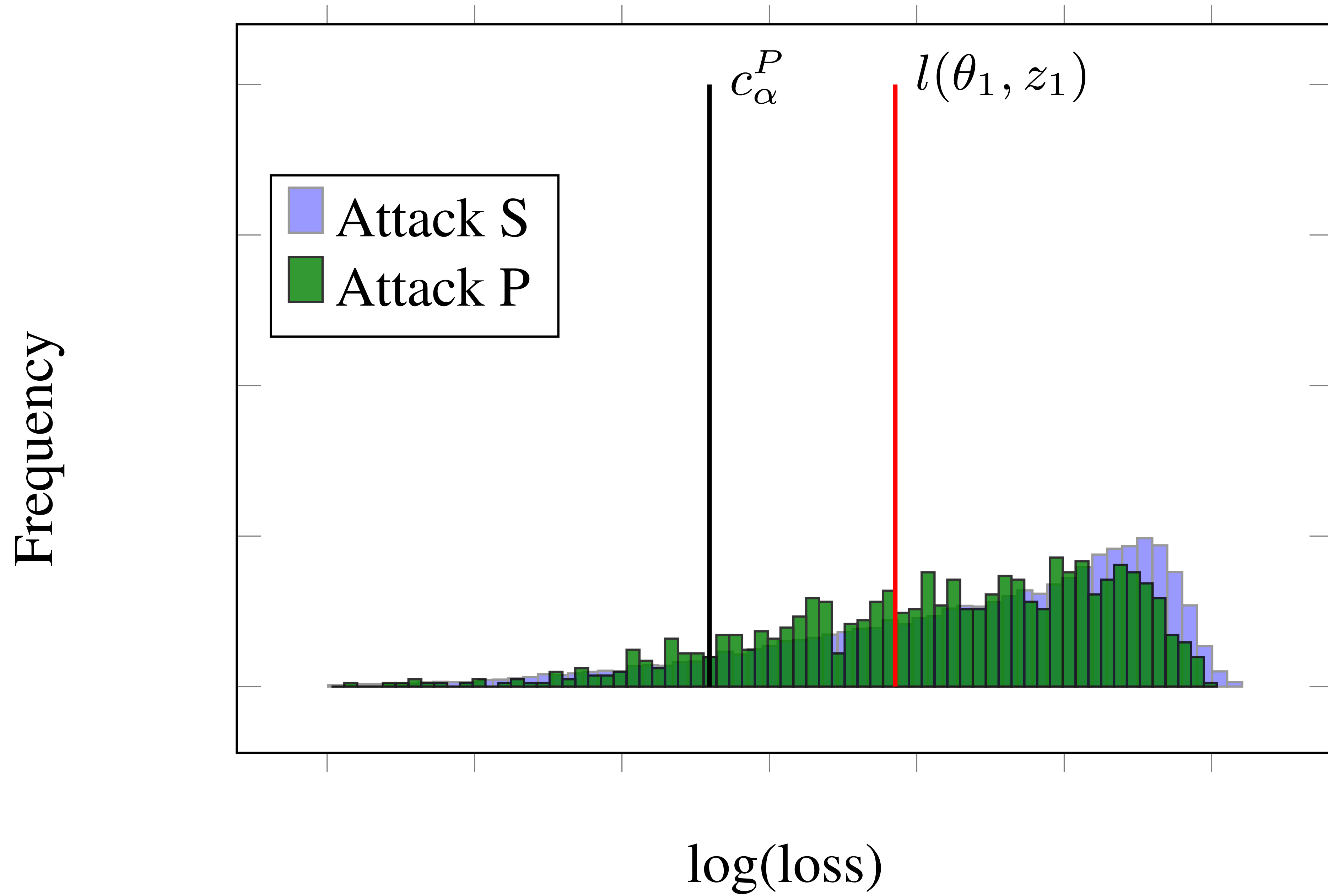
Reason for Leakage?

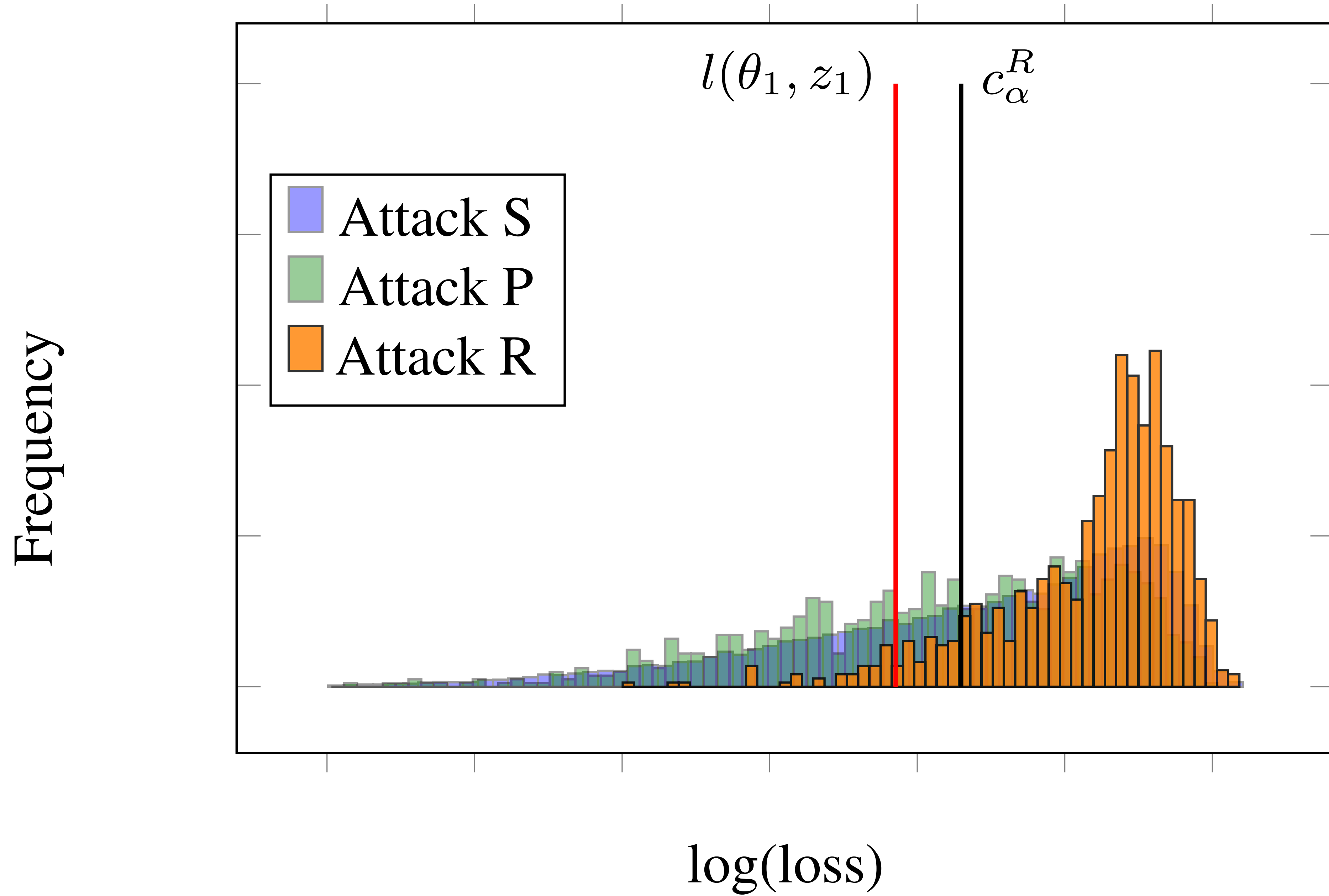
(Conditional) Memorization

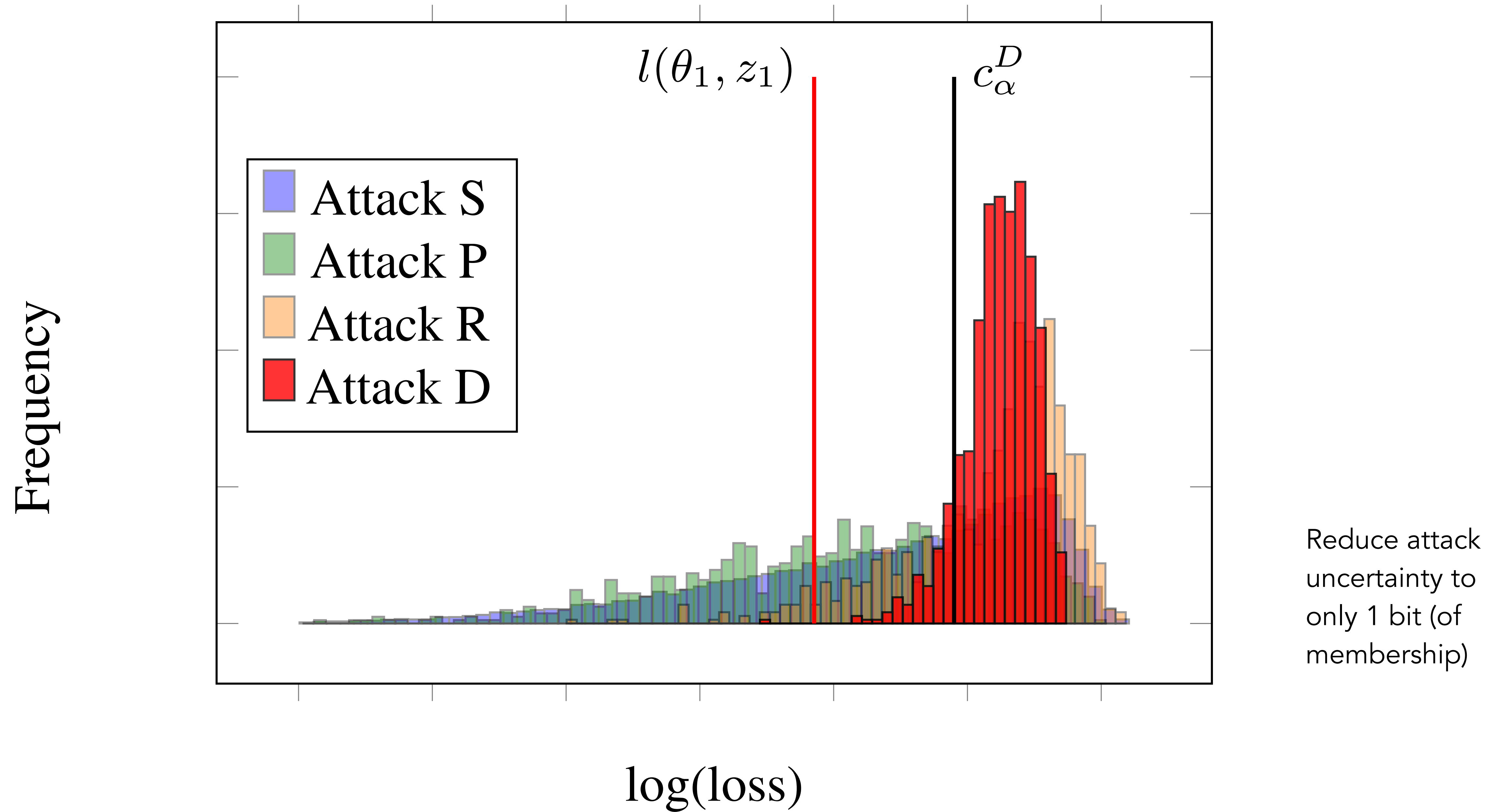
The behavior of models on a data point, given a specific set of samples as training data

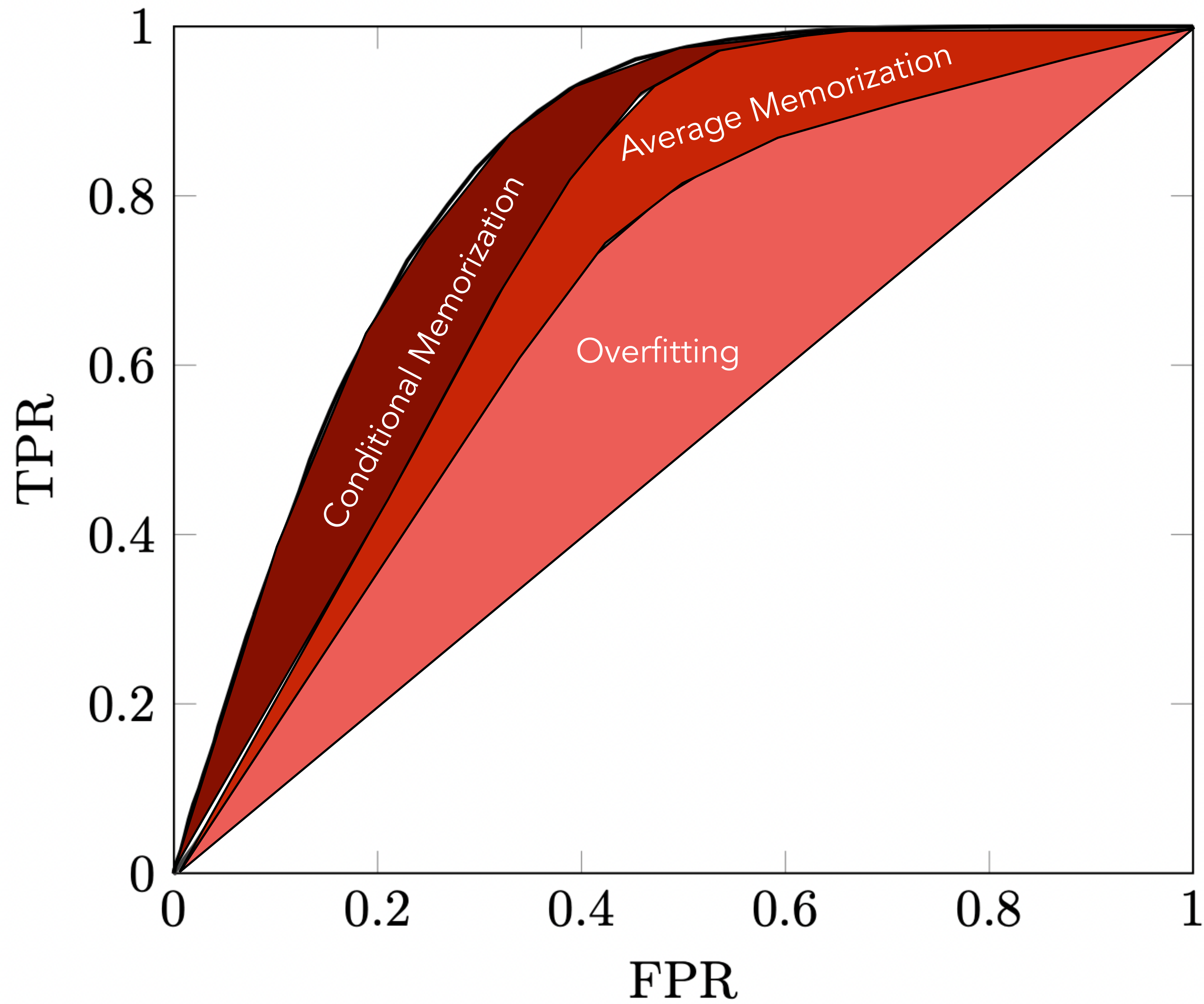












Auditing Data Privacy Using Privacy Meter

- 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
- Privacy Meter (privacy-meter.com) tool aids regulatory compliance, through a systematic method to audit data privacy for a wide range of machine learning algorithms

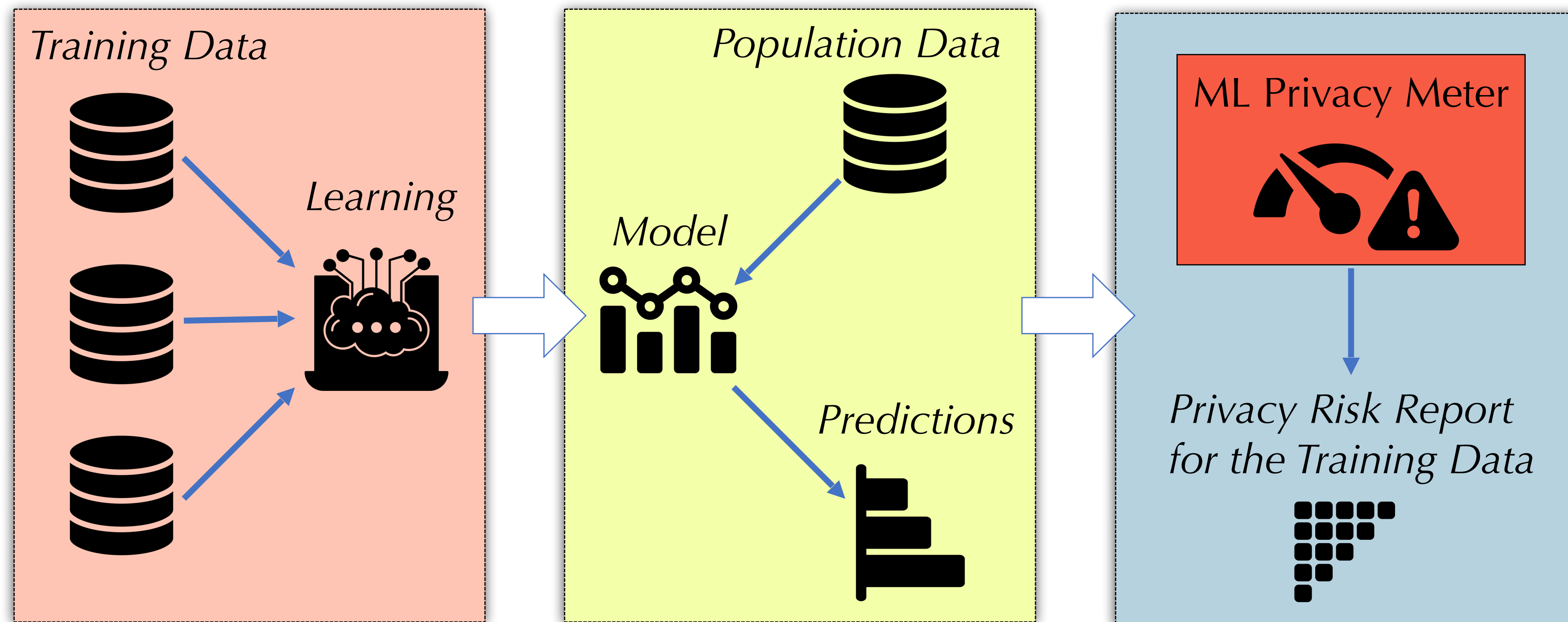
[Shokri, Stronati, Song, Shmatikov] Membership Inference Attacks against Machine Learning Models, SP'17

[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, Bindschaedler, Shokri] Enhanced Membership Inference Attacks against Machine Learning Models, CCS'22

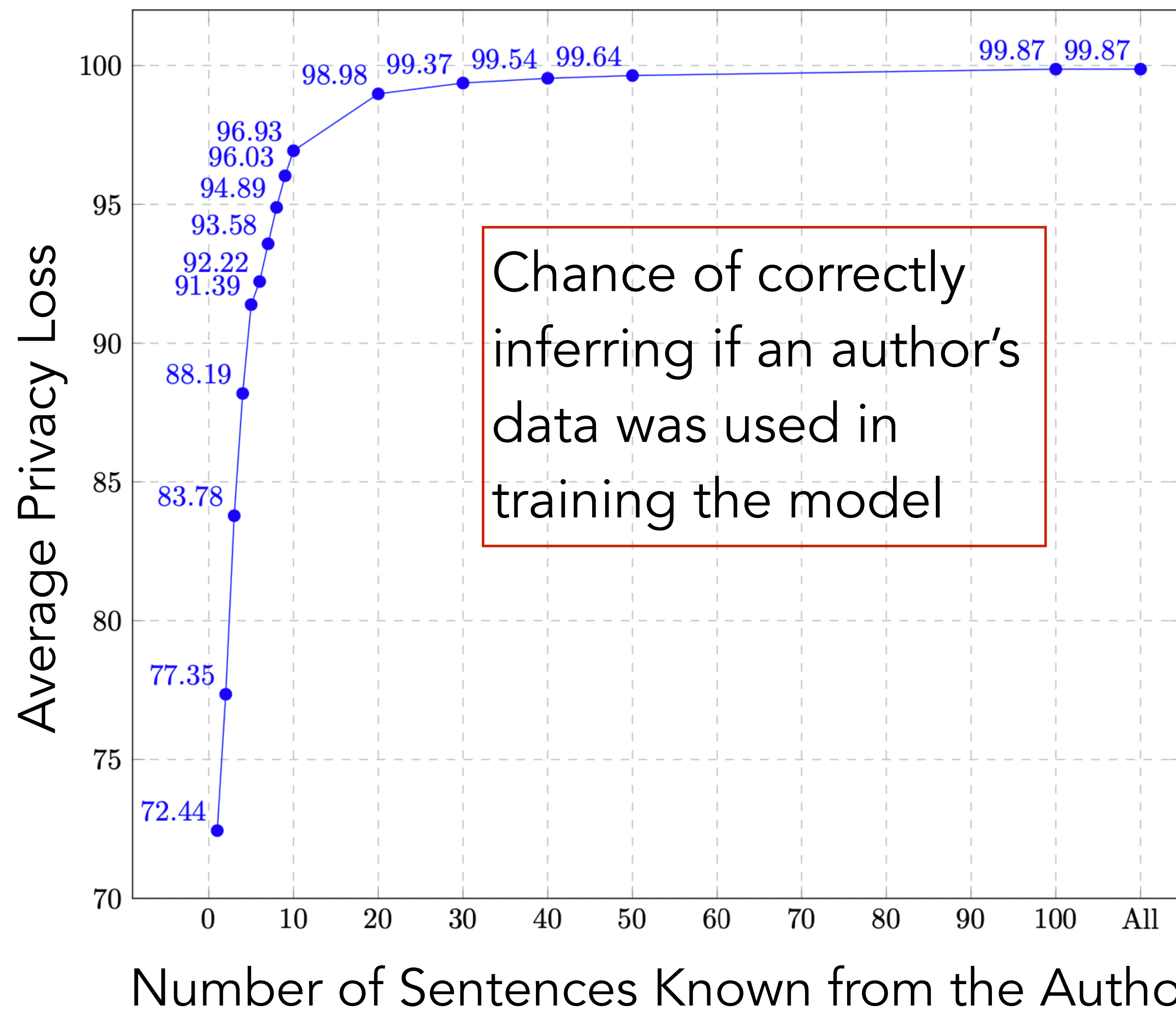
Privacy Meter

privacy-meter.com

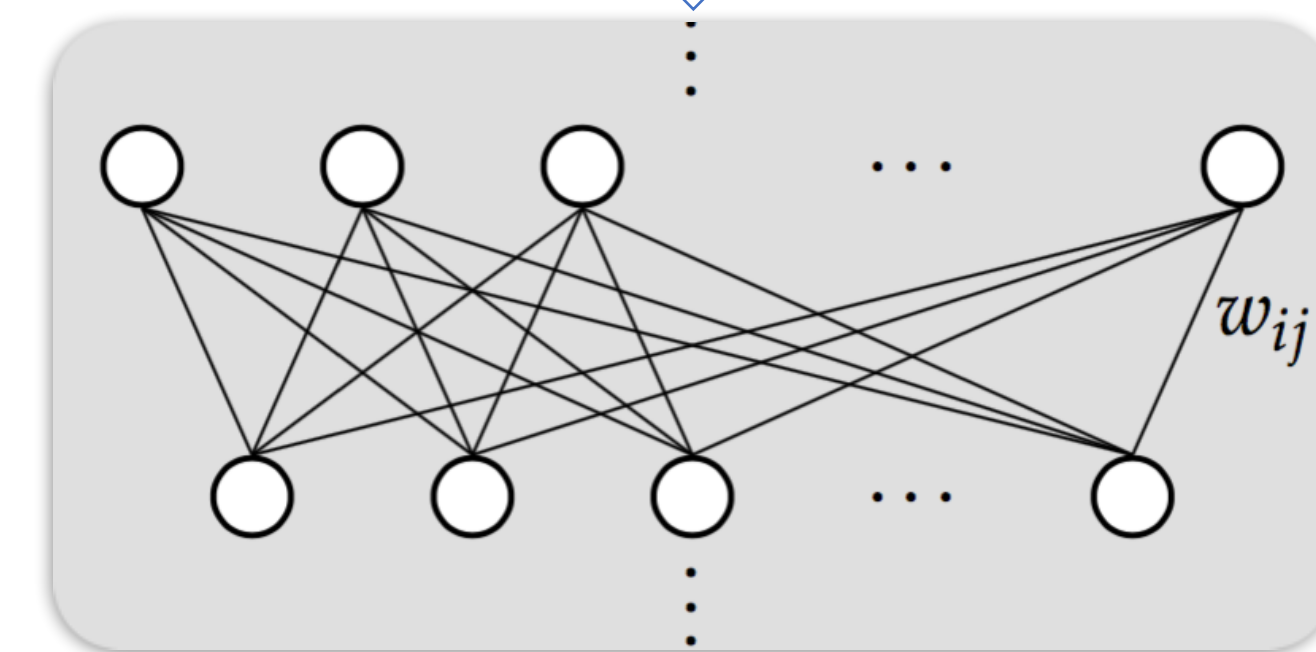
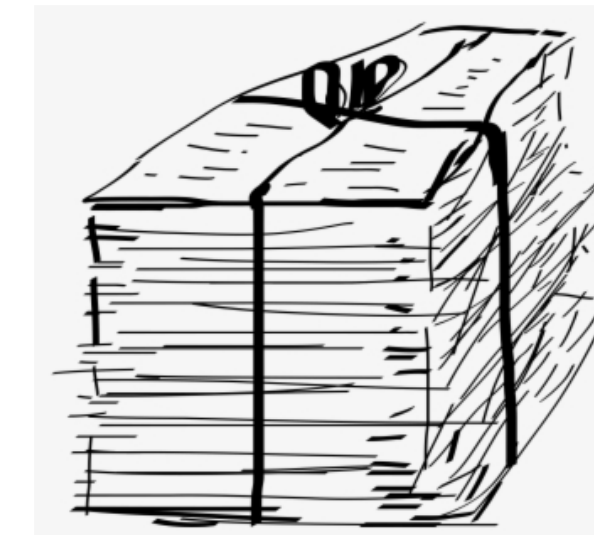


Privacy Meter is an open source tool that enables quantifying the privacy risks of statistical and machine learning models.

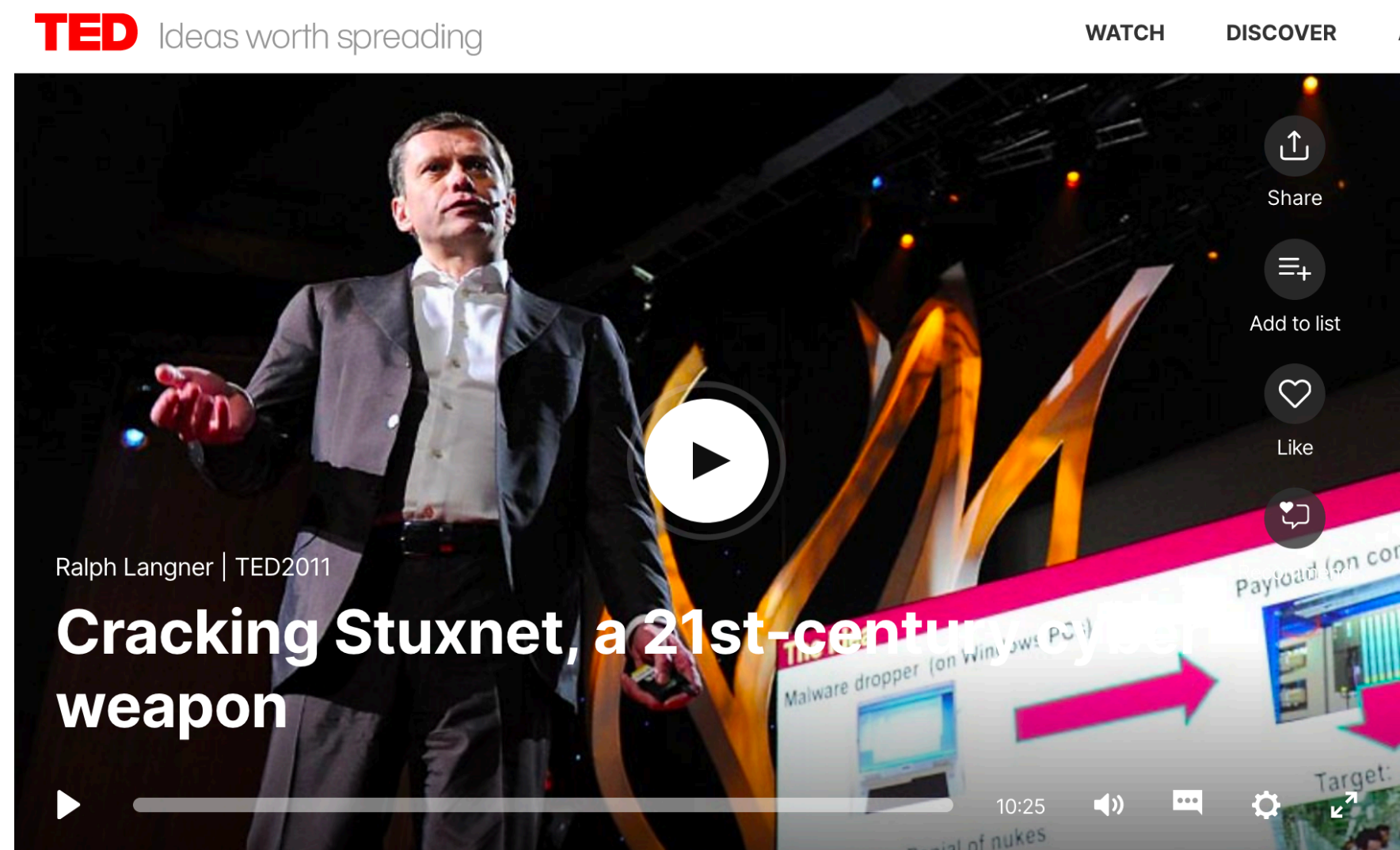
Example: Language Generative Model



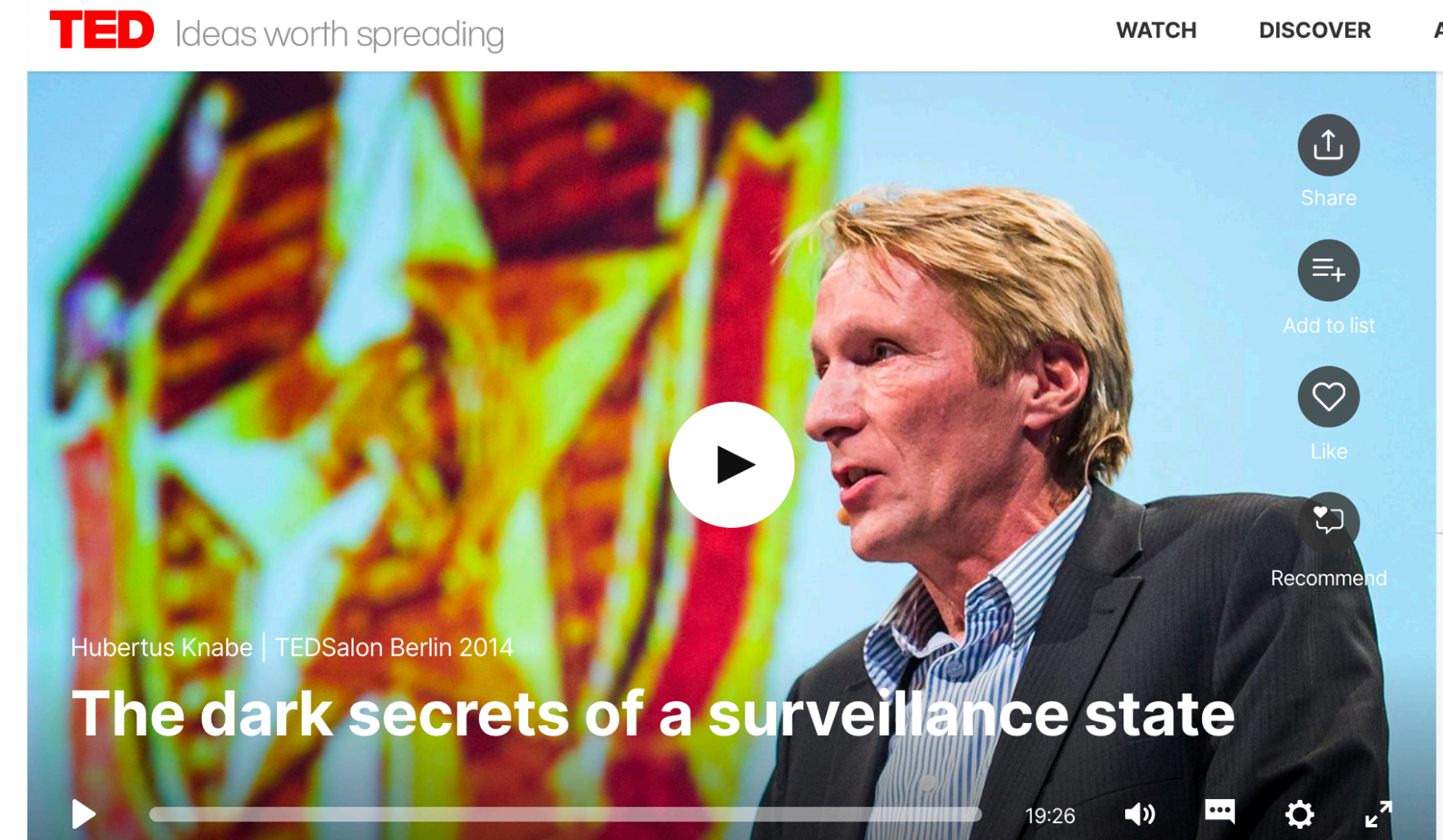
(Speaker Annotated TED talks)
SATED dataset



Examples of Vulnerable Training Data



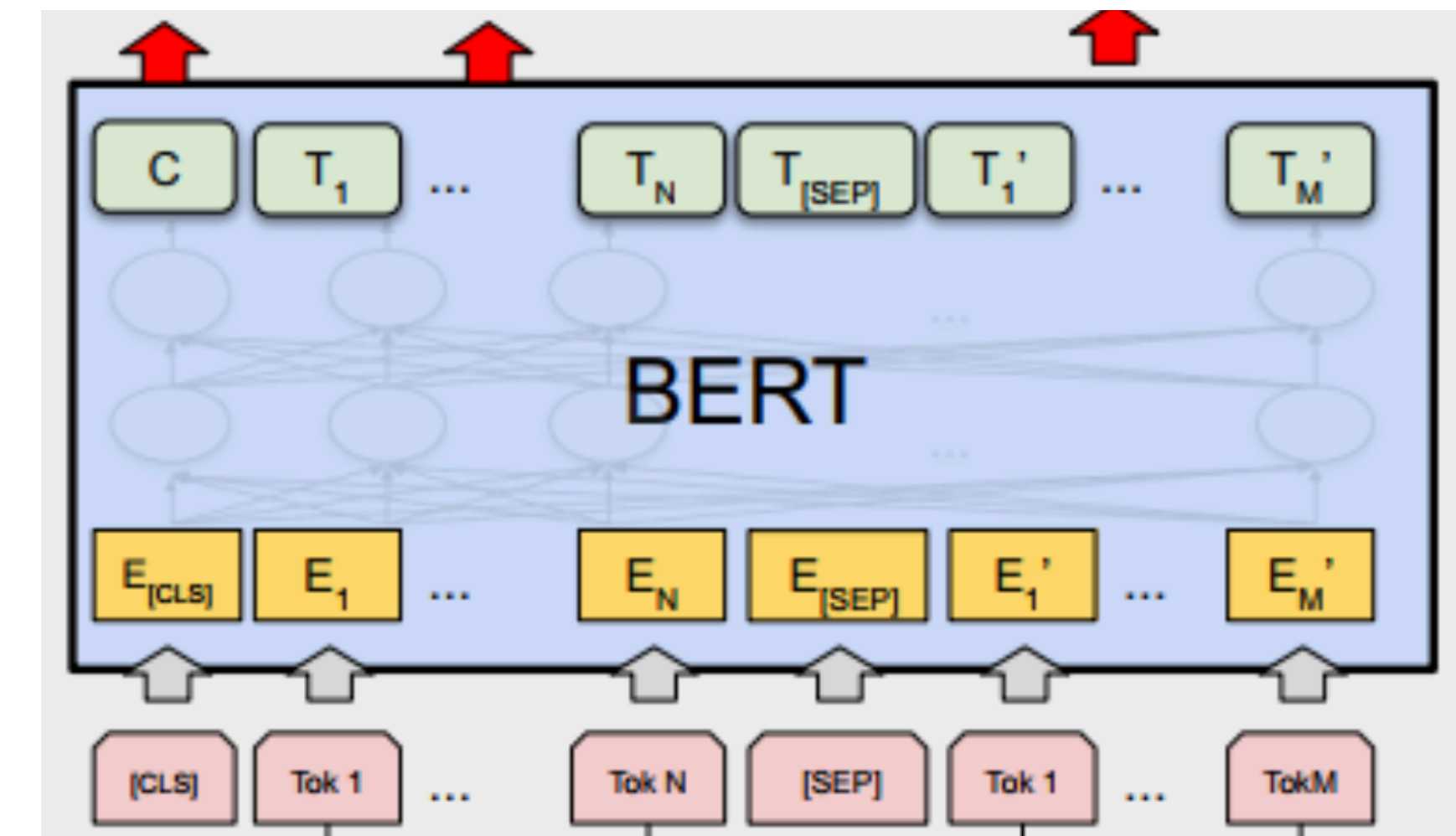
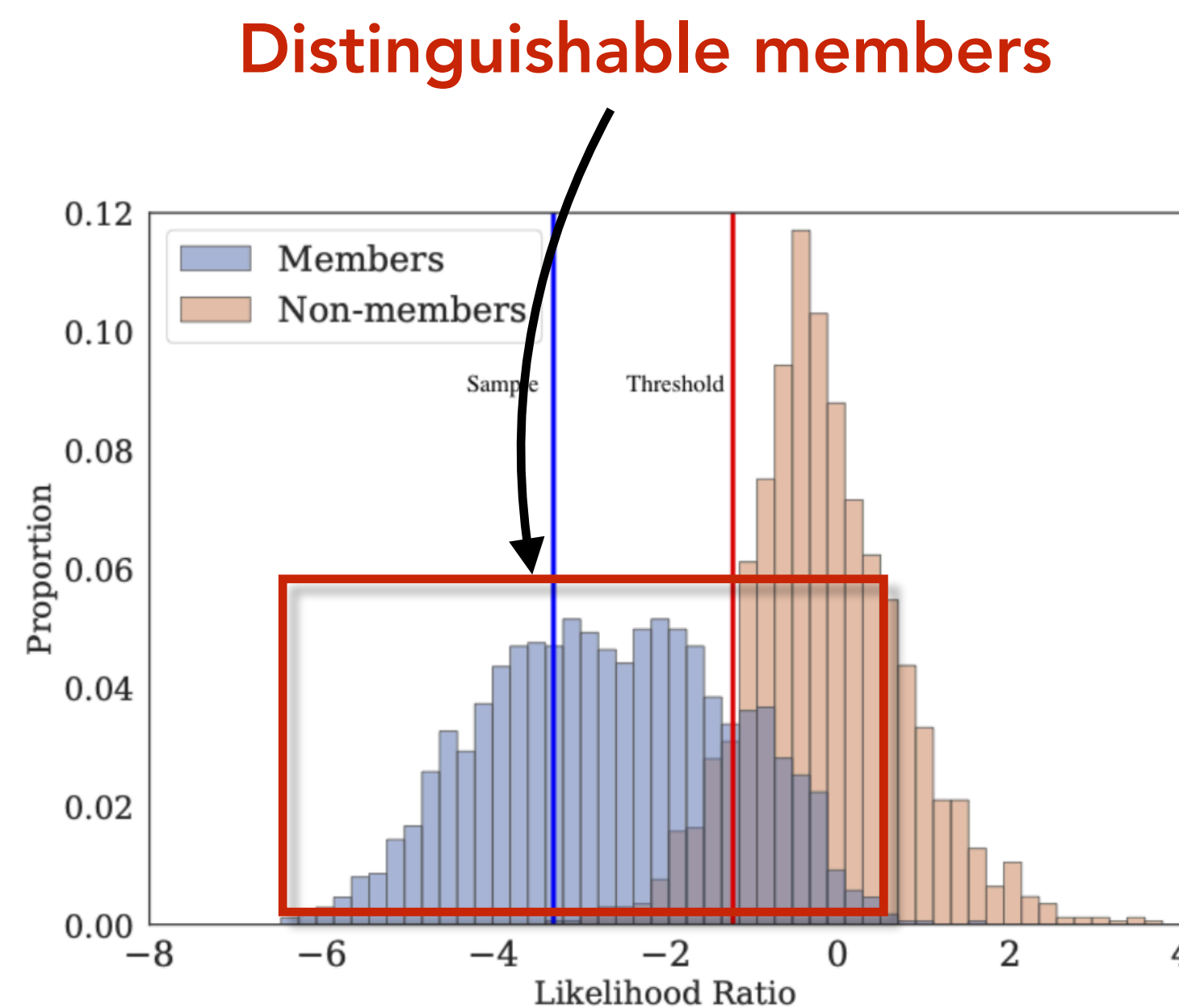
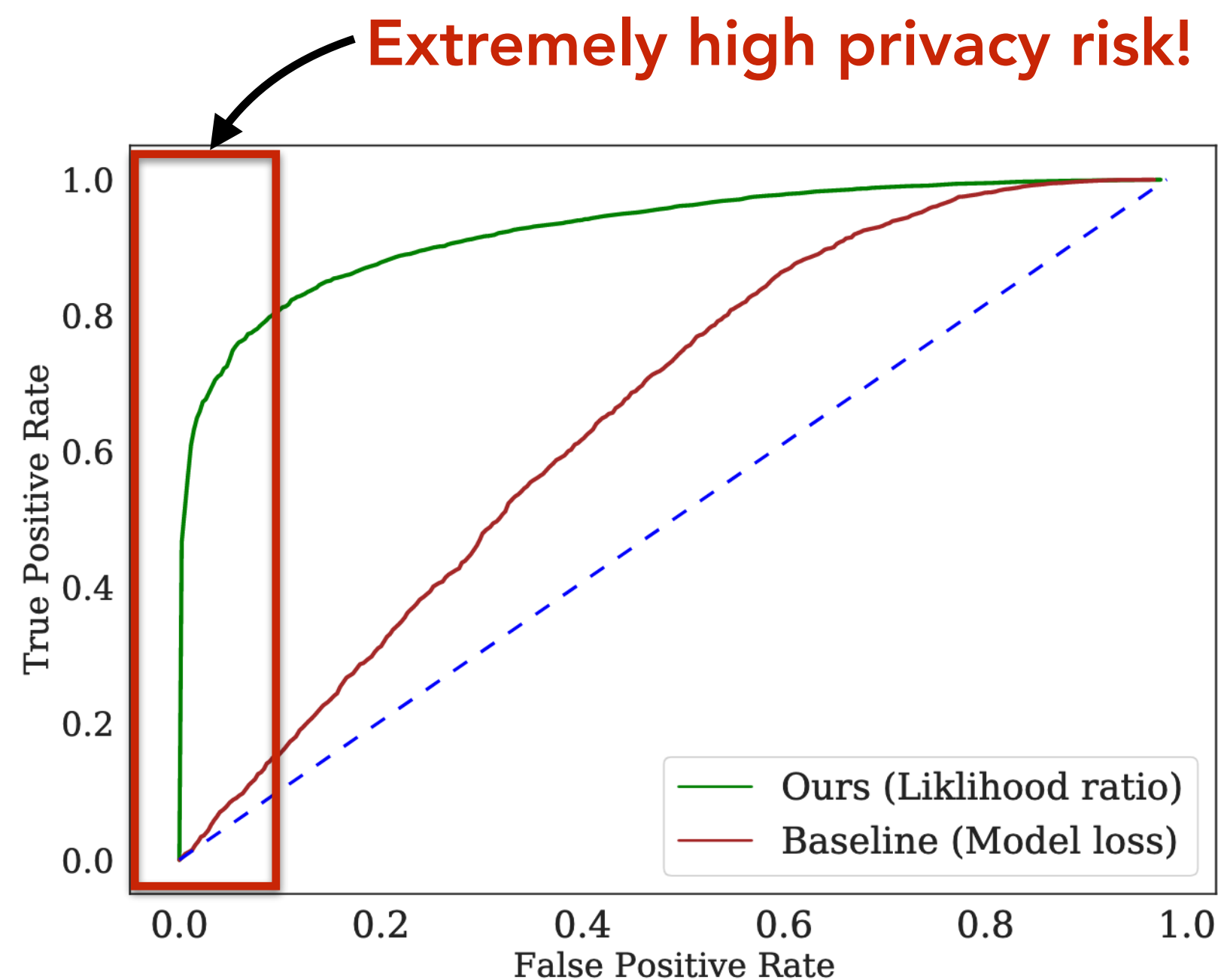
But it gets worse. And this is very important, what I'm talking about is very generic. It doesn't have anything to do, in specifics, with Stuxnet. It would work as well, for example, in a power plant or in a nuclear reactor. You don't have to deliver it -- as an attacker -- you don't have to deliver it. In 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, the German Democratic Republic, the GDR, in the East was merged with the West to form today's Germany. Among many countries, the East German secret police, known as the Stasi. Or the files were opened to the public, and historians such as me have written about how the GDR surveillance state functioned.

Example: Masked Language Models

- Members of the training set are identifiable: Presence of any document in a training dataset can be inferred very accurately using membership inference attacks



ClinicalBERT



Example: Image Classification Tasks

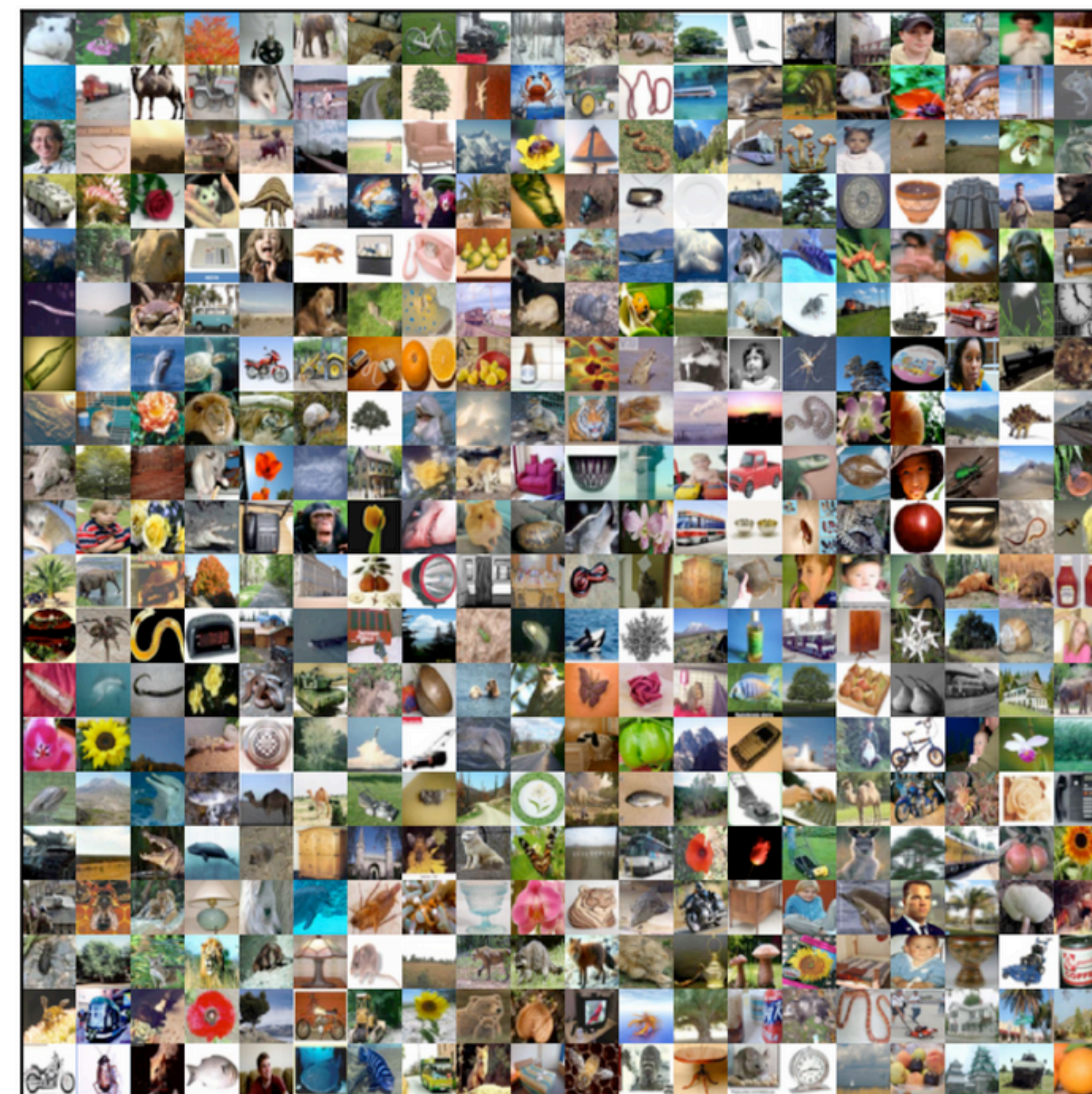
Model	Number of Parameters	Prediction (Test) Accuracy	Privacy Risk
AlexNet	2.47 million	44%	75.1%
ResNet	1.7 million	73%	64.3%
DenseNet	25.62 million	82%	74.3%

Large capacity

High generalizability

Low privacy

CIFAR100 Image classification

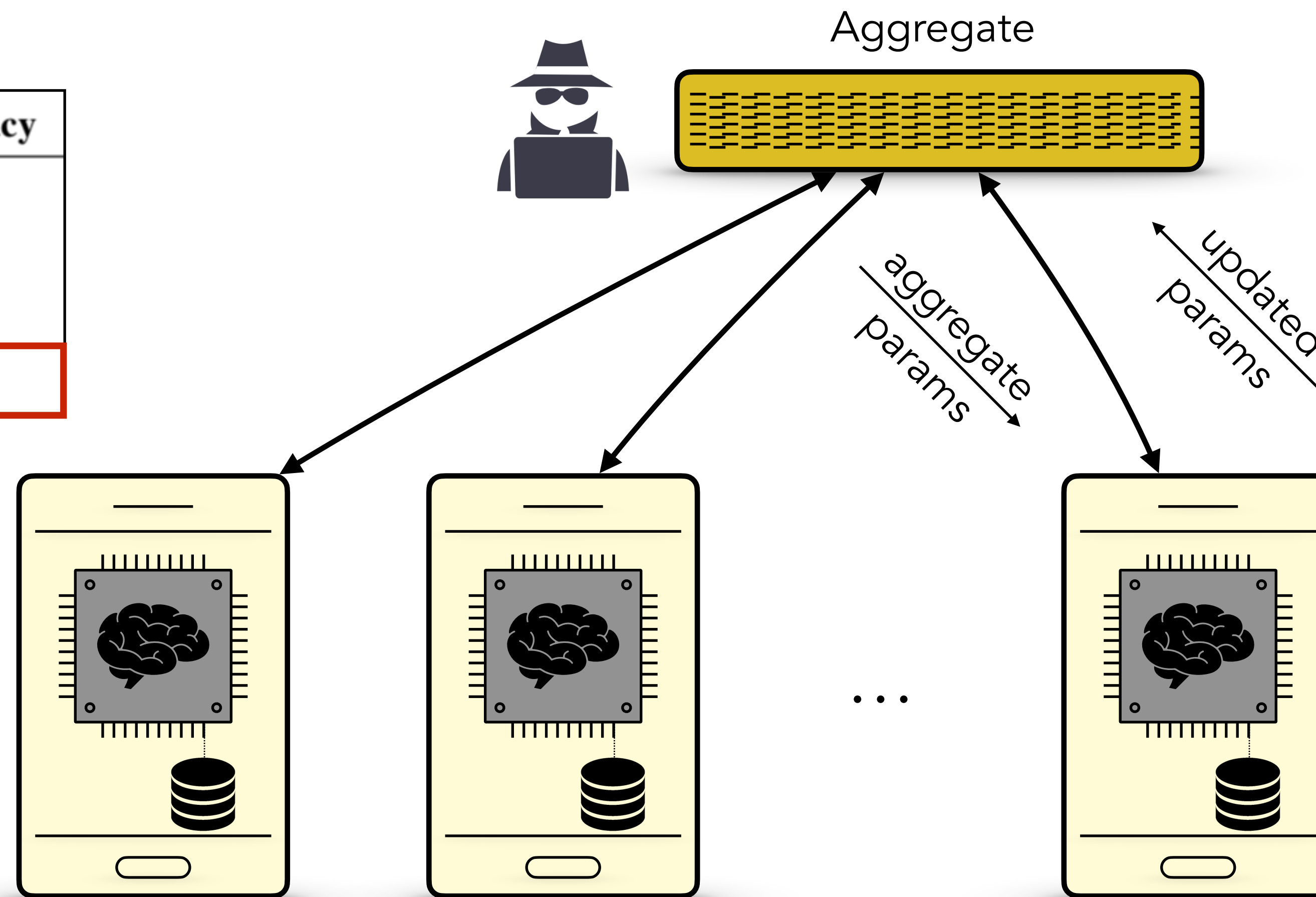


Example: Federated Learning

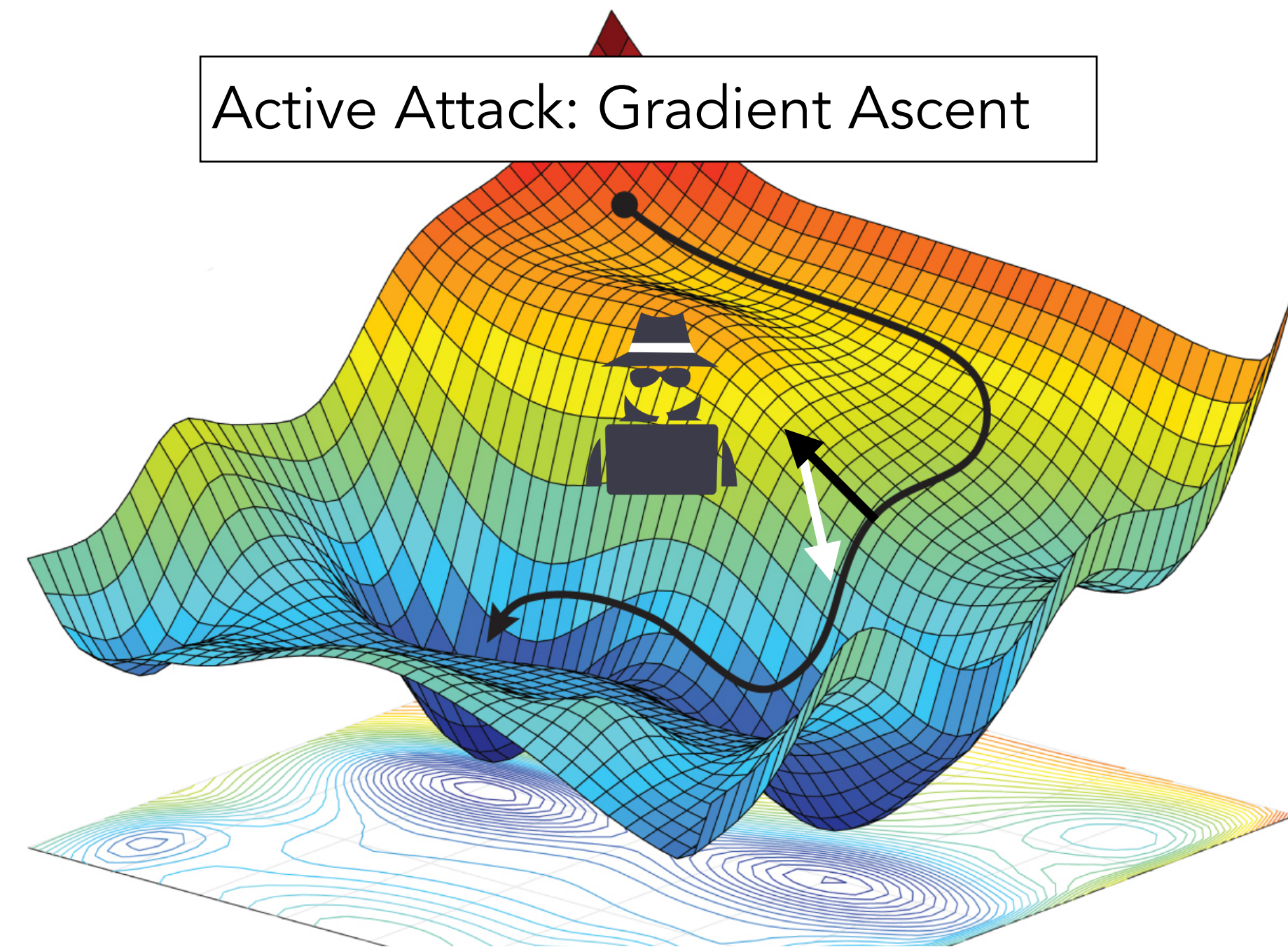
Adversary can observe multiple snapshots of the model

Observed Epochs	Attack Accuracy
5, 10, 15, 20, 25	57.4%
10, 20, 30, 40, 50	76.5%
50, 100, 150, 200, 250	79.5%
100, 150, 200, 250, 300	85.1%

CIFAR100-Alexnet

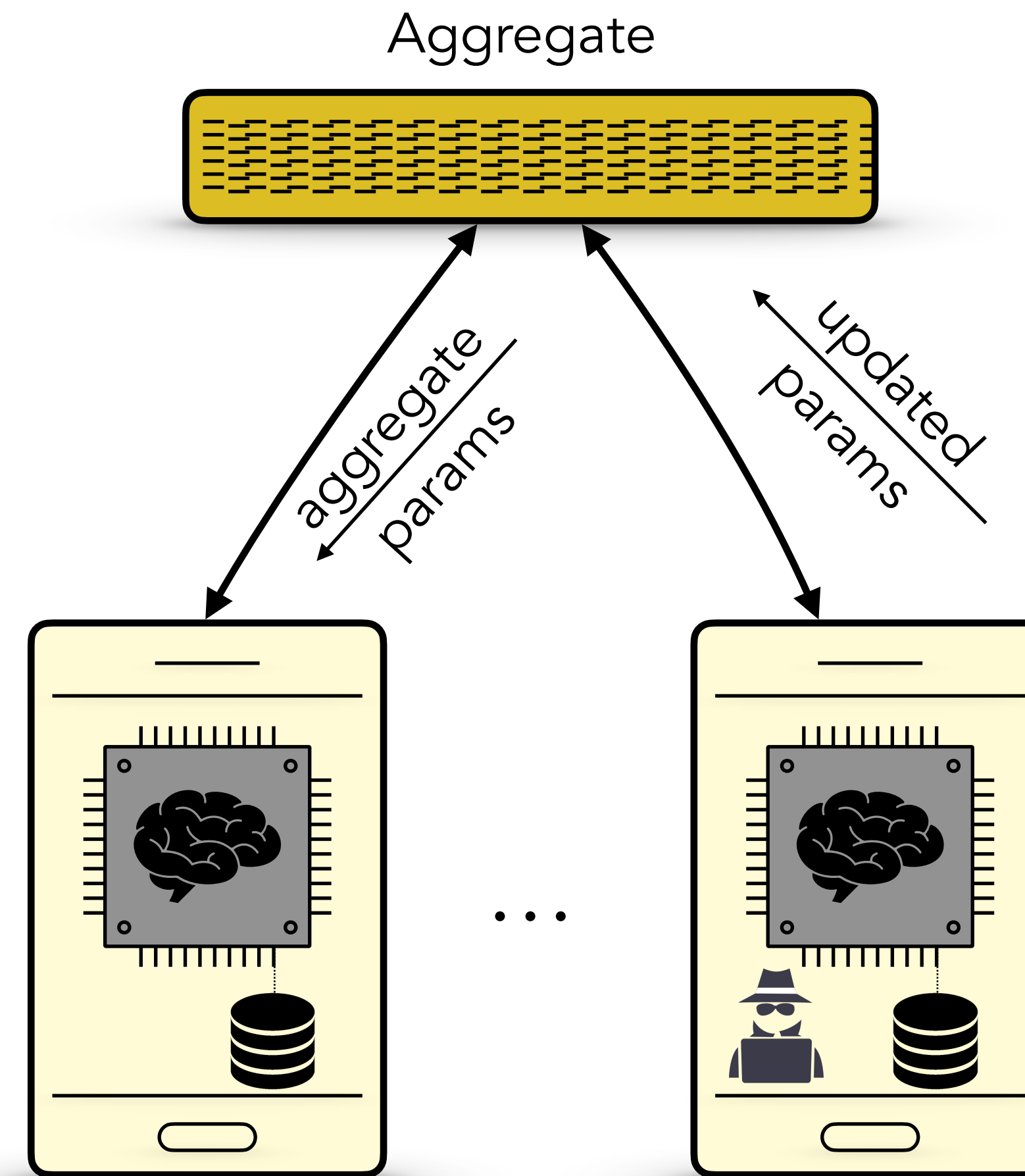


Decentralized (Federated) Learning

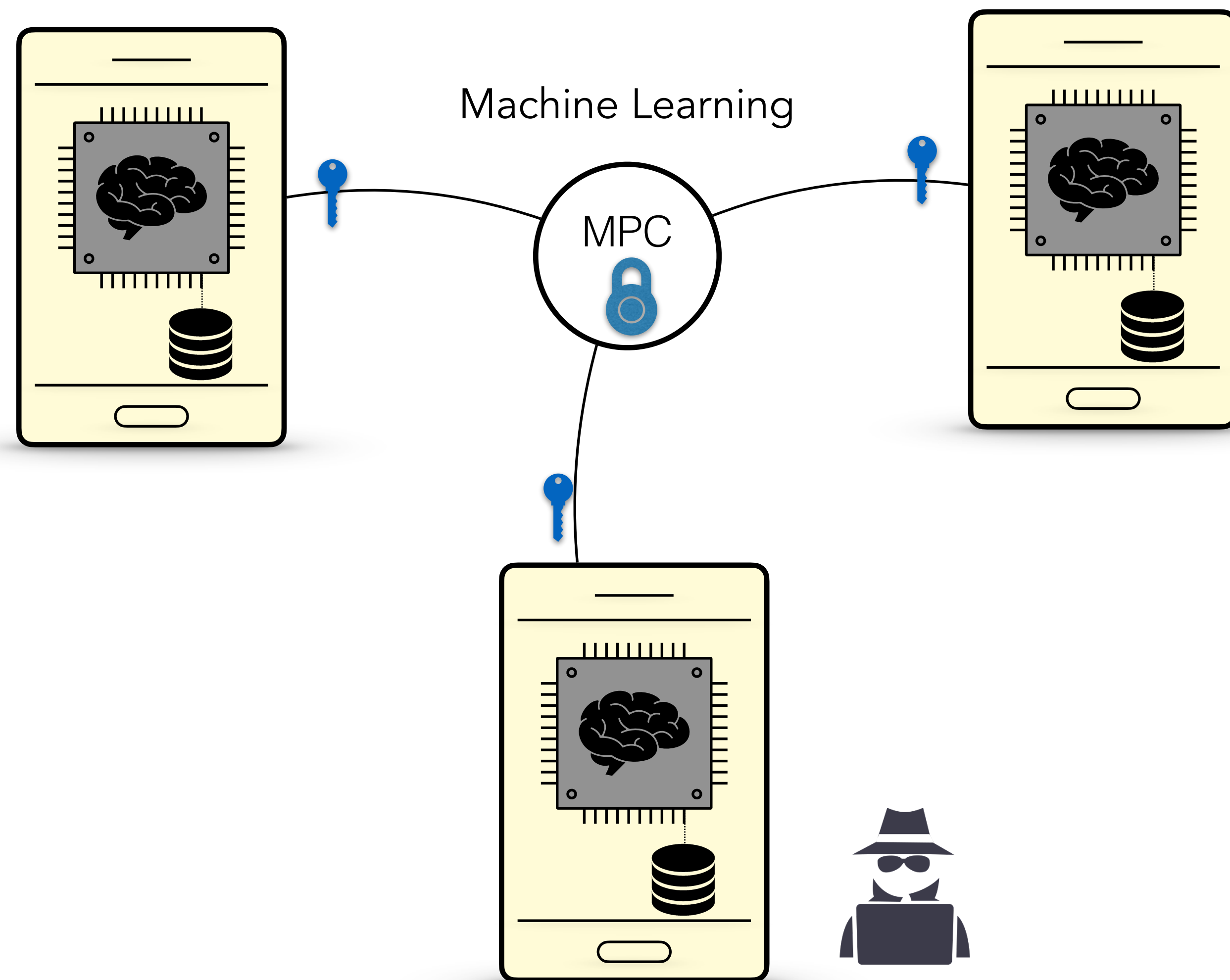


Increase loss on a particular data point x .

A participant corrects it back (by running gradient descent locally) only if x is part of its training set. => **membership leakage**

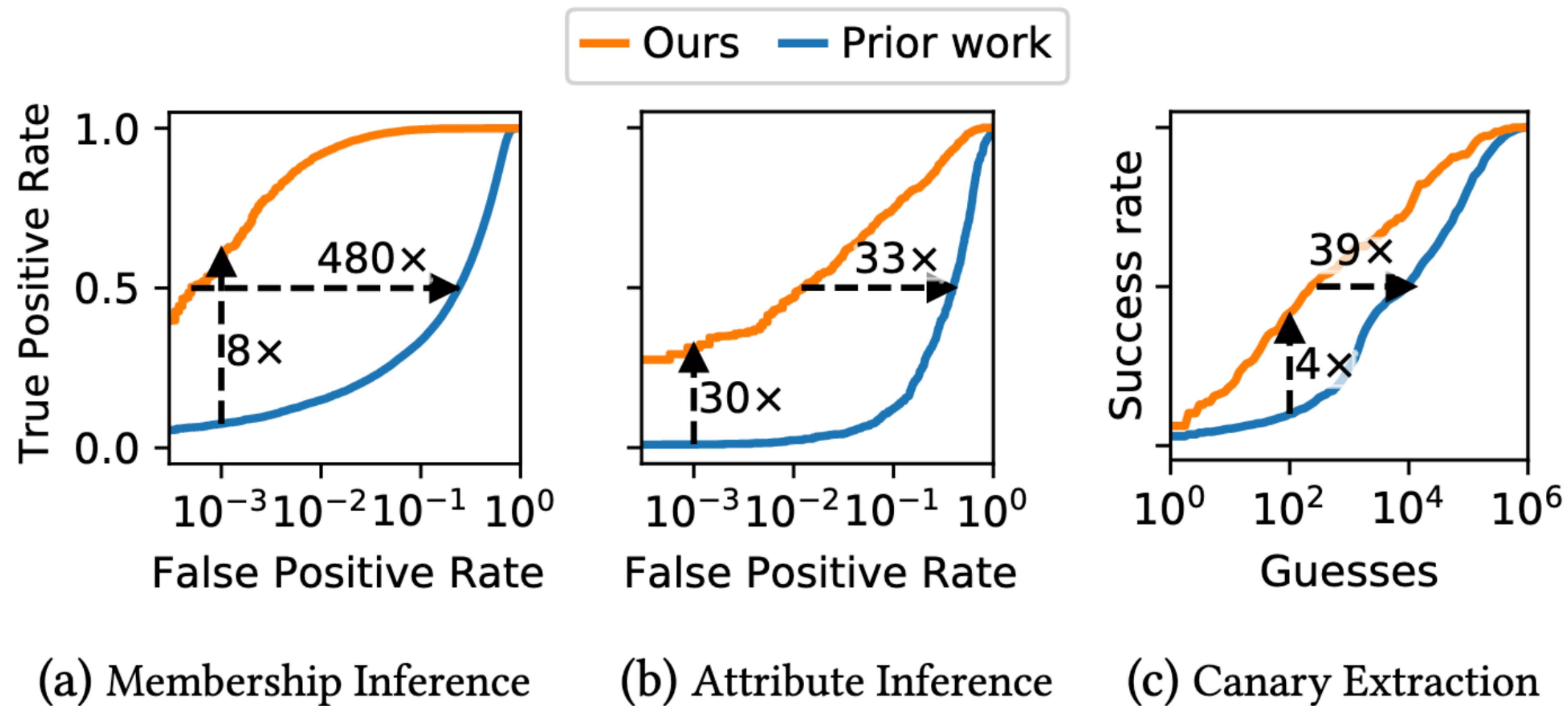


Example: Secure Multi-Party Computation



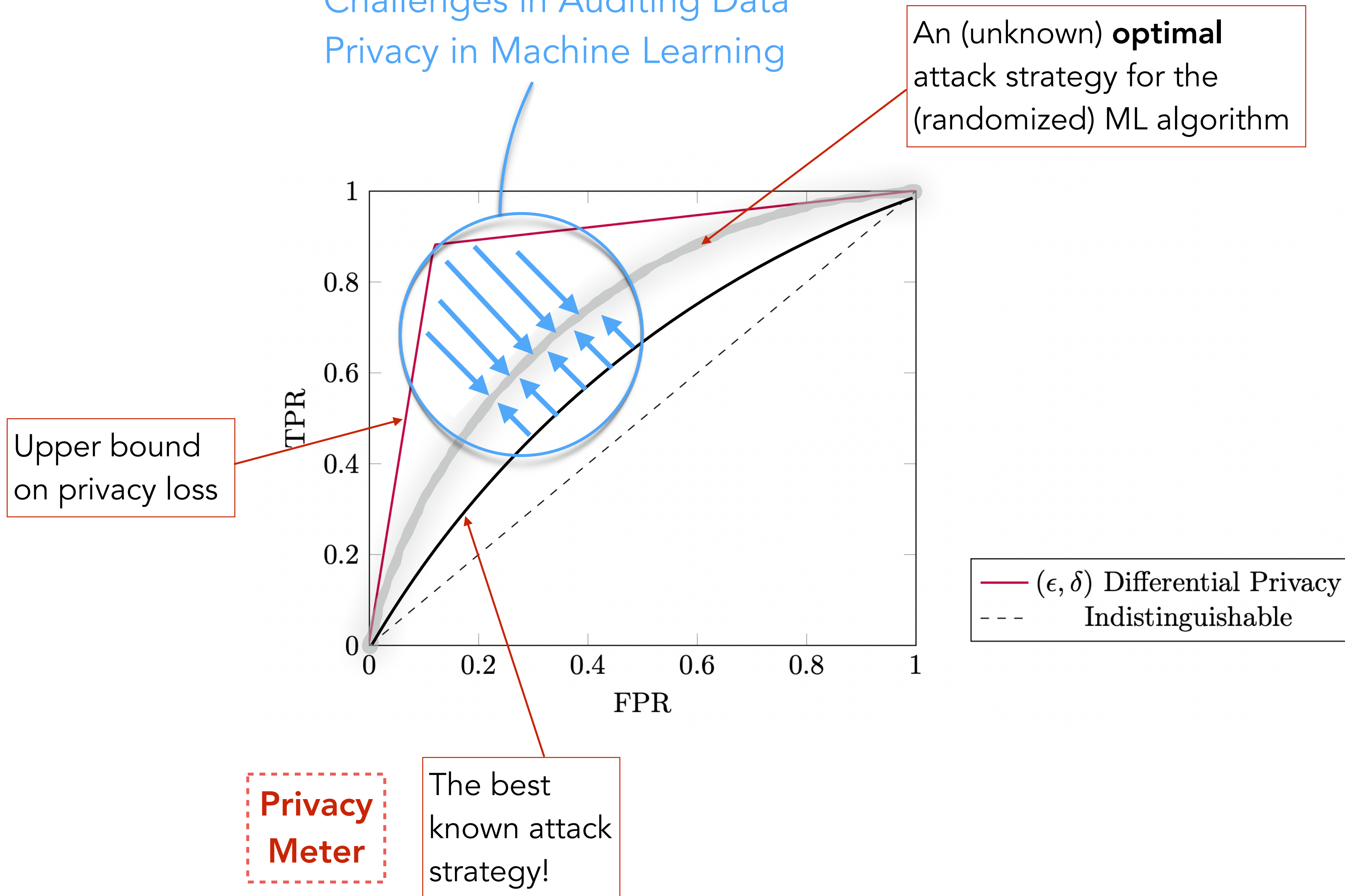
- No data is shared
- No entity can observe the intermediate steps of the computation
- The final model, however, is available to all parties
- New Attack: Adversary poisons his dataset to increase information leakage from other parties! **Exploit memorization.**

Example: Secure Multi-Party Computation



Conclusions

Challenges in Auditing Data Privacy in Machine Learning



Other challenges:

Alleviating the potential tension between privacy and

- Generalizability
- Robustness
- Fairness
- Explainability
- Scalability