

Data Privacy in Machine Learning

Reza Shokri

Data Privacy and Trustworthy ML Research Lab
National University of Singapore



reza@comp.nus.edu.sg



@rzshokri

Threats to Data Privacy

Threats to Data Privacy

- Unauthorized access to data, and data breaches
- Massive data collection

Threats to Data Privacy

Direct and intentional leakage

- Unauthorized access to data, and data breaches
- Massive data collection

Threats to Data Privacy

Direct and intentional leakage

- Unauthorized access to data, and data breaches
- Massive data collection

Indirect and unintentional leakage

Threats to Data Privacy

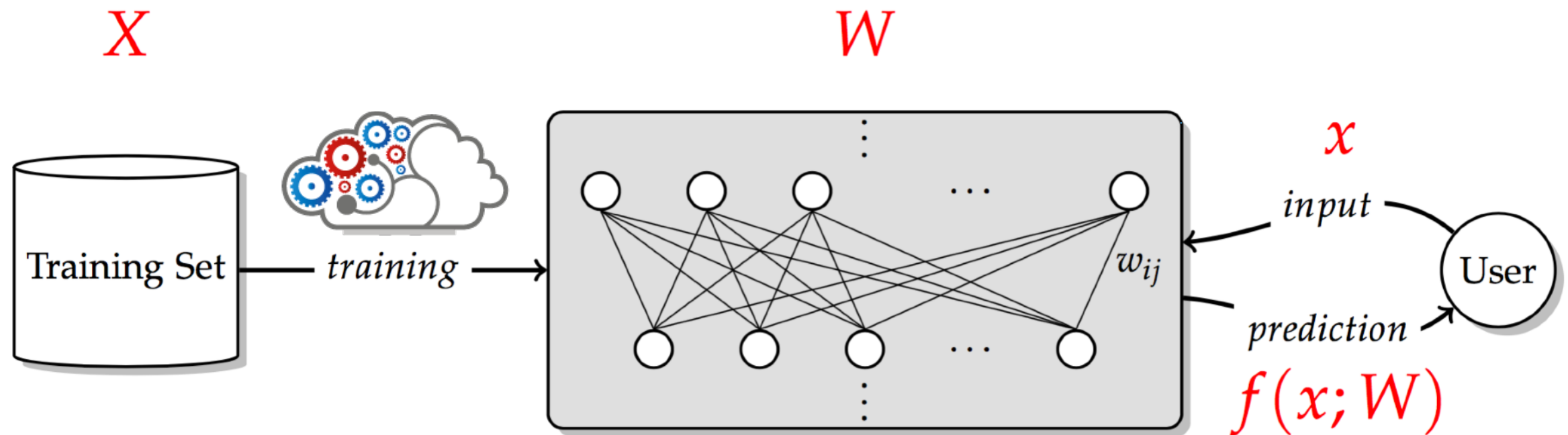
Direct and intentional leakage

- Unauthorized access to data, and data breaches
- Massive data collection

Indirect and unintentional leakage

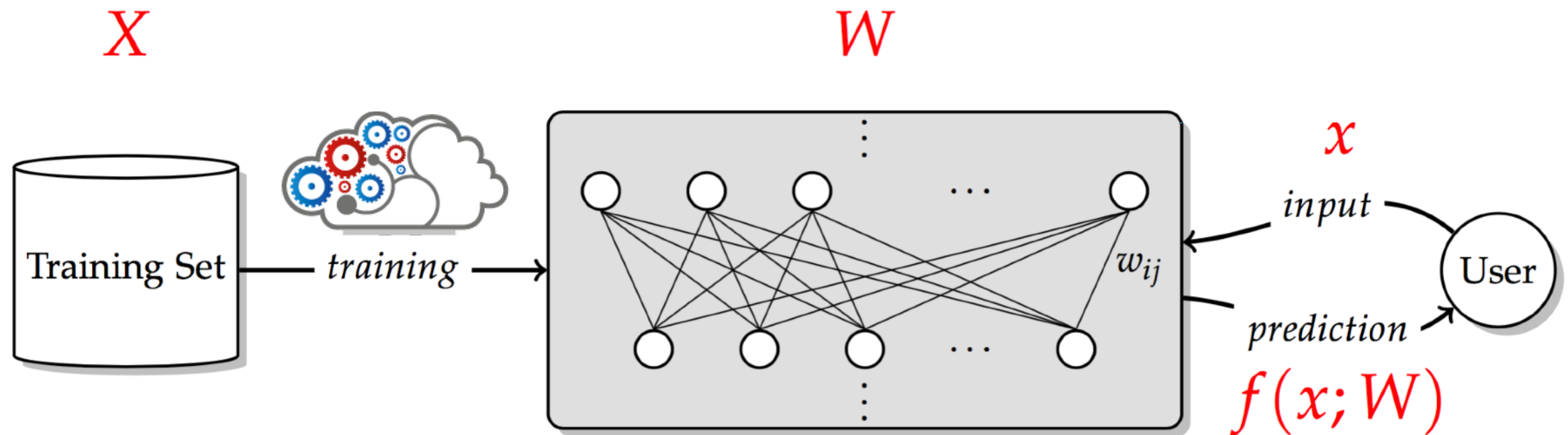
- Meta-data: Data about data
- Data correlated with data
- Computations on data

Privacy Risks in Machine Learning

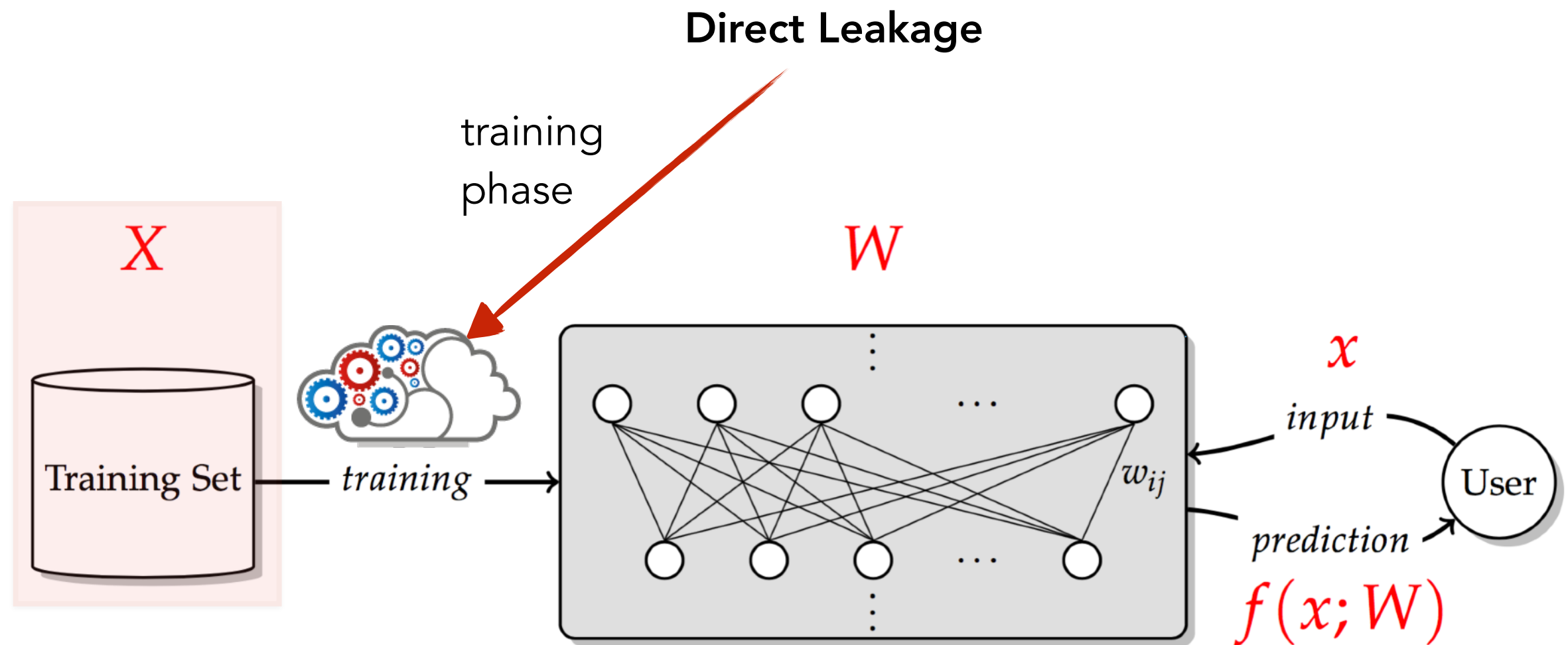


Privacy Risks in Machine Learning

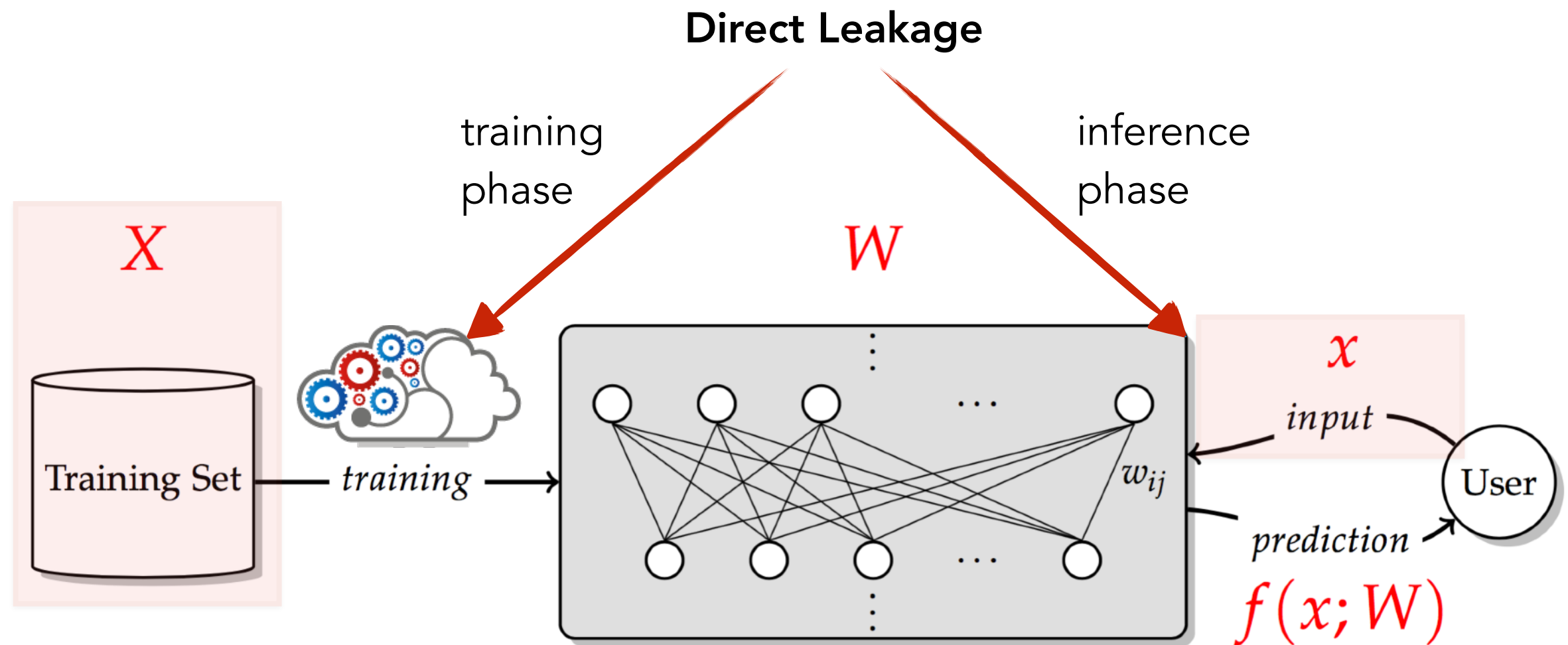
Direct Leakage



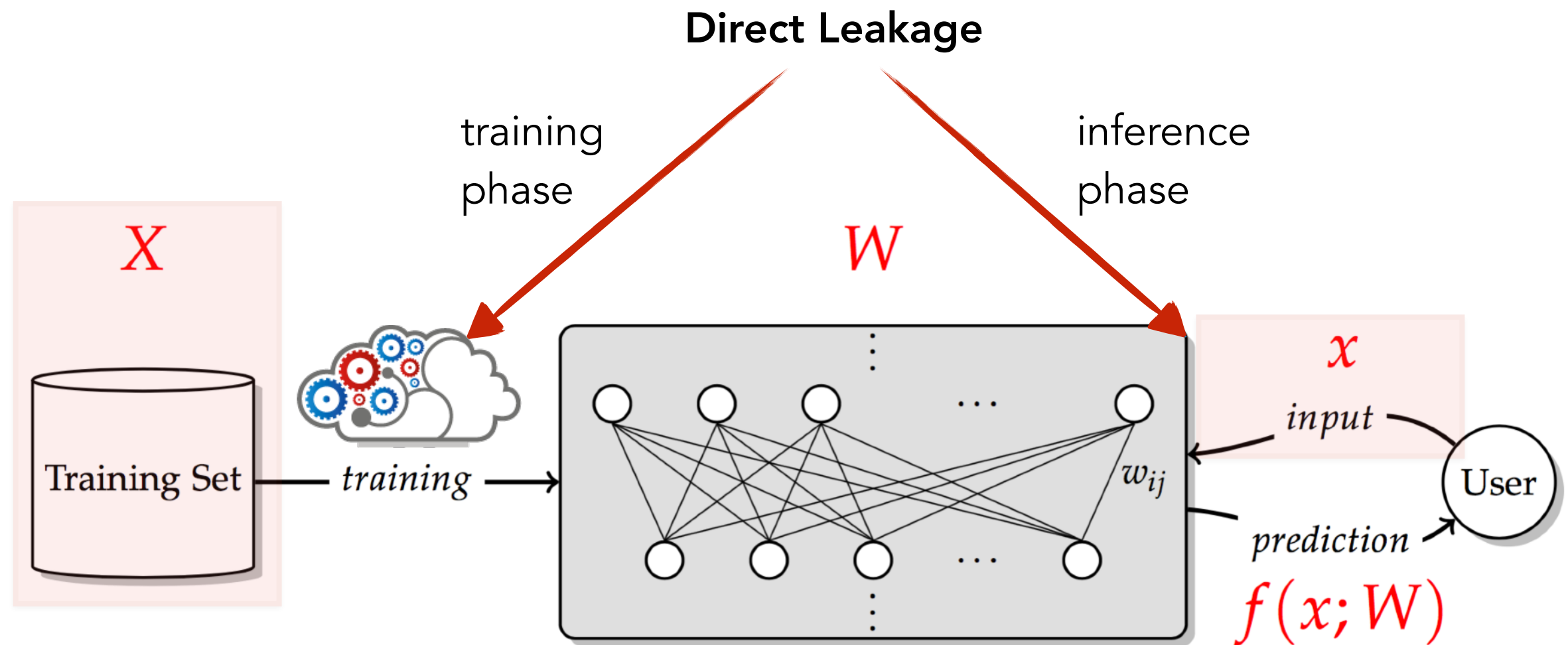
Privacy Risks in Machine Learning



Privacy Risks in Machine Learning

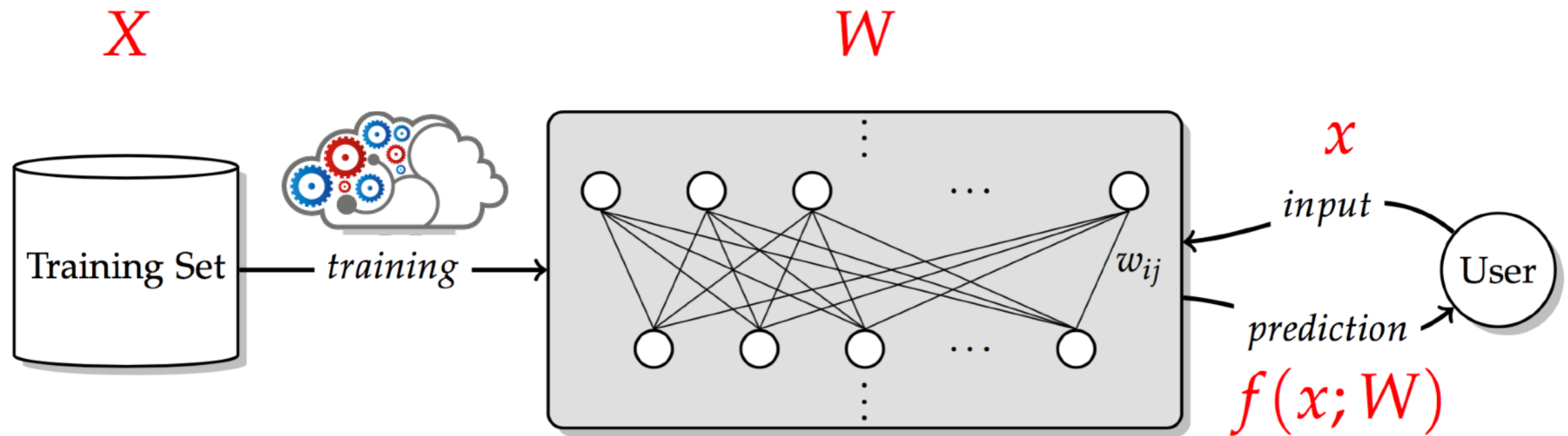


Privacy Risks in Machine Learning

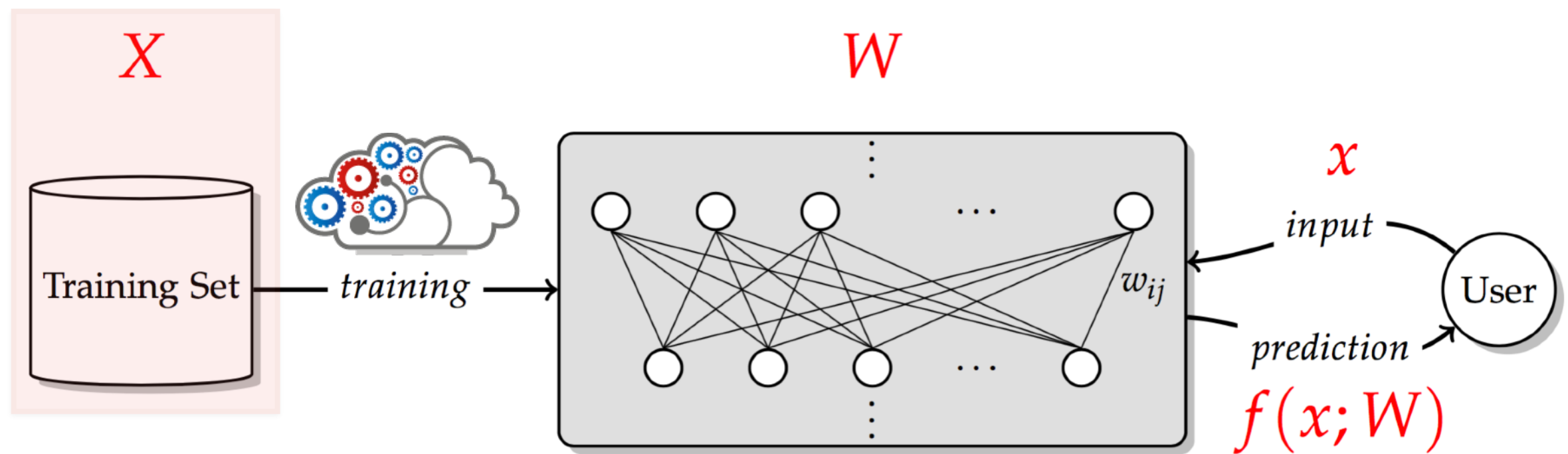


How to prevent leakage? Secure multi-party computation, homomorphic encryption, trusted hardware, ...

Privacy Risks in Machine Learning

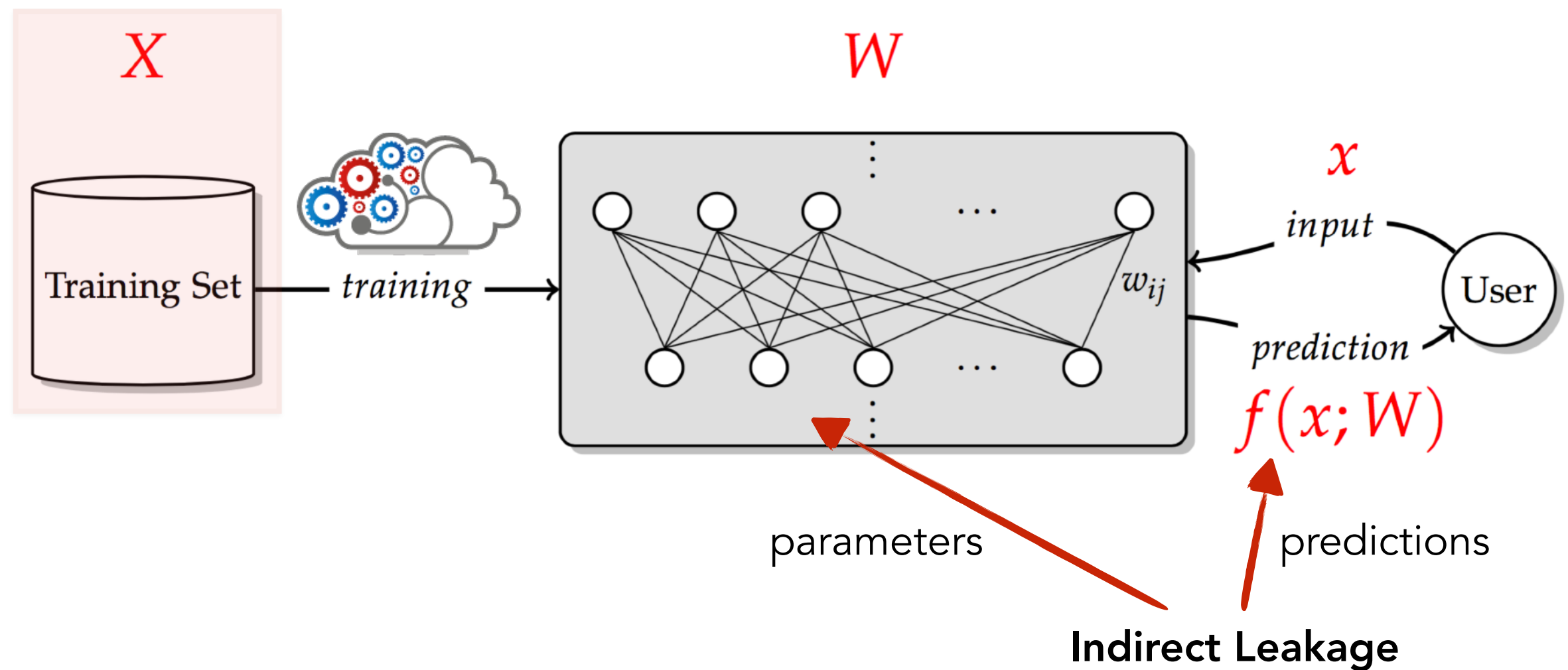


Privacy Risks in Machine Learning



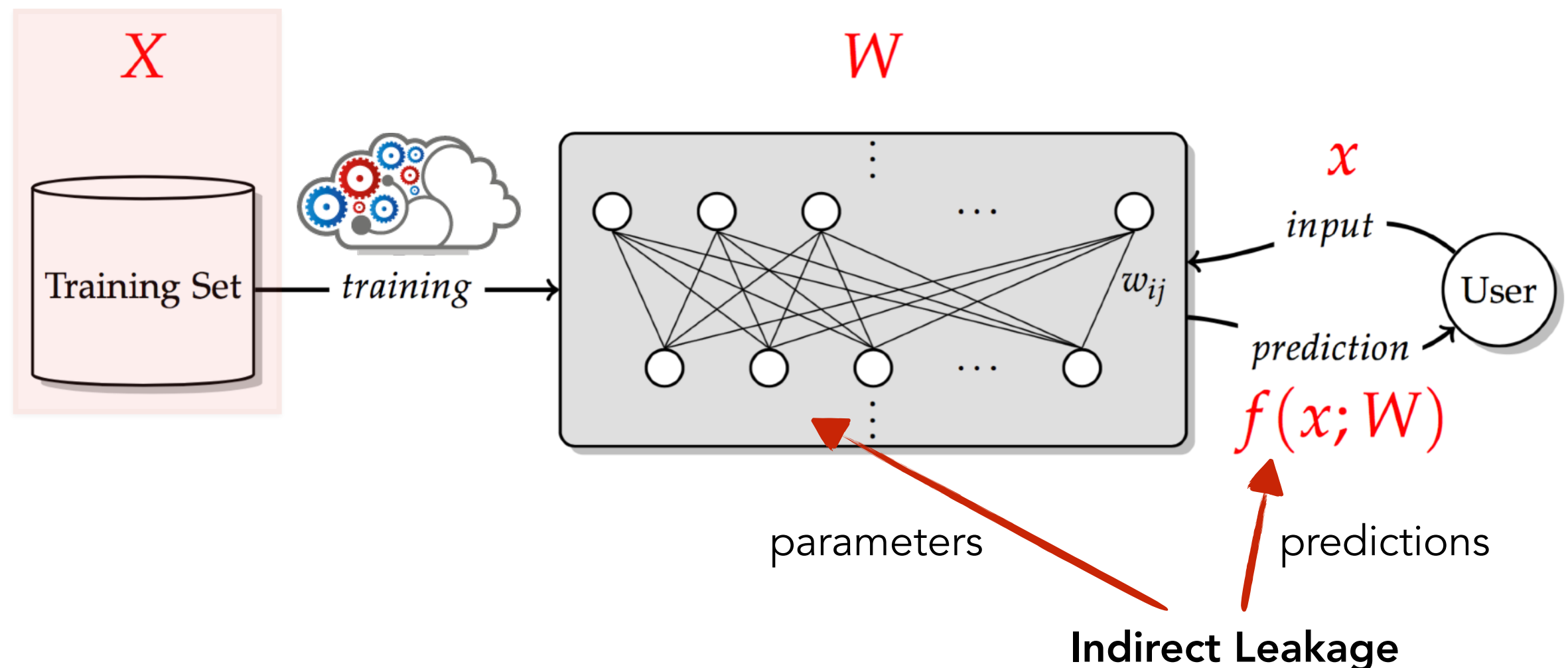
Indirect Leakage

Privacy Risks in Machine Learning



Privacy Risks in Machine Learning

What is leakage? Inferring information about members of X , beyond what can be learned about its underlying distribution

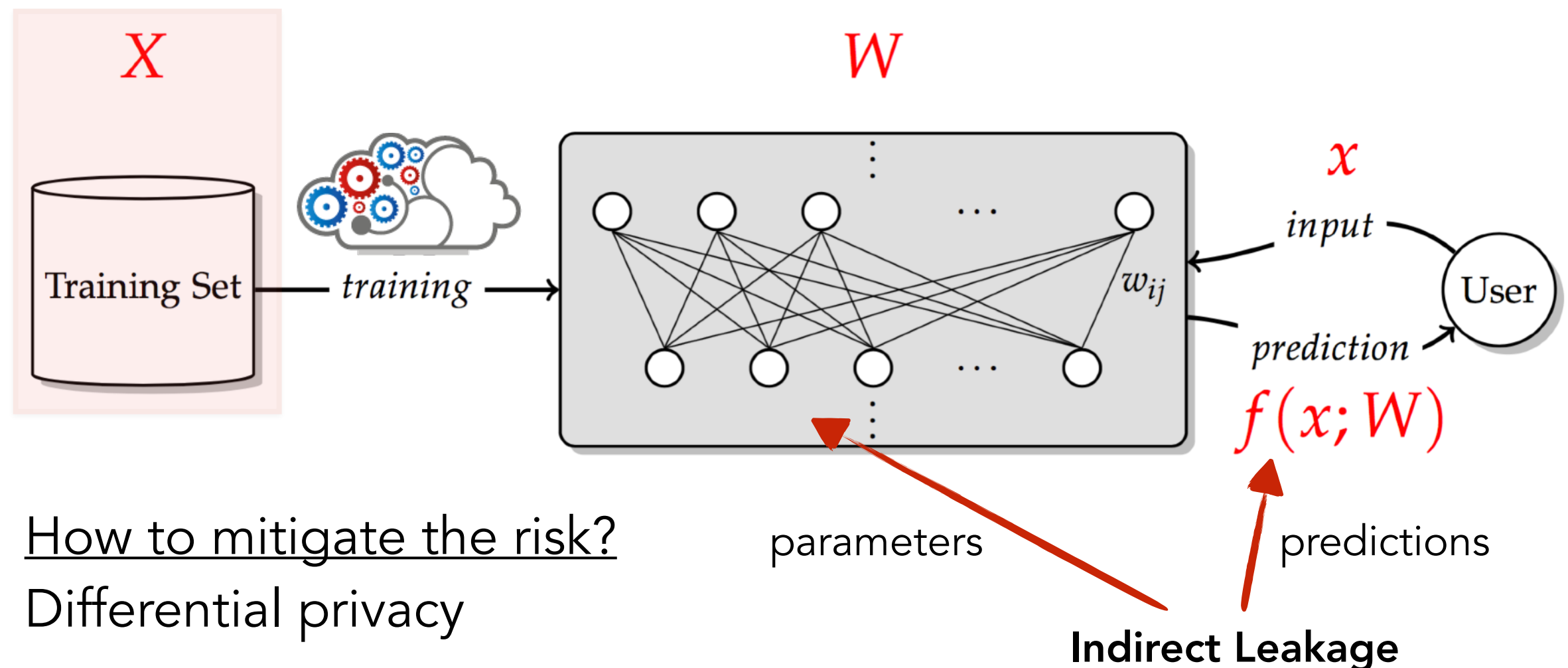


[Shokri, Stronati, Song, Shmatikov] Membership Inference Attacks against Machine Learning Models, SP'17

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

Privacy Risks in Machine Learning

What is leakage? Inferring information about members of X , beyond what can be learned about its underlying distribution

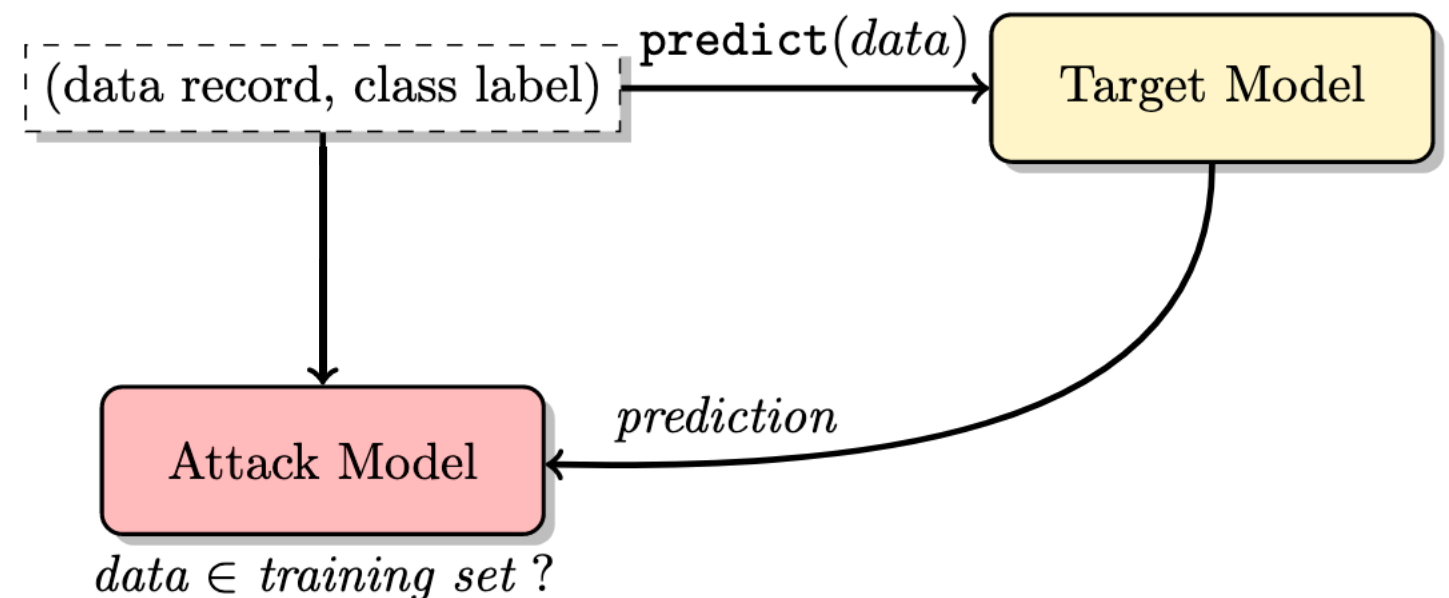


How to mitigate the risk?
Differential privacy

[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

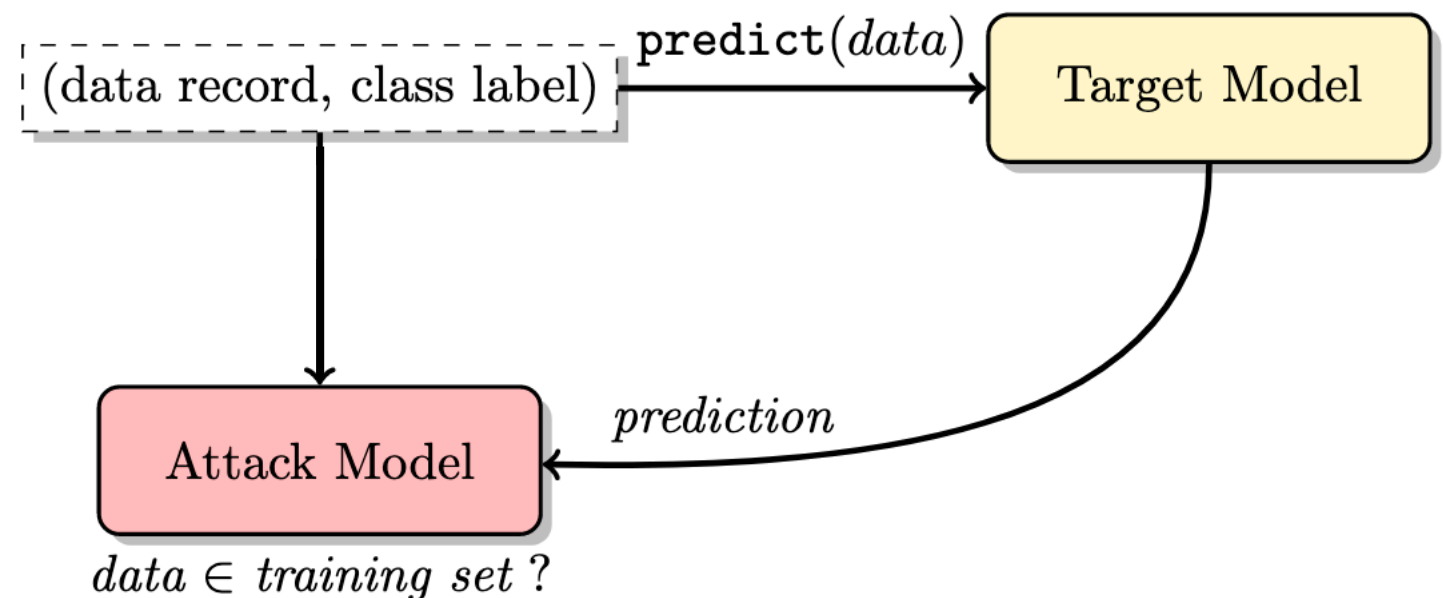
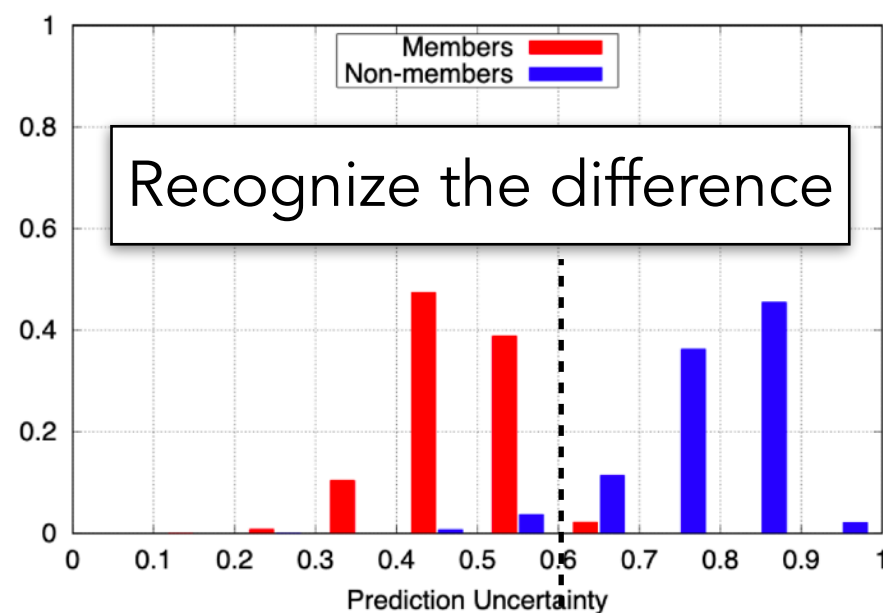
How to Quantify the Leakage?

- Indistinguishability game: Can an adversary distinguish between two models that are trained on two neighboring datasets (one includes an extra data point x)?
- Membership inference: Given a model, can an adversary infer whether data point x is part of its training set?



How to Quantify the Leakage?

- Indistinguishability game: Can an adversary distinguish between two models that are trained on two neighboring datasets (one includes an extra data point x)?
- Membership inference: Given a model, can an adversary infer whether data point x is part of its training set?

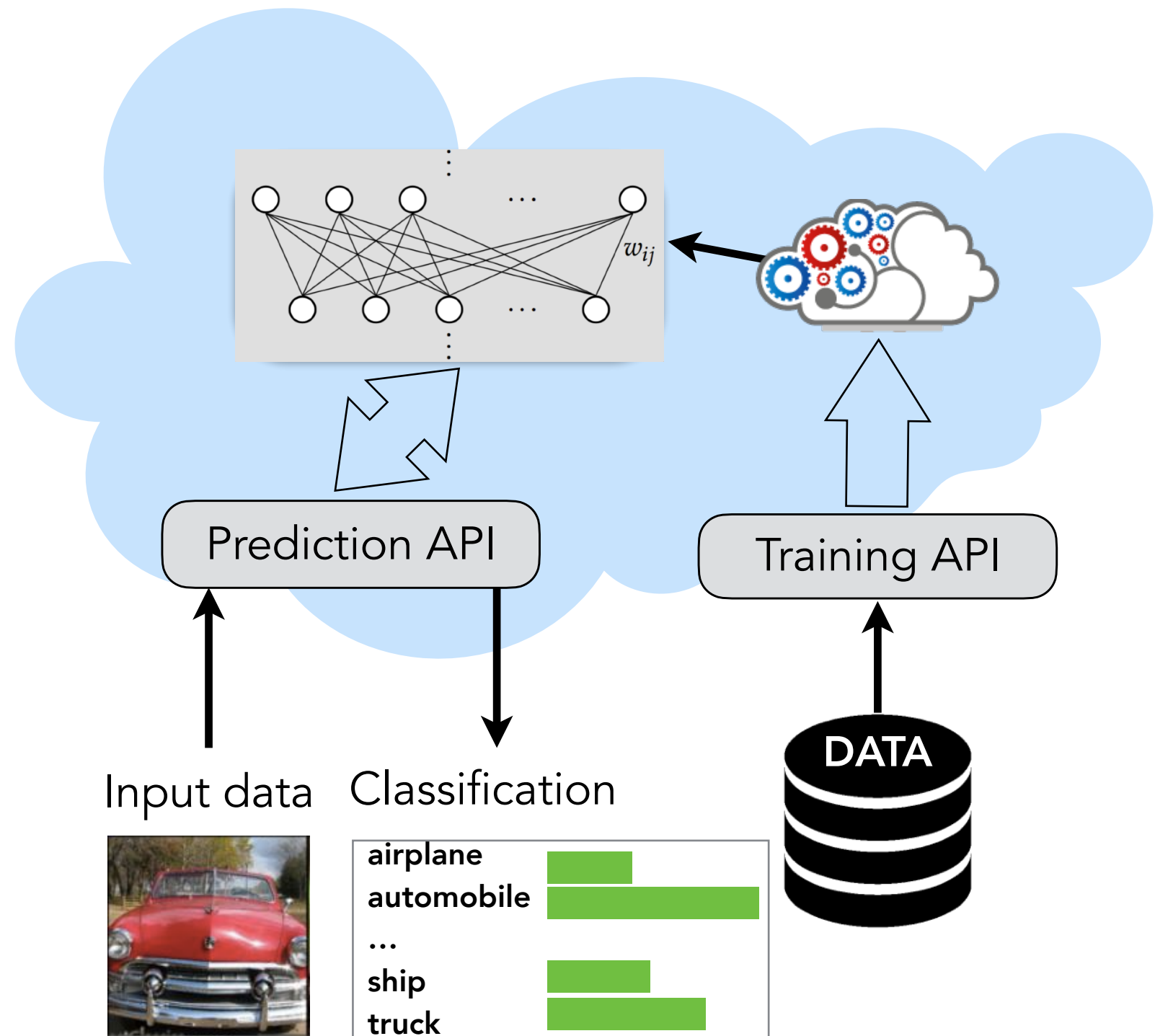


Membership Inference Attacks against Classification Models

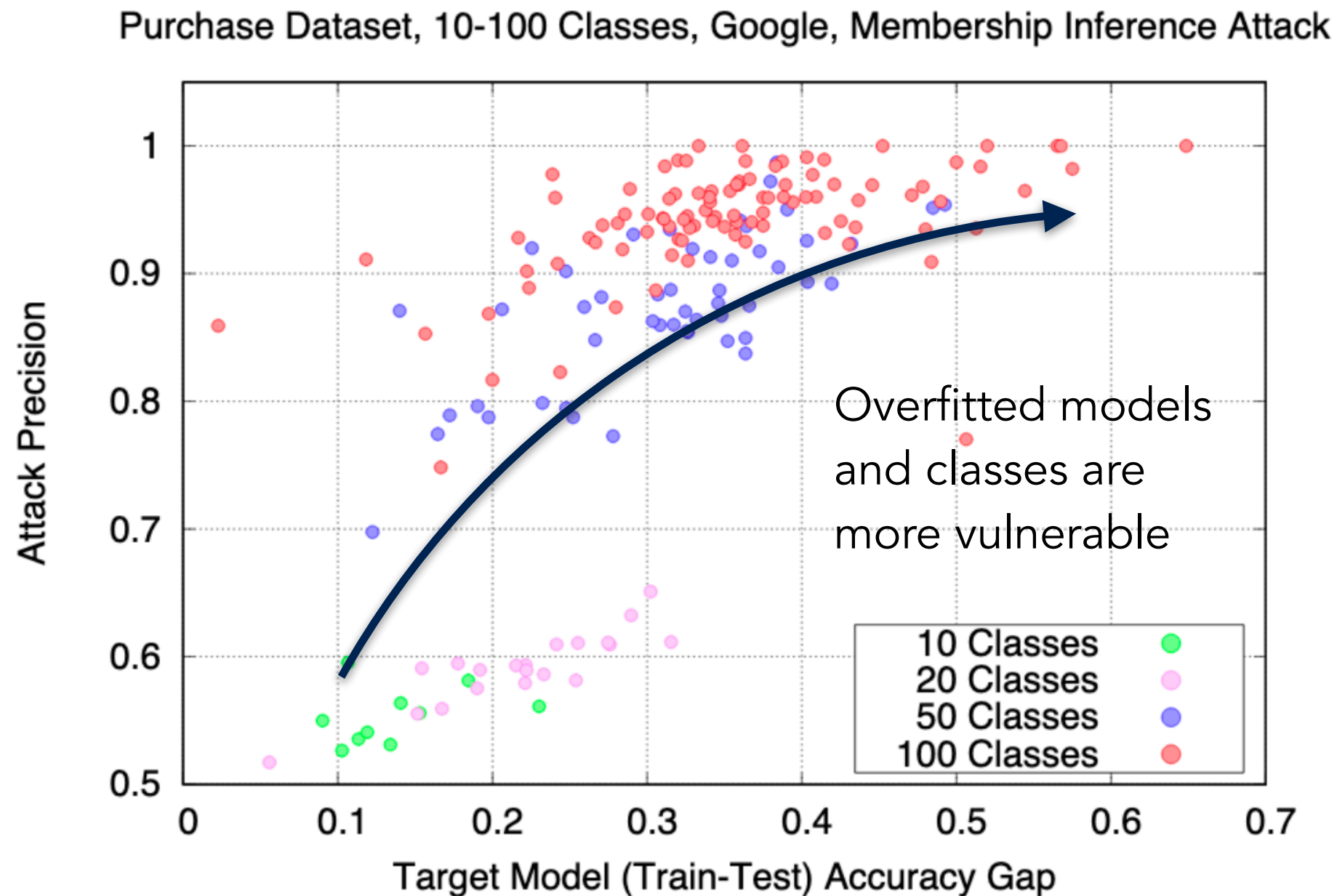
Machine Learning
as a Service



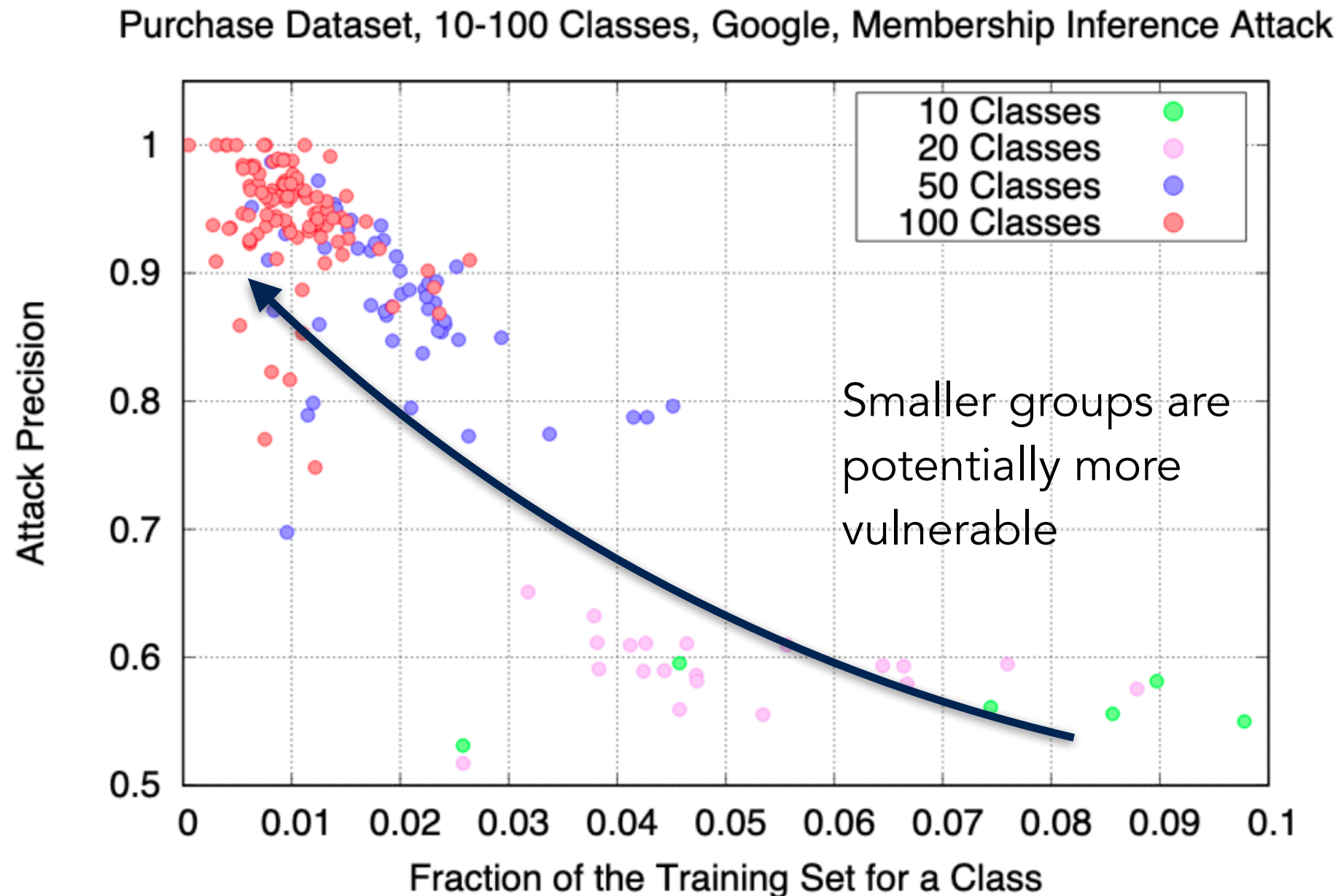
**Membership
Inference Attack**
Accuracy:
~ 90%



Privacy Leakage due to Overfitting



Disparate Privacy Vulnerability



White-box Privacy Analysis

- Leakage through parameters (white-box) vs. predictions (black-box)

Most accurate pre-trained models				Mem inference attack accuracy		
Dataset	Architecture	Train Accuracy	Test Accuracy	Black-box	White-box (Outputs)	White-box (Gradients)
CIFAR100	Alexnet	99%	44%			
CIFAR100	ResNet	89%	73%			
CIFAR100	DenseNet	100%	82%			

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

White-box Privacy Analysis

- Leakage through parameters (white-box) vs. predictions (black-box)

Most accurate pre-trained models				Mem inference attack accuracy		
Dataset	Architecture	Train Accuracy	Test Accuracy	Black-box	White-box (Outputs)	White-box (Gradients)
CIFAR100	Alexnet	99%	44%			
CIFAR100	ResNet	89%	73%			
CIFAR100	DenseNet	100%	82%			

High generalizability
to test data

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

White-box Privacy Analysis

- Leakage through parameters (white-box) vs. predictions (black-box)

Most accurate pre-trained models				Mem inference attack accuracy		
Dataset	Architecture	Train Accuracy	Test Accuracy	Black-box	White-box (Outputs)	White-box (Gradients)
CIFAR100	Alexnet	99%	44%	74.2%	74.6%	75.1%
CIFAR100	ResNet	89%	73%	62.2%	62.2%	64.3%
CIFAR100	DenseNet	100%	82%	67.7%	67.7%	74.3%

High generalizability
to test data

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

White-box Privacy Analysis

- Leakage through parameters (white-box) vs. predictions (black-box)

Most accurate pre-trained models				Mem inference attack accuracy		
Dataset	Architecture	Train Accuracy	Test Accuracy	Black-box	White-box (Outputs)	White-box (Gradients)
CIFAR100	Alexnet	99%	44%	74.2%	74.6%	75.1%
CIFAR100	ResNet	89%	73%	62.2%	62.2%	64.3%
CIFAR100	DenseNet	100%	82%	67.7%	67.7%	74.3%

High generalizability
to test data

Low privacy
(Significant leakage
through parameters)

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

White-box Privacy Analysis

- Leakage through parameters (white-box) vs. predictions (black-box)

Most accurate pre-trained models				Mem inference attack accuracy		
Dataset	Architecture	Train Accuracy	Test Accuracy	Black-box	White-box (Outputs)	White-box (Gradients)
CIFAR100	Alexnet	99%	44%	74.2%	74.6%	75.1%
CIFAR100	ResNet	89%	73%	62.2%	62.2%	64.3%
CIFAR100	DenseNet	100%	82%	67.7%	67.7%	74.3%

Large
capacity

High generalizability
to test data

Low privacy
(Significant leakage
through parameters)

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

White-box Privacy Analysis

- Leakage through parameters (white-box) vs. predictions (black-box)

Most accurate pre-trained models				Mem inference attack accuracy		
Dataset	Architecture	Train Accuracy	Test Accuracy	Black-box	White-box (Outputs)	White-box (Gradients)
CIFAR100	Alexnet	99%	44%	74.2%	74.6%	75.1%
CIFAR100	ResNet	89%	73%	62.2%	62.2%	64.3%
CIFAR100	DenseNet	100%	82%	67.7%	67.7%	74.3%

Large
capacity

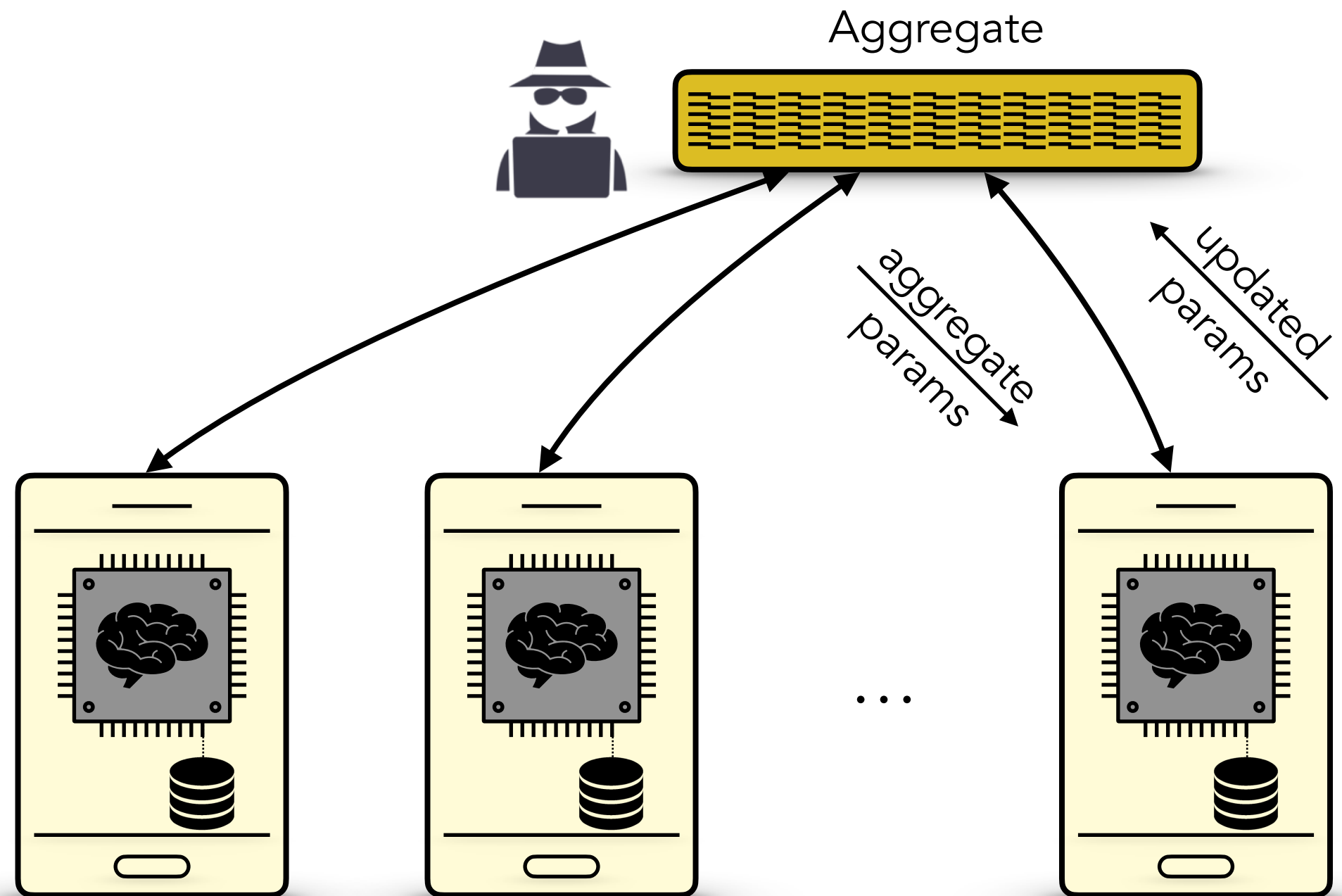
High generalizability
to test data

Low privacy
(Significant leakage
through parameters)

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

[Feldman] Does Learning Require **Memorization**? A Short Tale about a Long Tail, STOC'20

Decentralized (Federated) Learning



[Shokri and Shmatikov] Privacy-Preserving Deep Learning, CCS'15

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

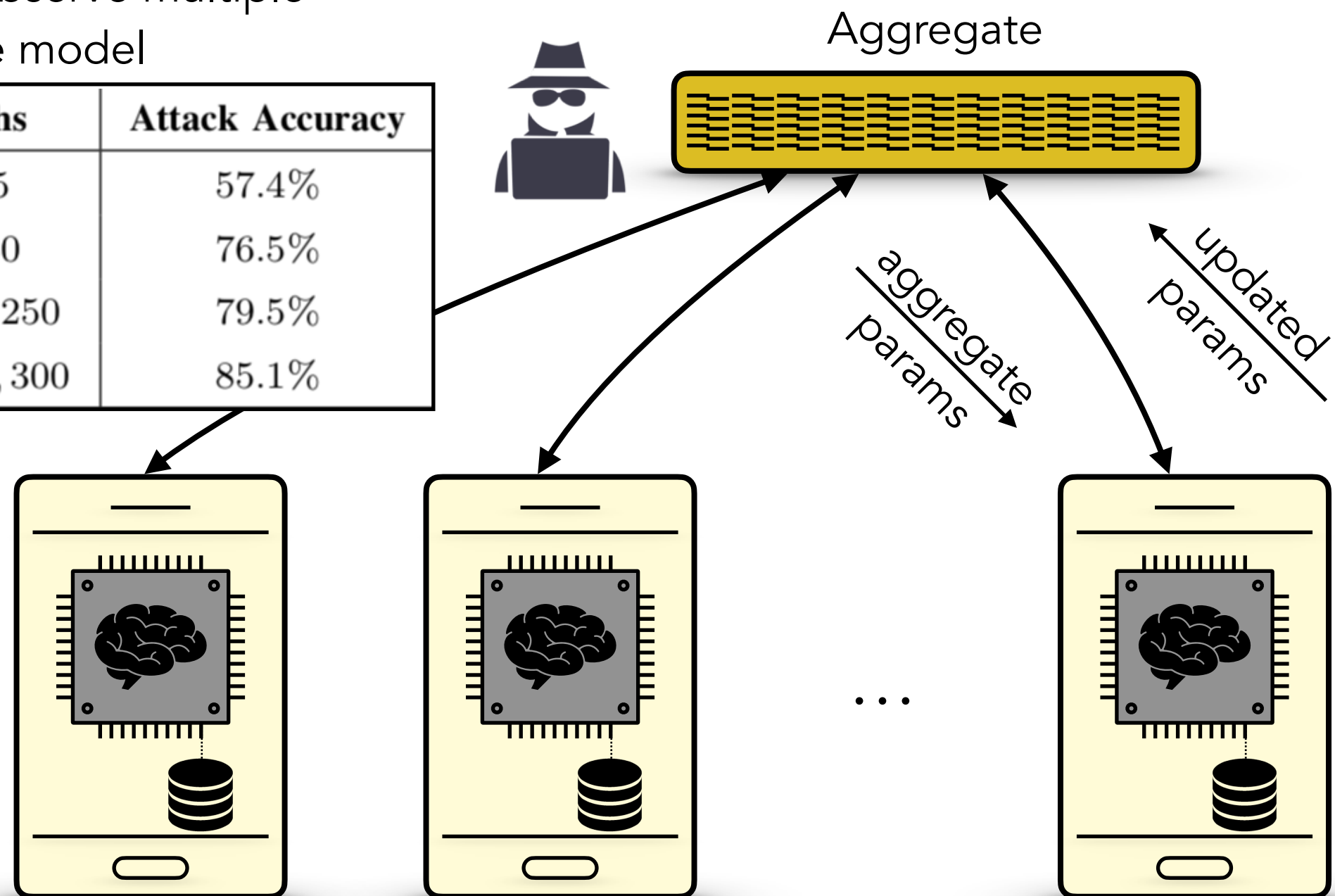
[Melis, Song, De Cristofaro, Shmatikov] Exploiting Unintended Feature Leakage in Collaborative Learning, SP'19

Decentralized (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



[Shokri and Shmatikov] Privacy-Preserving Deep Learning, CCS'15

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

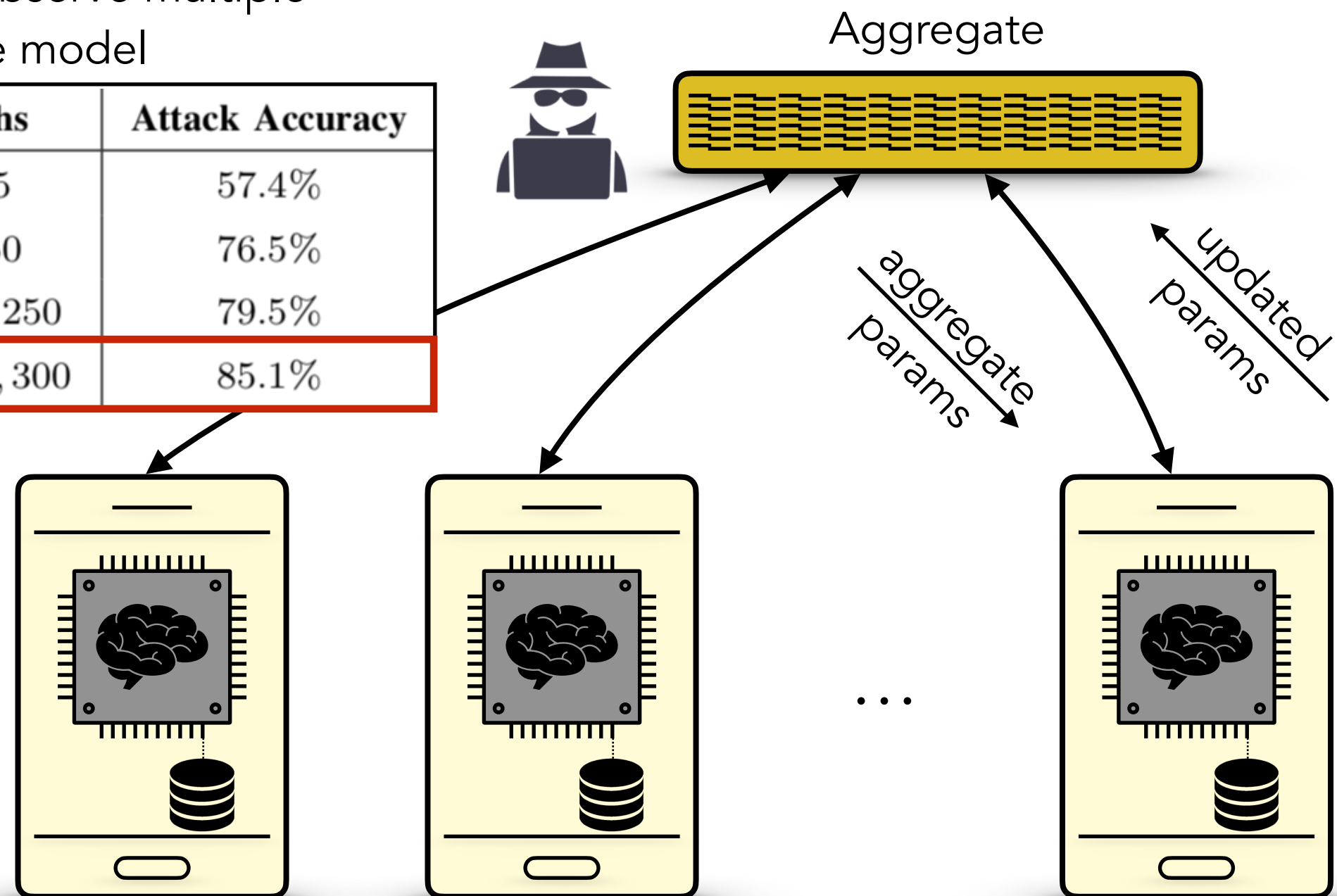
[Melis, Song, De Cristofaro, Shmatikov] Exploiting Unintended Feature Leakage in Collaborative Learning, SP'19

Decentralized (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

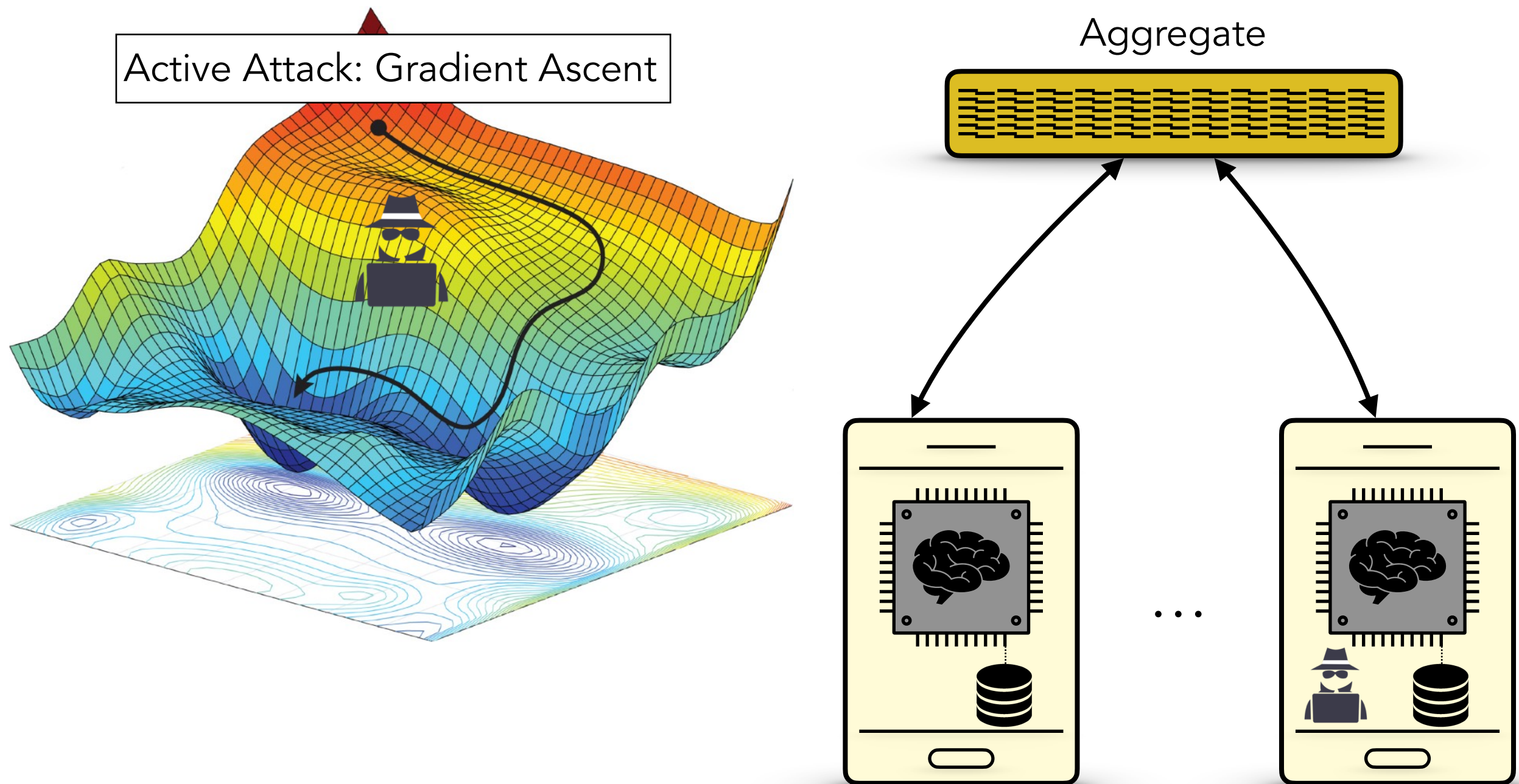


[Shokri and Shmatikov] Privacy-Preserving Deep Learning, CCS'15

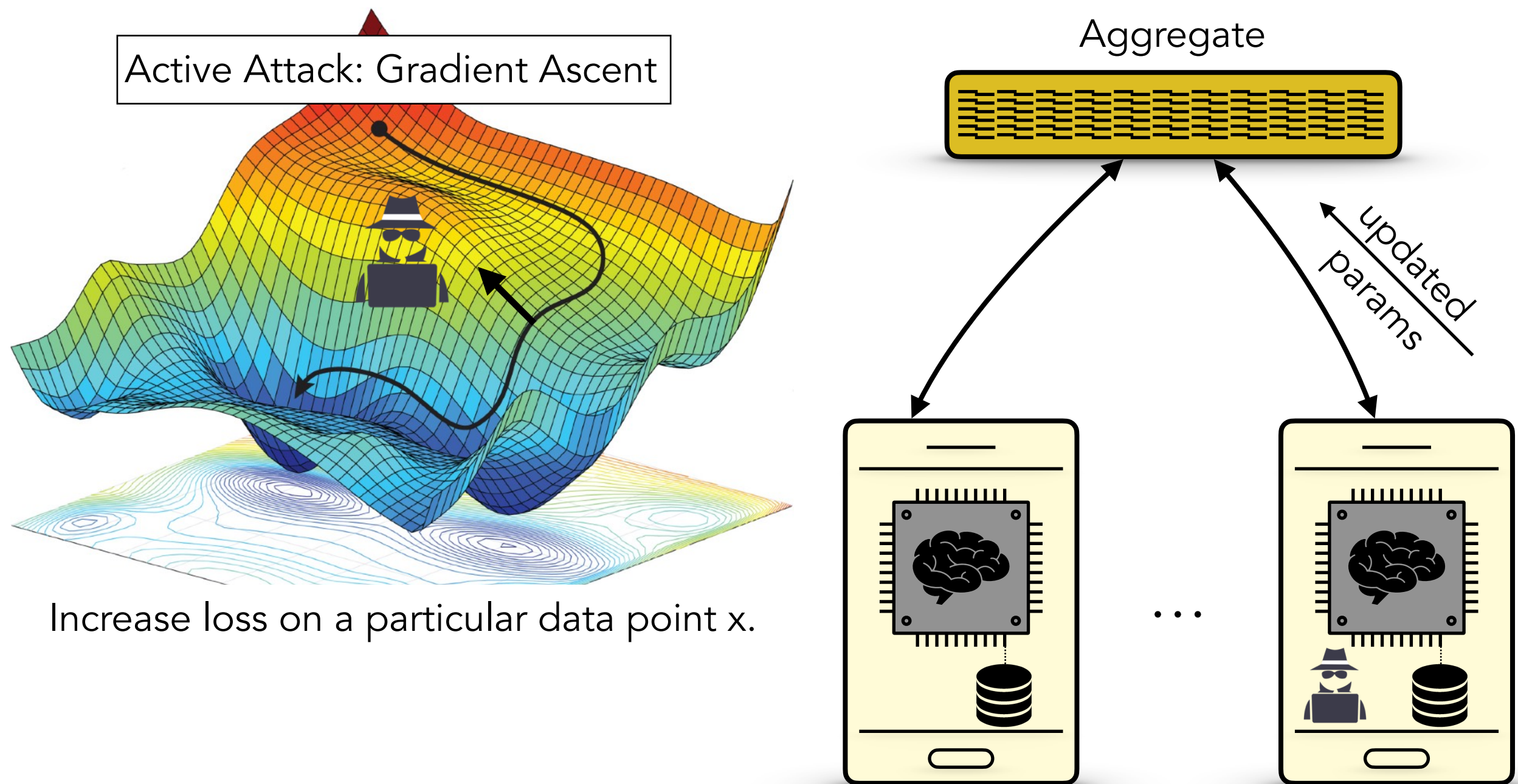
[Nasr, Shokri, Houmansadr] Comprehensive Privacy Analysis of Deep Learning: Passive and Active White-box Inference Attacks against Centralized and Federated Learning, SP'19

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

Decentralized (Federated) Learning

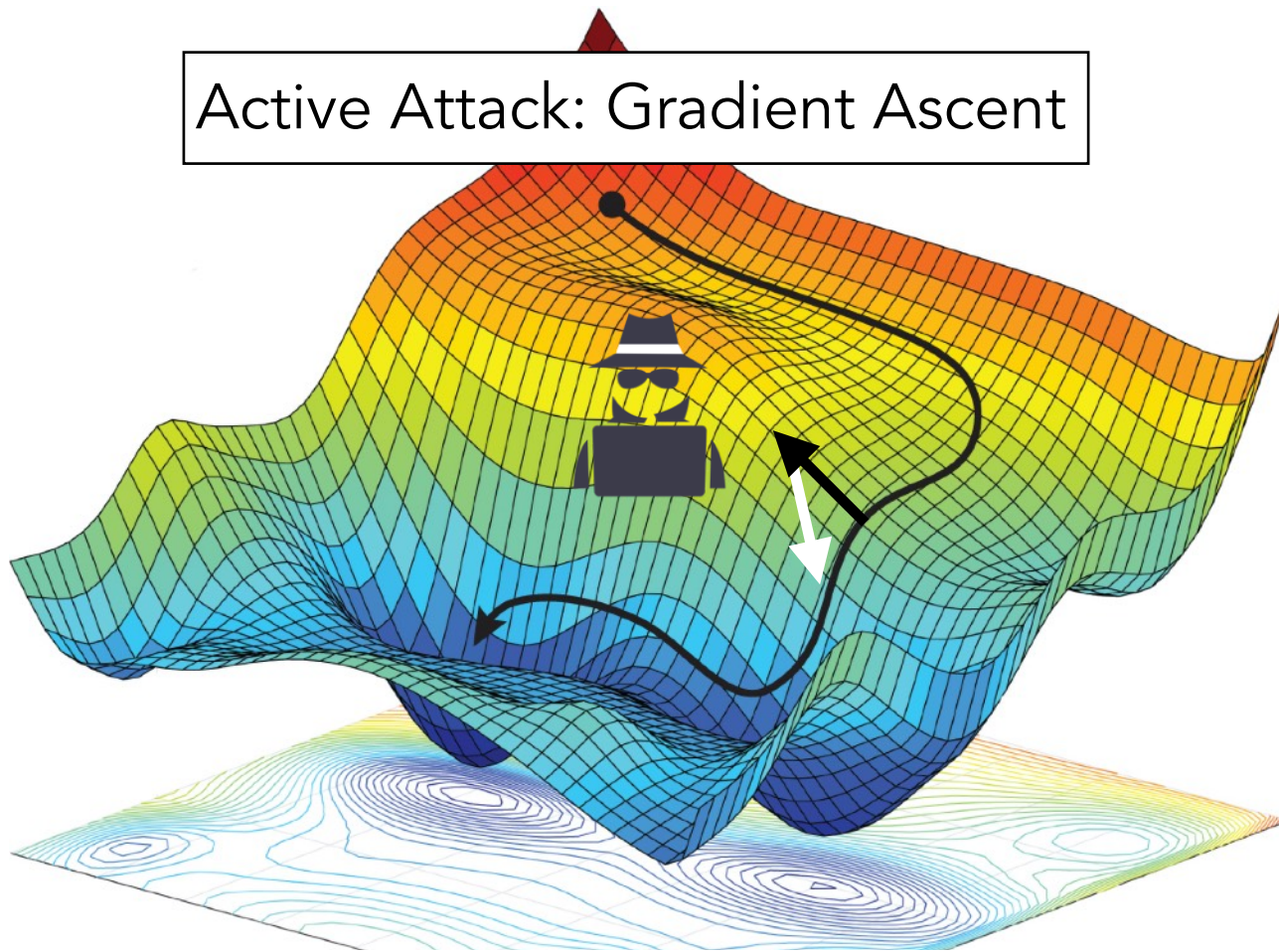


Decentralized (Federated) Learning



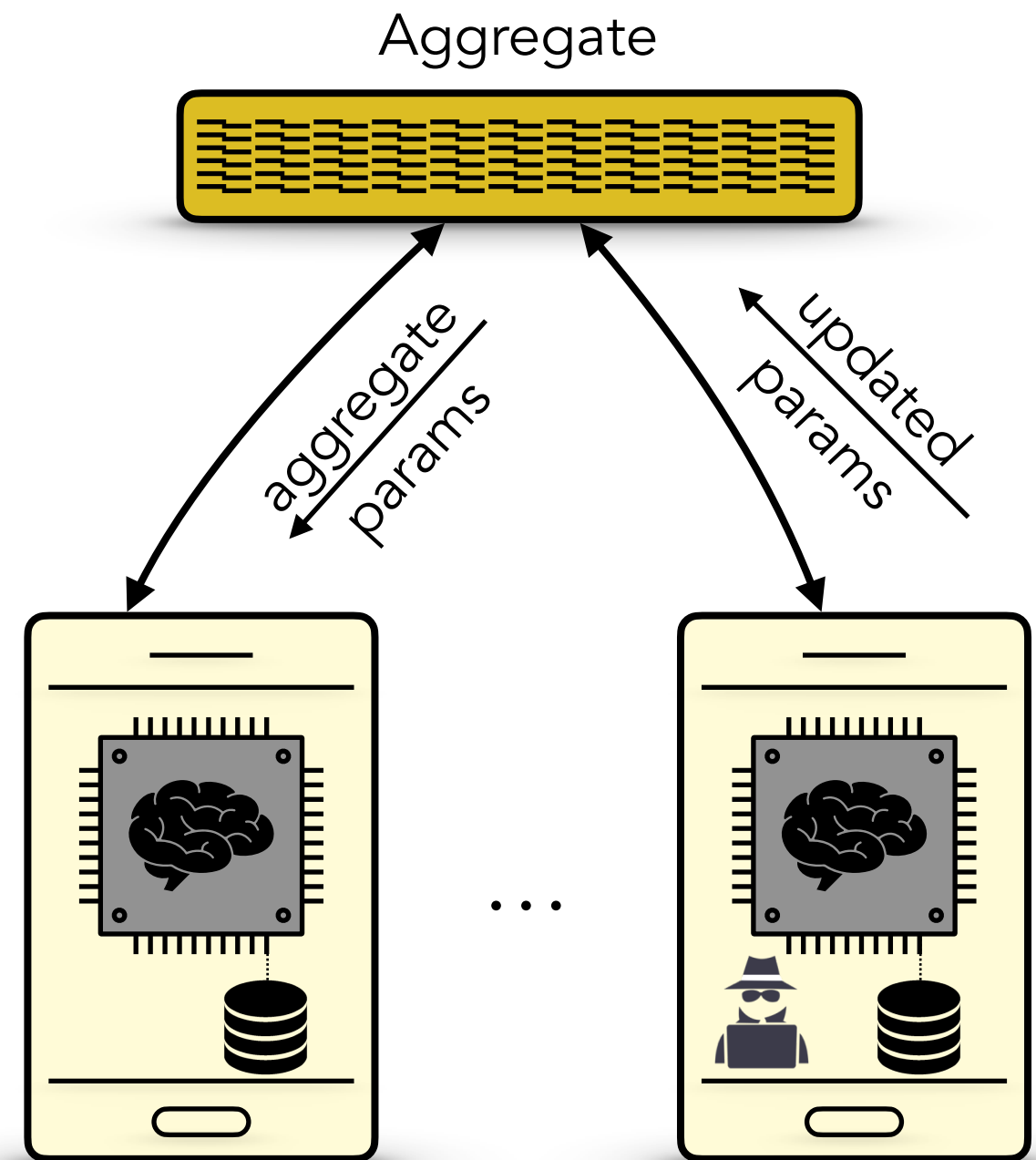
Decentralized (Federated) Learning

Active Attack: Gradient Ascent





Increase loss on a particular data point x .

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



AI Regulations - Data Protection

- "... membership inferences show that AI models can inadvertently contain personal data"  Information Commissioner's Office
- "Attacks that reveal confidential information about the data include membership inference whereby ..."  National Institute of Standards and Technology
U.S. Department of Commerce
- "... ensuring that privacy and personal data are adequately protected during the use of AI"
- "... ensuring that AI systems are resilient to overt attacks and subtle attacks that manipulate data or algorithms...."
- "...should consider the risks to data throughout the design, development, and operation of an AI system"

On Artificial Intelligence - A European Approach to excellence and trust - Feb 2020

The White House Memo on Guidance for Regulation of Artificial Intelligence Applications - Jan 2020

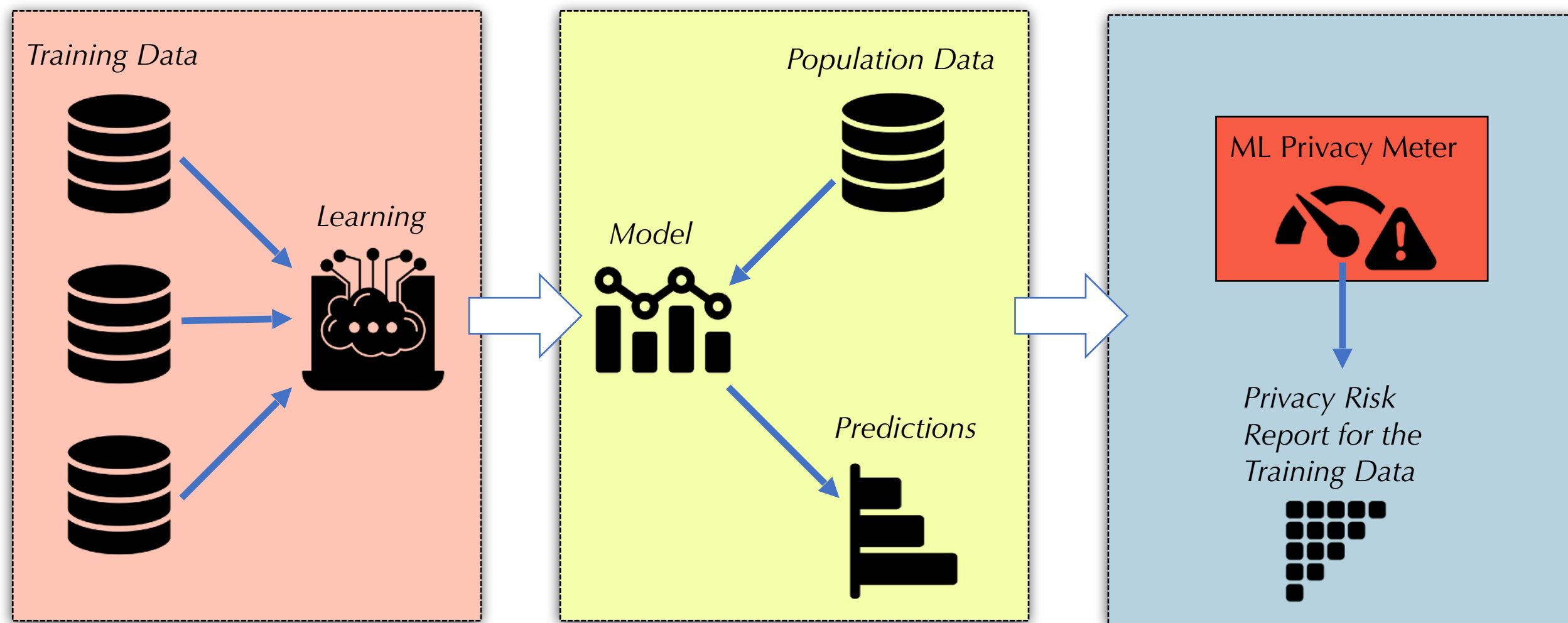
Guidance on the AI auditing framework Draft guidance for consultation. Information Commissioner's Office

A Taxonomy and Terminology of Adversarial Machine Learning. Draft NISTIR 8269

Data Protection Impact Assessment



Tool: ML Privacy Meter

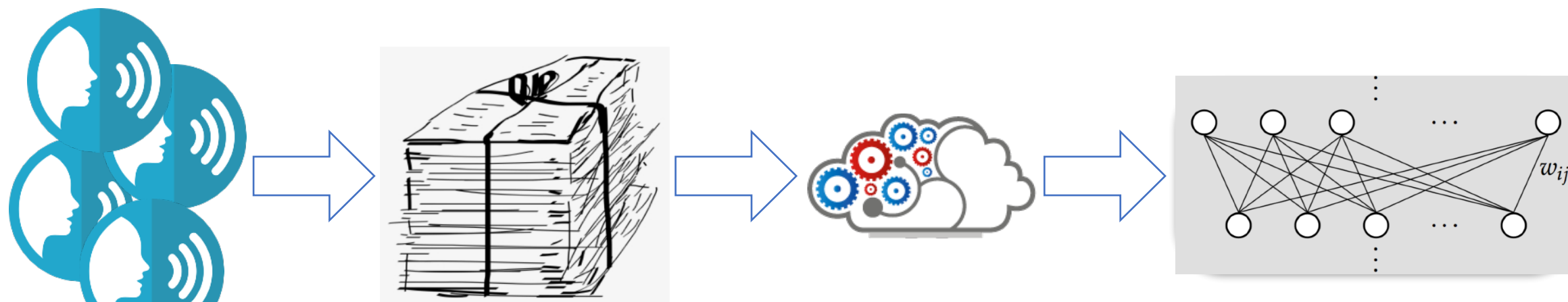


ML Privacy Meter is a Python library (`ml_privacy_meter`) that enables quantifying the privacy risks of machine learning models. https://github.com/privacytrustlab/ml_privacy_meter

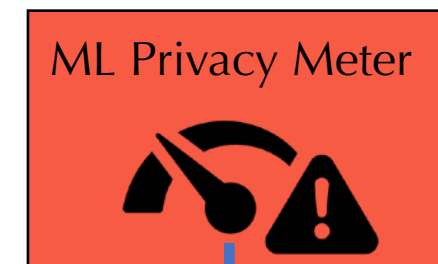


ML Privacy Meter

Example: NLP Models



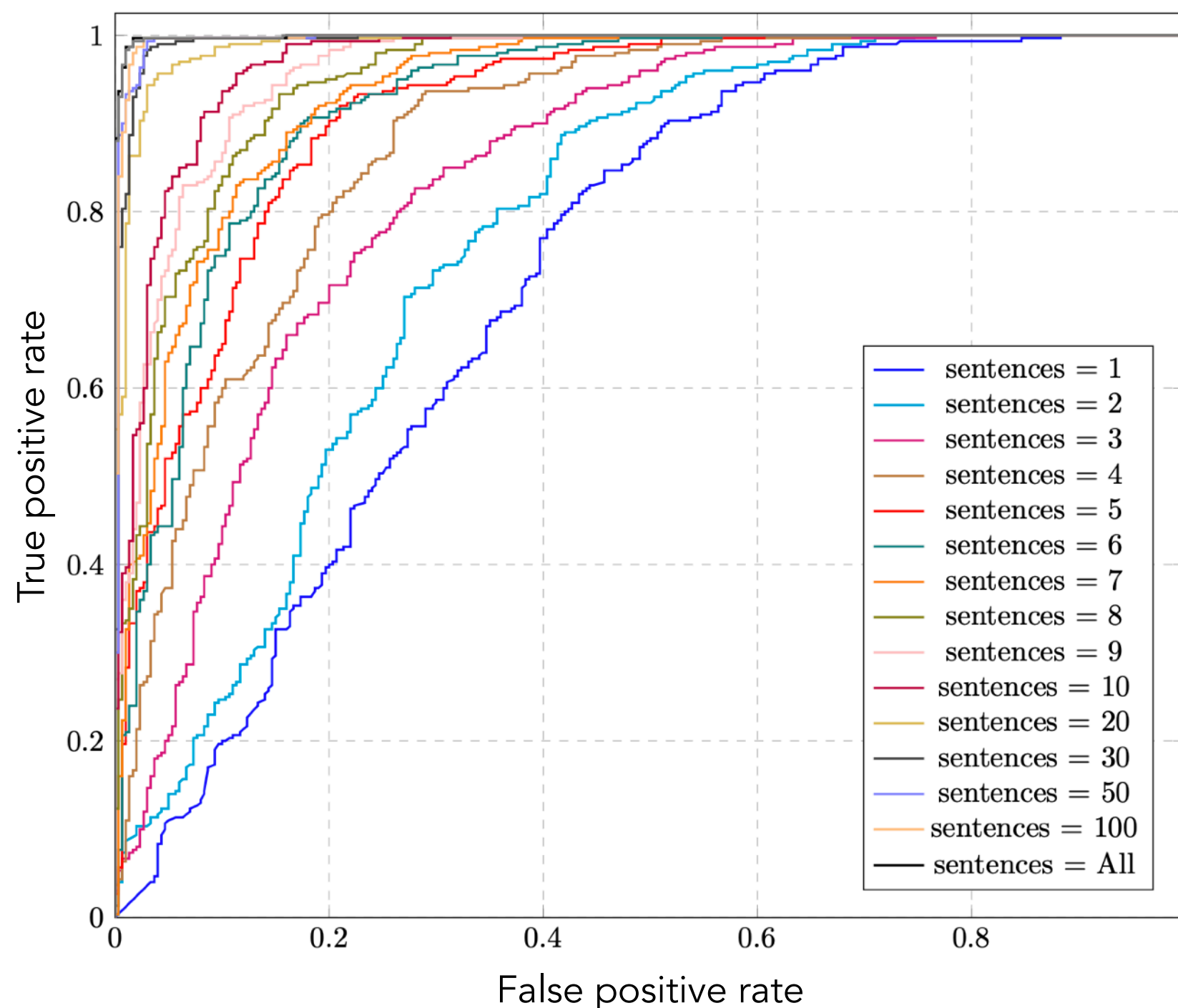
- How much does the model leak about the sentences of a particular author/speaker? What about the membership of the author in the training set (based on known samples)?
- Which samples are leaked?



*Privacy Risk
Report for the
Training Data*

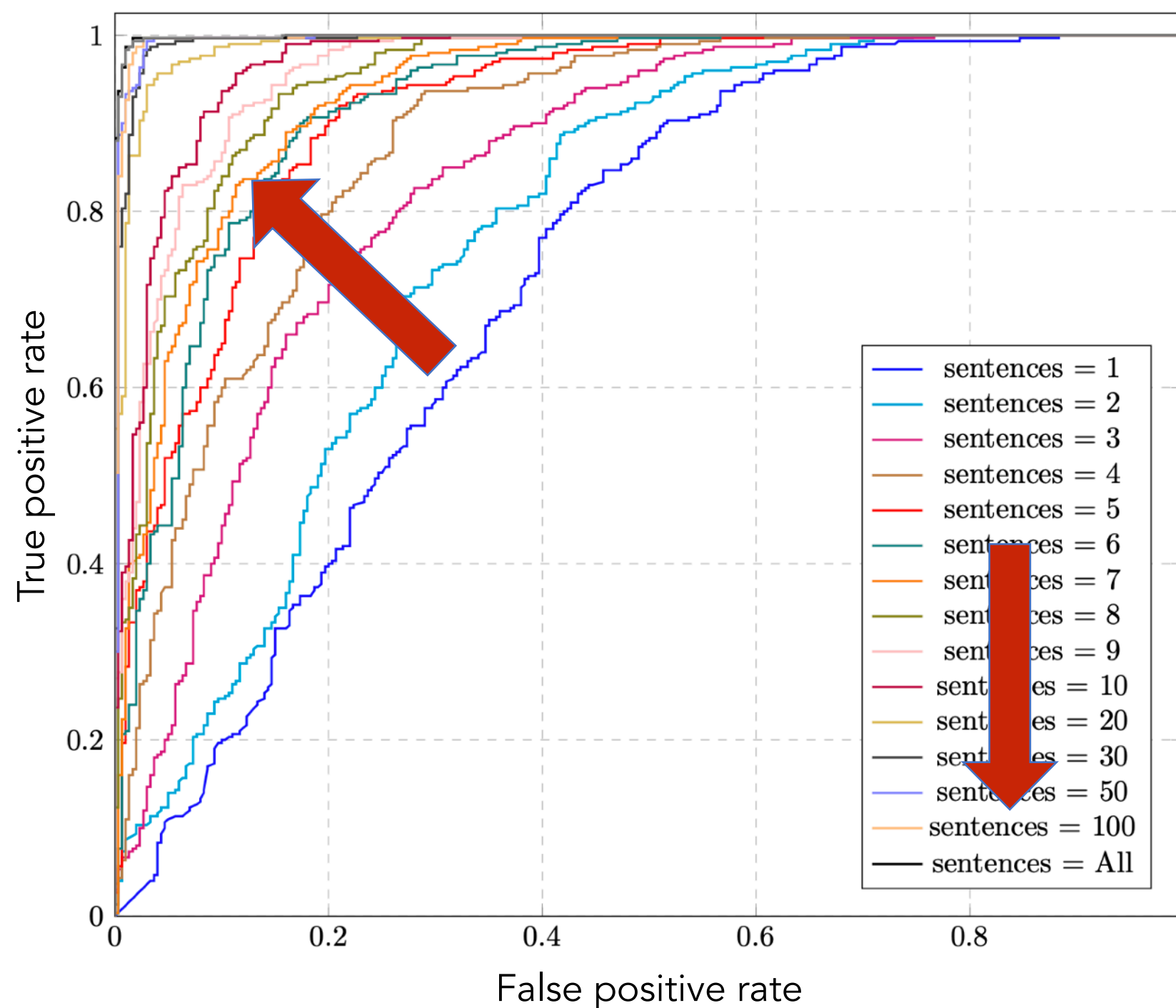


Membership Inference



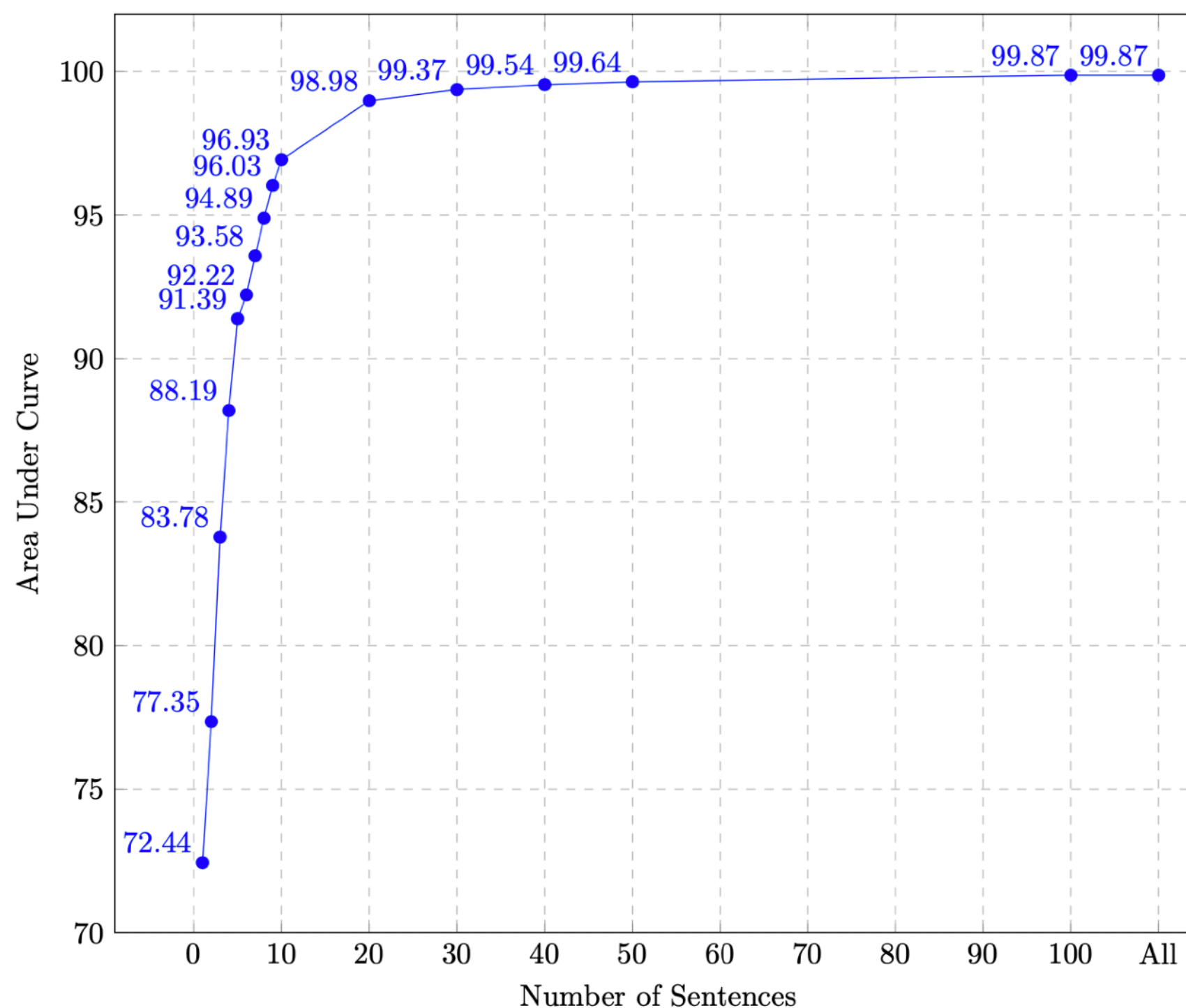
SATED (Speaker
Annotated TED
talks) dataset

Membership Inference



SATED (Speaker
Annotated TED
talks) dataset

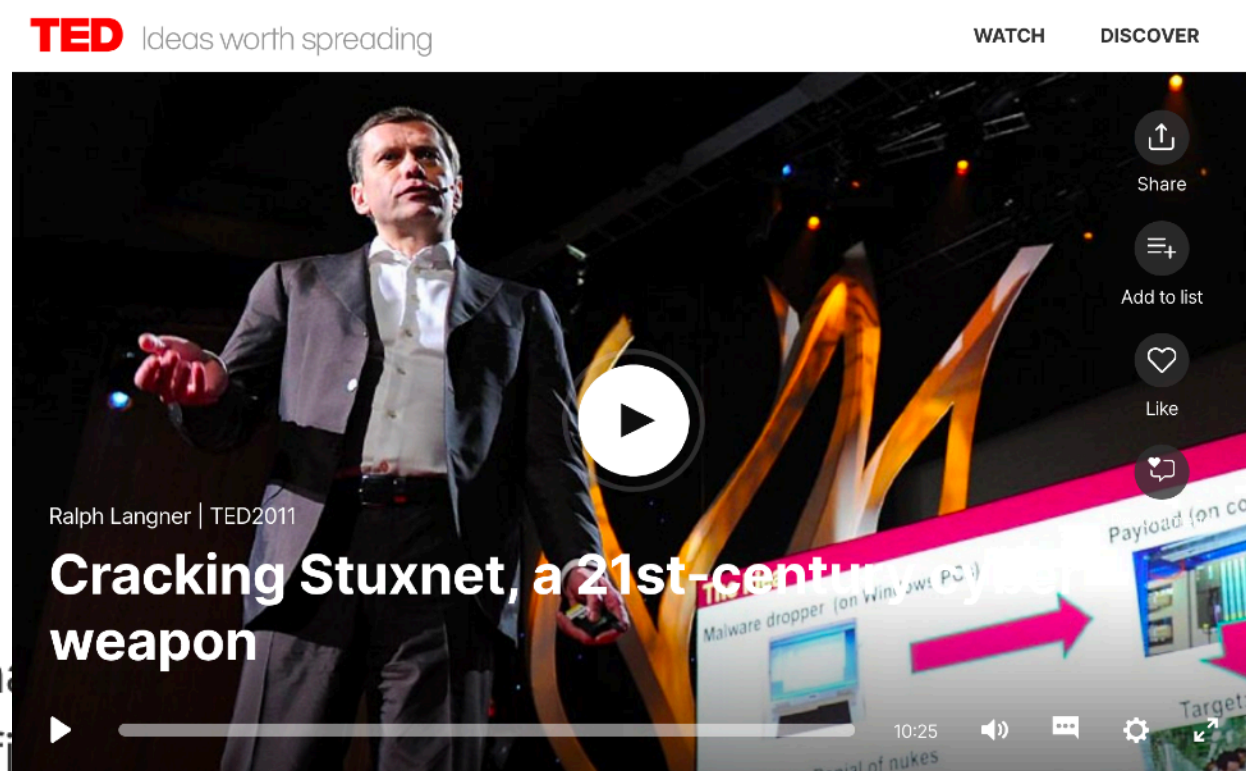
Membership Inference



SATED (Speaker Annotated TED talks) dataset

[Maddi] https://github.com/privacytrustlab/ml_privacy_meter based on [Song, Shmatikov] Auditing Data Provenance in Text-Generation Models, KDD'19

Examples of Vulnerable Training Data



But it gets worse. And this is very important, which is generic. It doesn't have anything to do, in specific, with a power plant. It would work as well, for example, in a power plant or in an automobile factory. It is generic. And you don't have -- as an attacker -- you don't have to deliver this payload by a USB stick, as we saw it in the case of Stuxnet. You could also use conventional worm technology for spreading. Just spread it as

Chris Anderson: I've got a question. Ralph, it's been quite widely reported that people assume that Mossad is the main entity behind this. Is that your opinion?

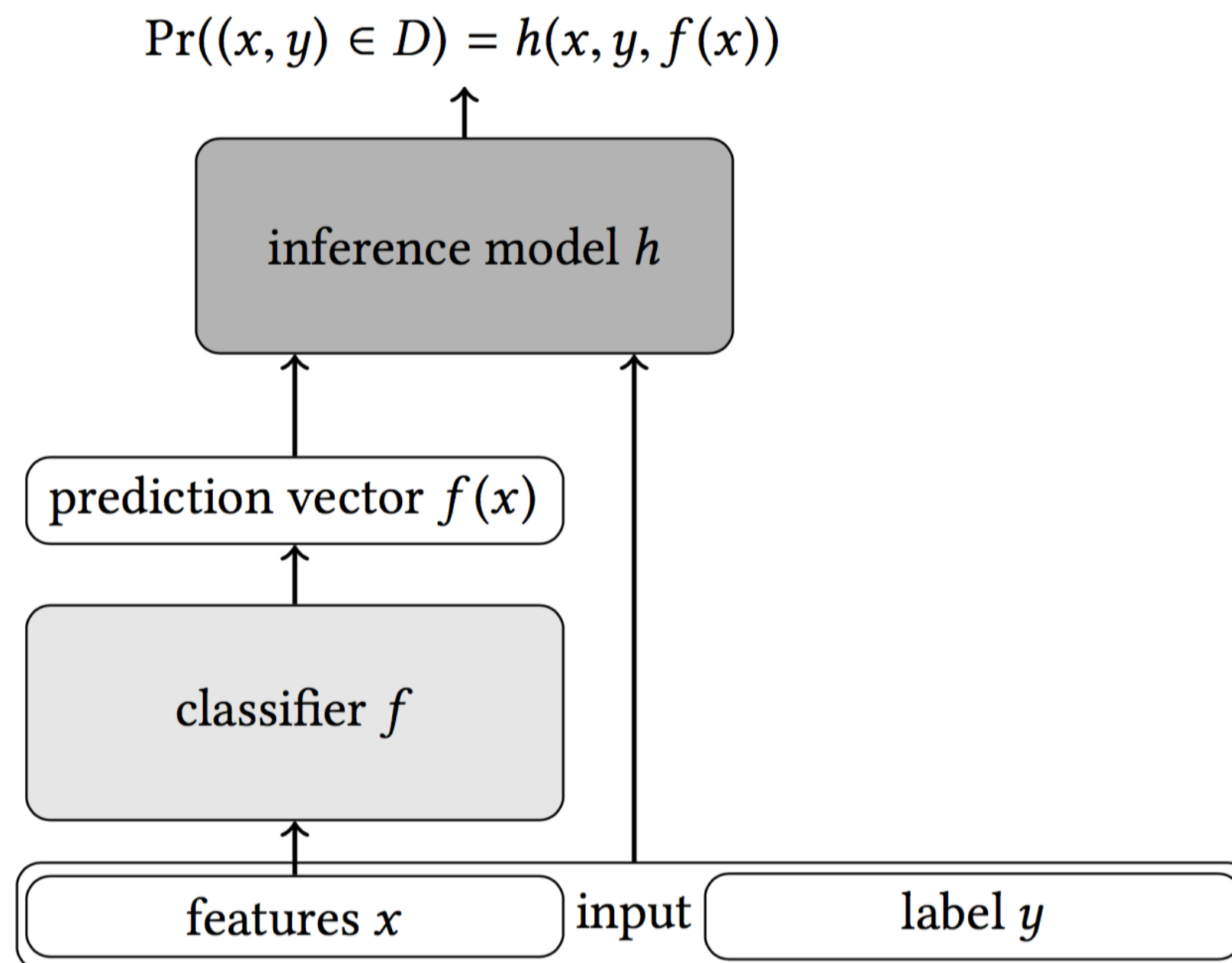
Ralph Langner: Okay, you really want to hear that? Yeah. Okay. My opinion is that the Mossad is involved, but that the leading force is not Israel. So the leading force behind that is the cyber superpower. There is only one, and that's the United States -- fortunately, fortunately. Because otherwise, our problems would even be bigger.

Examples of Vulnerable Training Data

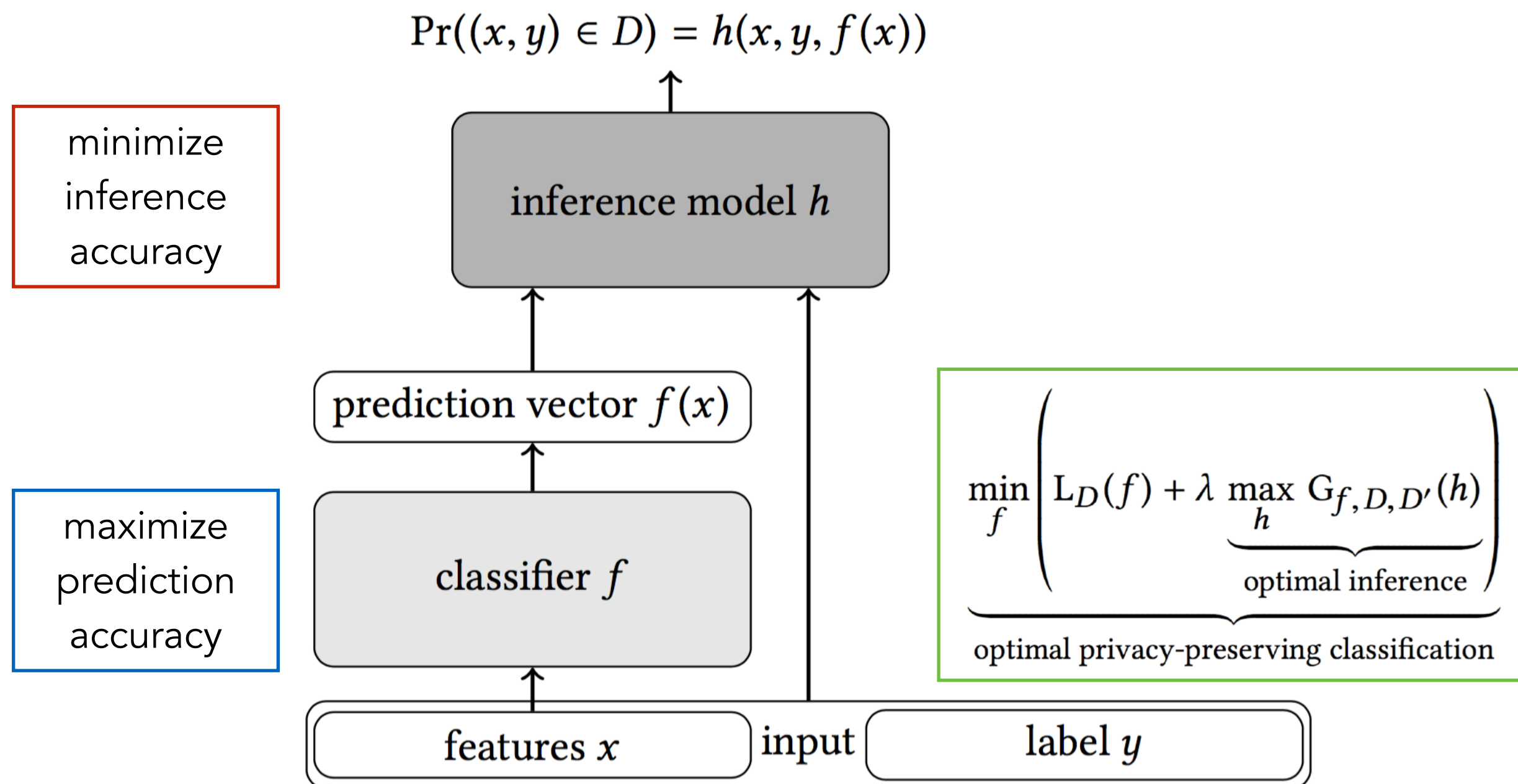


This year, Germany is celebrating the 25th anniversary of the peaceful revolution in East Germany. In 1989, the Communist regime was moved away, the Berlin Wall came down, and one year later, the German Democratic Republic, the GDR, in the East was unified with the Federal Republic of Germany in the West to found today's Germany. Among many other things, Germany inherited the archives of the East German secret police, known as the Stasi. Only two years after its dissolution, its documents were opened to the public, and historians such as me started to study these documents to learn more about how the GDR surveillance state functioned.

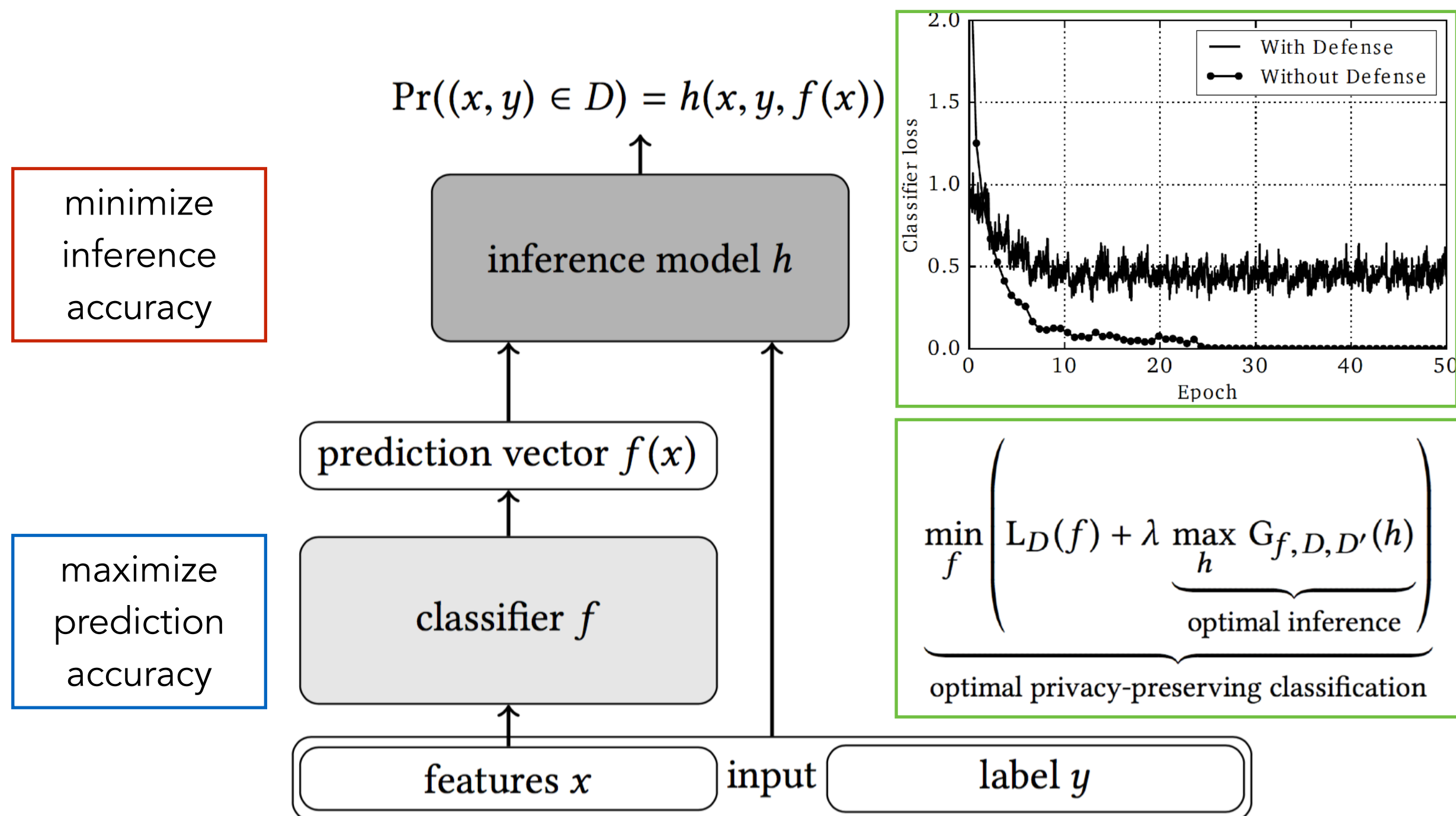
Privacy as a Learning Objective



Privacy as a Learning Objective



Privacy as a Learning Objective



Privacy and Generalization

	Without defense			With defense		
Dataset	Training accuracy	Testing accuracy	Attack accuracy	Training accuracy	Testing accuracy	Attack accuracy
Purchase100	100%	80.1%	67.6%	92.2%	76.5%	51.6%
Texas100	81.6%	51.9%	63%	55%	47.5%	51.0%
CIFAR100- Alexnet	99%	44.7%	53.2%	66.3%	43.6%	50.7%
CIFAR100- DenseNET	100%	70.6%	54.5%	80.3%	67.6%	51.0%



Smaller gap



Random guess

Bound the Worst-case Privacy Loss

- Differential Privacy: Ensure the indistinguishability between two models which are trained on two neighboring datasets.
- Randomize the training algorithm to bound the privacy loss

mechanism (randomized model)

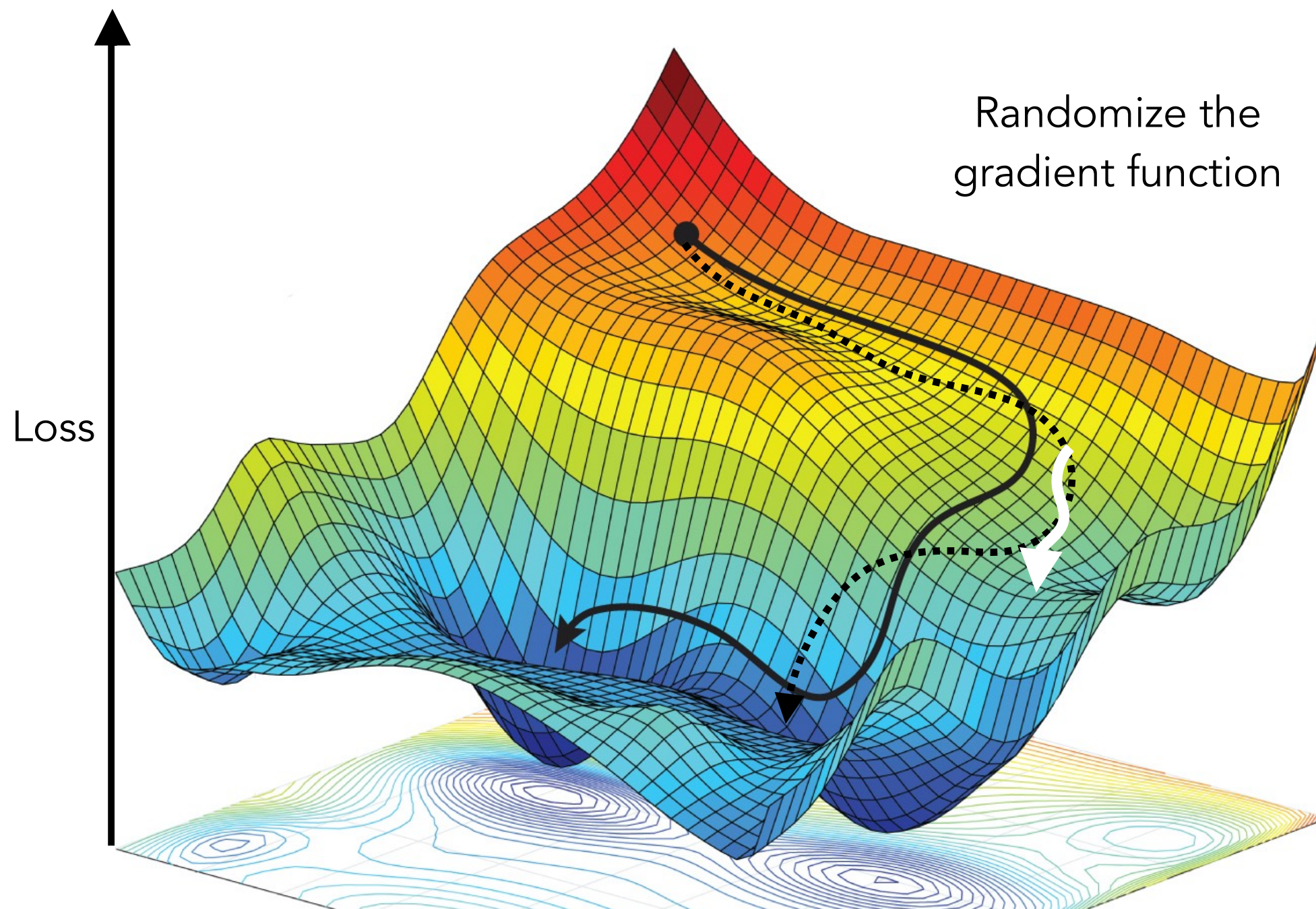
privacy loss

observables

$$\mathcal{L}_{\mathcal{M}(x) || \mathcal{M}(y)}^{(\xi)} = \ln \left(\frac{\Pr[\mathcal{M}(x) = \xi]}{\Pr[\mathcal{M}(y) = \xi]} \right)$$

input (dataset)

DP Stochastic Gradient Descent



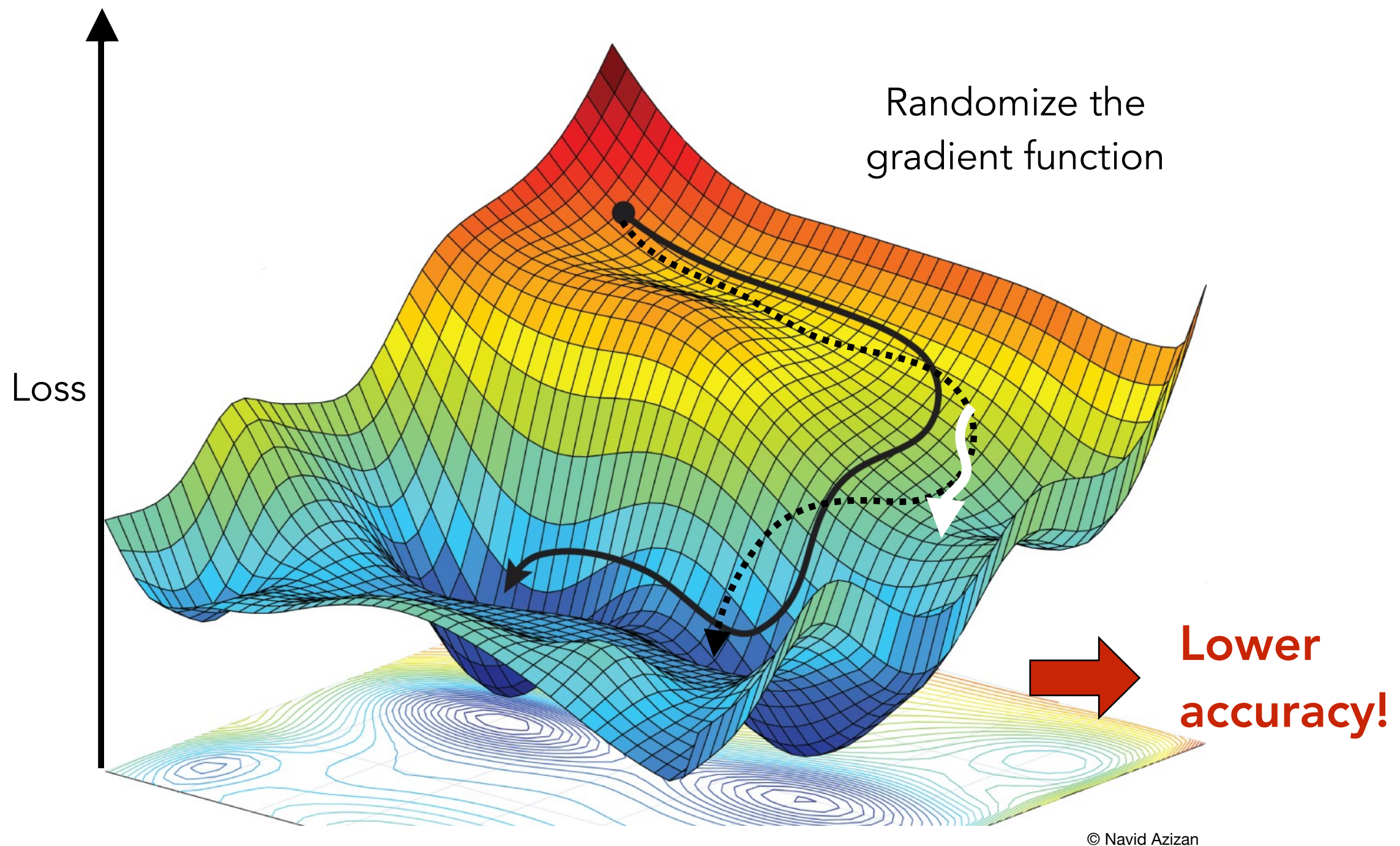
© Navid Azizan

[Bassily, Smith, Thakurta] Private empirical risk minimization: Efficient algorithms and tight error bounds, FOCS'14

[Shokri and Shmatikov] Privacy-Preserving Deep Learning, CCS'15

[Abadi, et al.] Deep learning with differential privacy. CCS'16.

DP Stochastic Gradient Descent



[Bassily, Smith, Thakurta] Private empirical risk minimization: Efficient algorithms and tight error bounds, FOCS'14

[Shokri and Shmatikov] Privacy-Preserving Deep Learning, CCS'15

[Abadi, et al.] Deep learning with differential privacy. CCS'16.

Causes of Performance Loss

- Computation of total privacy loss is not exact (i.e., the upper bound of the privacy loss (epsilon) is not tight). By overestimating the privacy loss, the added noise is larger than what is really needed to achieve the same true level of privacy

Causes of Performance Loss

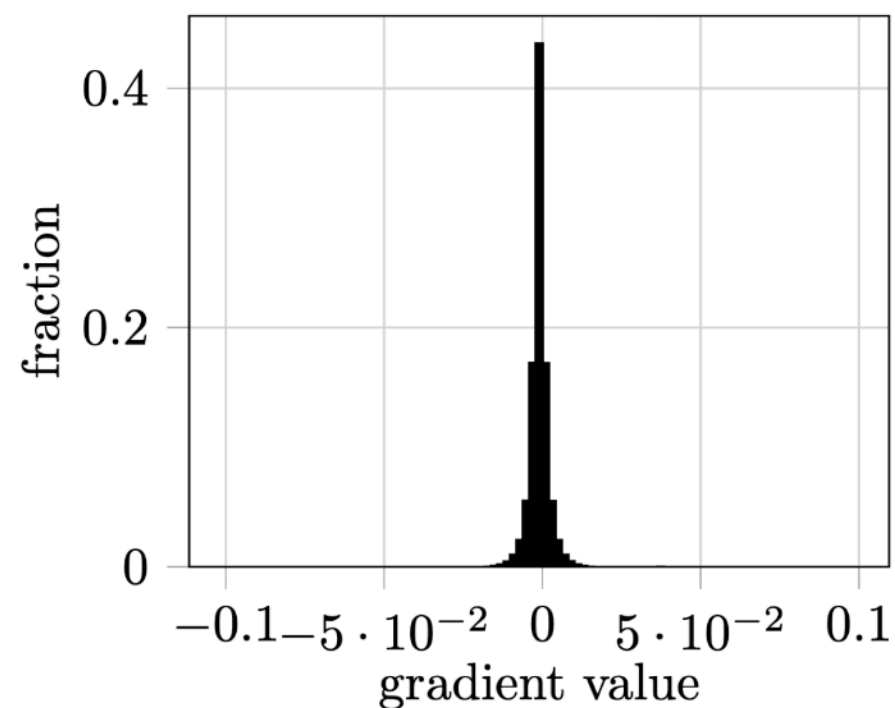
- Computation of total privacy loss is not exact (i.e., the upper bound of the privacy loss (epsilon) is not tight). By overestimating the privacy loss, the added noise is larger than what is really needed to achieve the same true level of privacy
- Gaussian mechanism is not a utility-preserving mechanism for DP SGD

Causes of Performance Loss

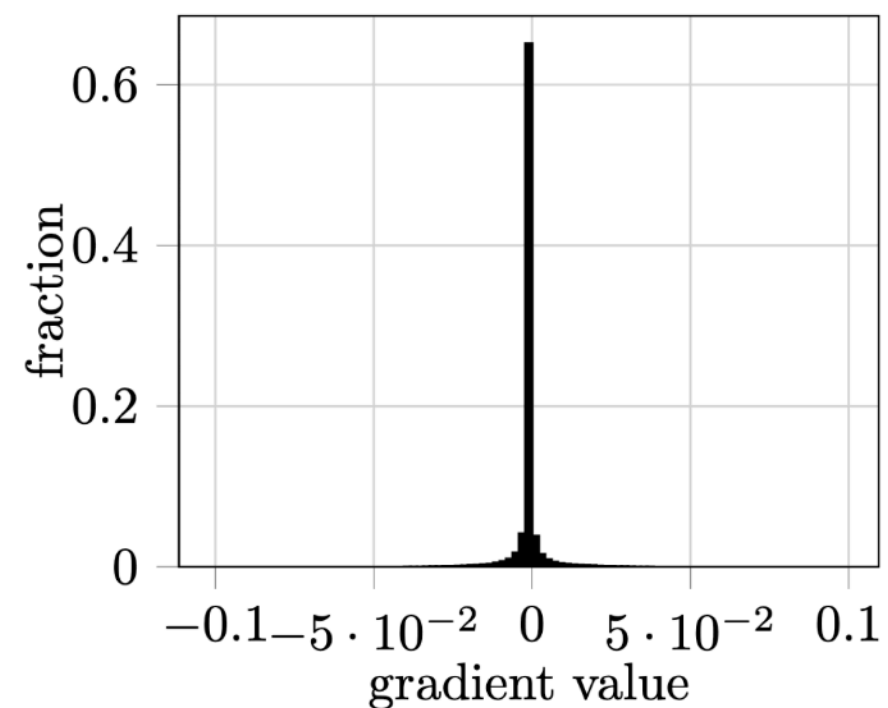
- Computation of total privacy loss is not exact (i.e., the upper bound of the privacy loss (epsilon) is not tight). By overestimating the privacy loss, the added noise is larger than what is really needed to achieve the same true level of privacy
- Gaussian mechanism is not a utility-preserving mechanism for DP SGD
- All randomized gradient vectors are treated equally (but, the signal to noise ratio is not the same across all, and their influence on the parameter vector should not be the same)

Observation

- Gradients follow a symmetric distribution, concentrated around zero



(a) CIFAR



(b) MNIST

- The DP noise would dominate the gradient values

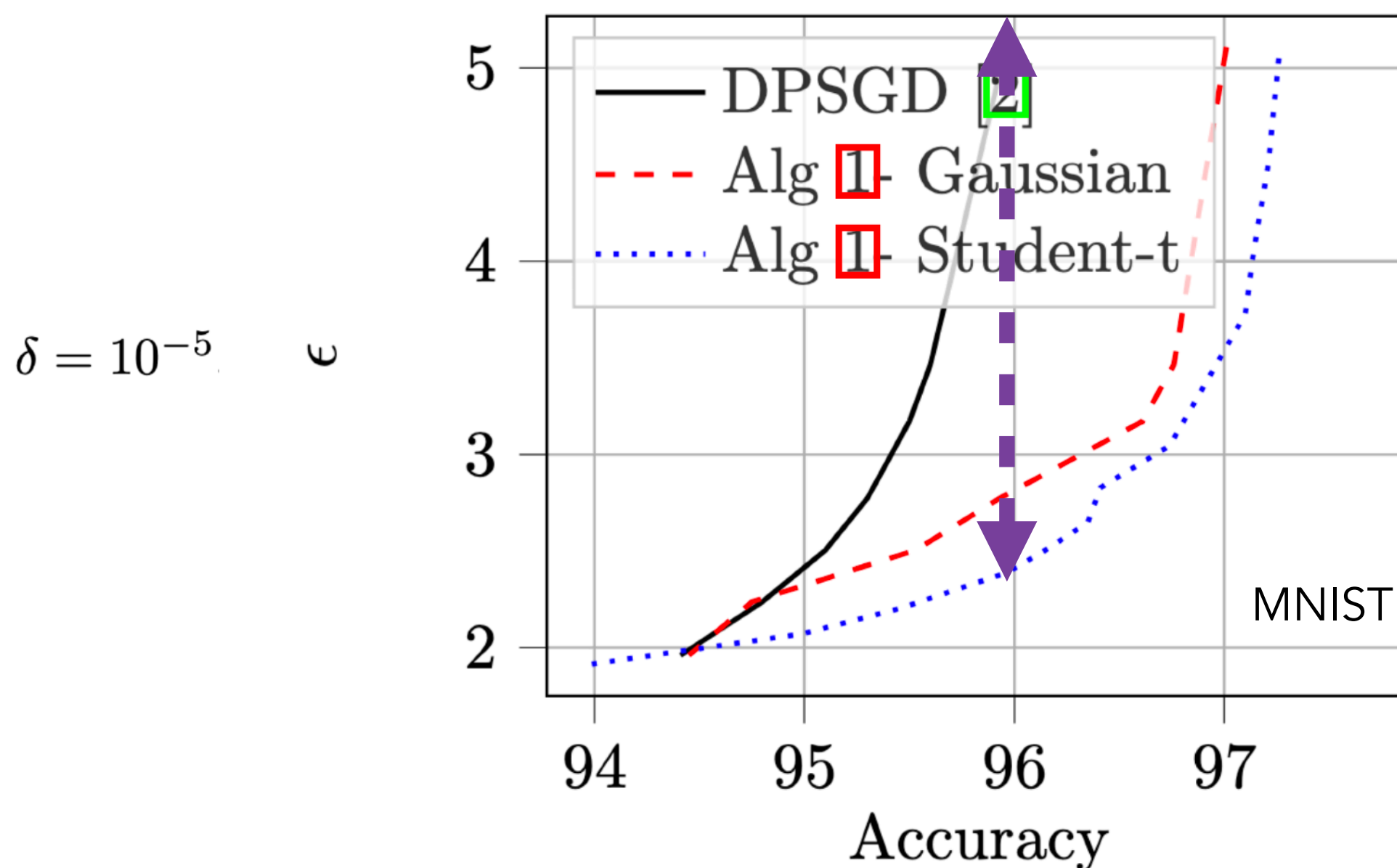
Gradient Coding and De-noising

- Randomize gradients using a student-t distribution
 - To compute DP parameters, encode gradient values into a finite number of samples from a Gaussian distribution

Gradient Coding and De-noising

- Randomize gradients using a student-t distribution
 - To compute DP parameters, encode gradient values into a finite number of samples from a Gaussian distribution
- Weighted update of model parameters
 - Lower the weight if noise dominates the signal

Significant Privacy Improvement for the Same Test Accuracy





Data Privacy and Trustworthy ML Research Lab

National University of Singapore

reza@nus.edu.sg

 @rzshokri

<https://www.comp.nus.edu.sg/~reza/>