

# CS489/698

# Privacy, Cryptography, Network and Data Security

---

Privacy-Preserving Machine Learning

Fall 2024, Tuesday/Thursday 02:30pm-03:50pm

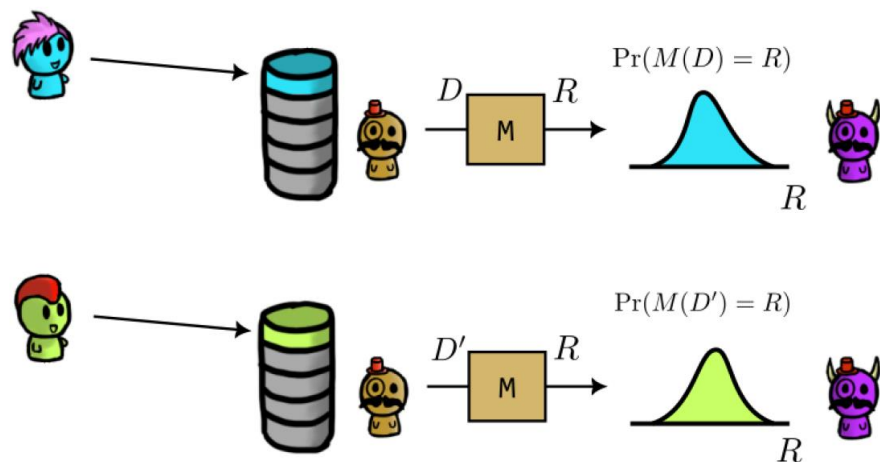
# Prologue: a couple more DP properties

---

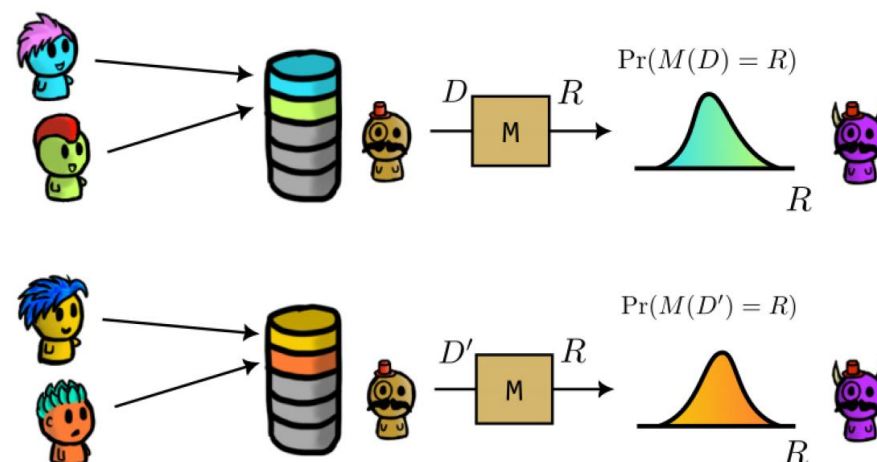
# Recap on Group privacy

**Group privacy:** Let  $M: \mathcal{D} \rightarrow \mathcal{R}$  be a mechanism that provides  $\epsilon$ -DP for  $D, D'$  that differ in one entry. Then, it provides  $k\epsilon$ -DP for datasets  $D, D'$  that differ in  $k$  entries.

If this is  $\epsilon$ -DP....



... then this is  $2\epsilon$ -DP



# Group privacy with $(\epsilon, \delta)$ -DP

---

- For approximate DP,  $\delta$  gets an additional factor of  $ke^{(k-1)\epsilon}$  :

**Group privacy:** Let  $M: \mathcal{D} \rightarrow \mathcal{R}$  be a mechanism that provides  $(\epsilon, \delta)$ -DP for  $D, D'$  that differ in one entry. Then, it provides  $(k\epsilon, ke^{(k-1)\epsilon}\delta)$ -DP for datasets  $D, D'$  that differ in  $k$  entries.

# Renyi Differential Privacy

---

- Differential privacy is a very ambitious privacy guarantee, that protects against a worst-case adversary that potentially knows  $D$  and  $D'$ , and for all possible outputs of the mechanism.
- $\epsilon$  and  $\delta$  provide a very **limited** and **pessimistic** description of the differences between  $\Pr(M(D) \in S)$  and  $\Pr(M(D') \in S)$ .
- There are other *relaxed* notions of DP that capture other nuances between these distributions, allowing for a tighter analysis.
  - Relaxes how much we care about the worst case (**sometimes very unlikely**)
  - A popular one is **Renyi Differential Privacy**

# Renyi Differential Privacy

- To introduce Renyi DP we need to know Renyi Divergence

**Renyi Divergence:** given two probability distributions  $P$  and  $Q$ , the Renyi divergence of order  $\alpha > 1$  is

$$D_{\alpha}(P||Q) = \frac{1}{\alpha - 1} \log \mathbb{E}_{x \sim Q} \left( \frac{P(x)}{Q(x)} \right)^{\alpha} = \frac{1}{\alpha - 1} \log \left( \int_x P(x)^{\alpha} Q(x)^{1-\alpha} dx \right)$$

- The logarithm measures how much **more** or **less** likely an event  $x$  is under  $P$  compared to  $Q$ .  
(divergence is intended to measure **how much information is lost** when using  $Q(x)$  instead of  $P(x)$  )

# Renyi Differential Privacy

**Renyi Divergence:** given two probability distributions  $P$  and  $Q$ , the Renyi divergence of order  $\alpha = \infty$  (defined by its limit) is

$$D_{\infty}(P||Q) = \sup_x \log \left( \frac{P(x)}{Q(x)} \right) = \log \left( \max_x \frac{P(x)}{Q(x)} \right)$$

**Renyi Divergence:** given two probability distributions  $P$  and  $Q$ , the Renyi divergence of order  $\alpha = 1$  is the Kullback-Leibler divergence

$$D_1(P||Q) = \mathbb{E}_{x \sim P} \log \left( \frac{P(x)}{Q(x)} \right) = \int P(x) \log \left( \frac{P(x)}{Q(x)} \right) dx$$

(note that  $x \sim P$  in the expectation)

# Renyi Differential Privacy – DP connection

**Renyi DP:** a mechanism  $M: \mathcal{D} \rightarrow \mathcal{R}$  is  $(\epsilon, \alpha)$ -RDP (also read as “ $\epsilon$ -RDP of order  $\alpha$ ”) if, for any neighboring datasets  $D, D'$  it holds that

$$D_{\alpha}(M(D) || M(D')) \leq \epsilon$$

- Recall that, when  $\alpha = \infty$ , then the divergence (defined by its limit) is:

$$D_{\infty}(M(D) || M(D')) = \sup_x \log \left( \frac{\Pr(M(D) \in x)}{\Pr(M(D') \in x)} \right)$$

- In that case, it is easy to see that  $(\epsilon, \infty)$ -RDP is equivalent to  $\epsilon$ -DP



# Renyi Differential Privacy: Properties

---

**RDP Sequential Composition:** if  $M_1$  is  $(\alpha, \epsilon_1)$ -RDP and  $M_2$  is  $(\alpha, \epsilon_2)$ -RDP, then the sequential composition  $(M_1, M_2)$  satisfies  $(\alpha, \epsilon_1 + \epsilon_2)$ -RDP

**RDP to DP:** if  $M$  is  $(\alpha, \epsilon)$ -RDP, then it is also  $\left(\epsilon + \frac{\log\left(\frac{1}{\delta}\right)}{\alpha-1}, \delta\right)$ -DP

# Renyi Differential Privacy: Example

---

Let's consider the following probability distributions:

- $P = \{P(x_1) = 0.6, P(x_2) = 0.4\}$ ;  $Q = \{Q(x_1) = 0.5, Q(x_2) = 0.5\}$

**For  $\alpha = 1$  (KL Divergence):**

**For  $\alpha > 1$  (e.g.,  $\alpha=2$ ):**

**For  $\alpha = \infty$  :**

# Renyi Differential Privacy: Example

---

Let's consider the following probability distributions:

- $P = \{P(x_1) = 0.6, P(x_2) = 0.4\}$ ;  $Q = \{Q(x_1) = 0.5, Q(x_2) = 0.5\}$

**For  $\alpha = 1$  (KL Divergence):**

$$D_1(P||Q) = P(x_1)\log\left(\frac{P(x_1)}{Q(x_1)}\right) + P(x_2)\log\left(\frac{P(x_2)}{Q(x_2)}\right) \approx 0.1092 - 0.0892 = 0.02$$

**For  $\alpha > 1$  (e.g.,  $\alpha=2$ ):**

$$D_2(P||Q) = \frac{1}{2-1} \log (P(x_1)^2 Q(x_1)^0 + P(x_2)^2 Q(x_2)^0) = \log(0.36 + 0.16) \approx -0.653.$$

**For  $\alpha = \infty$  :**

$$D_\infty(P||Q) = \log\left(\max_x \frac{P(x)}{Q(x)}\right) \approx 0.182 \quad \text{For } x_1 \left(\frac{0.6}{0.5}\right) = 1.2 \quad \text{For } x_2 \left(\frac{0.4}{0.5}\right) = 0.8$$

# Many other variations...

- An SOK from 2020

Name & references		
$(\epsilon, \delta)$ -approximate DP [52]	$(D, t, \epsilon)$ -per-instance DP [162]	$(\Theta, \epsilon, \delta)$ -active PK DP [11, 14, 35]
$(\epsilon, \delta)$ -probabilistic DP [20, 124, 127]	$(\mathcal{R}, \epsilon)$ -generic DP [105]	$(\Theta, \epsilon, \delta)$ -passive PK DP [35]
$\epsilon$ -Kullback-Leiber Pr [9, 31]	$(G, \mathcal{I}_Q, \epsilon)$ -blowfish Pr [84, 86]	$(\Theta, \Phi, \epsilon)$ -pufferfish Pr [106]
$(\alpha, \epsilon)$ -Rényi DP [128]	$\epsilon$ -adjacency-relation div. DP [97]	$(\Theta, \epsilon, \delta)$ -distribution Pr [98]
$\epsilon$ -mutual-information DP [31]	$\Psi$ -personalized DP [59, 76, 94, 118]	$(d, \Theta, \epsilon)$ -extended DnPr [98]
$(\mu, \tau)$ -mean concentrated DP [58]	$\Psi$ -tailored DP/ $\epsilon(\cdot)$ -outlier Pr [120]	$(f, \Theta, \epsilon)$ -divergence DnPr [97]
$(\xi, \rho)$ -zero concentrated DP [19]	$(\pi, \gamma, \epsilon)$ -random DP [83]	$(d, f, \Theta, \epsilon)$ -ext. div. DnPr [97]
$(f, \epsilon)$ -divergence DP [9]	$d_{\mathcal{D}}$ -Pr [22]	$(\Theta, \epsilon)$ -positive membership Pr [114]
$\epsilon$ -unbounded DP [105]	$(\epsilon, \gamma)$ -distributional Pr [141, 177]	$(\Theta, \epsilon, \delta)$ -adversarial Pr [139]
$\epsilon$ -bounded/attribute/bit DP [105]	$(\epsilon(\cdot), \delta(\cdot))$ -endogenous DP [107]	$(\Theta, \epsilon)$ -aposteriori noiseless Pr [14]
$(c, \epsilon)$ -group DP [49]	$(d_{\mathcal{D}}, \epsilon, \delta)$ -pseudo-metric DP [36]	$\epsilon$ -semantic Pr [69, 96]
$\epsilon$ -free lunch Pr [105]	$(\theta, \epsilon, \gamma, \delta)$ -typical Pr [10]	$(\text{Agg}, \epsilon)$ -zero-knowledge Pr [72]
$(R, c, \epsilon)$ -dependent DP [116]	$(\Theta, \epsilon)$ -on average KL Pr [164]	$(\Theta, \Gamma, \epsilon)$ -coupled-worlds Pr [11]
$(P, \epsilon)$ -one-sided DP [42]	$(f, d, \epsilon)$ -extended divergence DP [97]	$(\Theta, \Gamma, \epsilon, \delta)$ -inference-based CW Pr [11]
$(D, \epsilon)$ -individual DP [149]	$(\mathcal{R}, M)$ -general DP [103]	$\epsilon_{\kappa}$ -SIM-computational DP [129]
	$(\Theta, \epsilon)$ -noiseless Pr [14, 44]	$\epsilon_{\kappa}$ -IND-computational DP [129]
	$(\Theta, \epsilon)$ -distributional DP [11, 35]	$(\text{Agg}, \epsilon)$ -computational ZK Pr [72]

# Primer on Machine Learning

---

# Machine learning: quick primer

---

- For simplicity, we will focus on a *classification problem* with *supervised learning*.
  - Unsupervised or Reinforcement learning are other types
- We have a training set  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$  with  $n$  samples. Given a sample  $(x_i, y_i)$ ,  $x_i$  are the *features* and  $y_i$  is its *label*.

# Machine learning: quick primer

---

- For simplicity, we will focus on a *classification problem* with *supervised learning*.
  - Unsupervised or Reinforcement learning are other types
- We have a training set  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$  with  $n$  samples. Given a sample  $(x_i, y_i)$ ,  $x_i$  are the *features* and  $y_i$  is its *label*.
- We want to produce a function  $f: \mathcal{X} \rightarrow \mathcal{Y}$  that can *predict* a sample's label from its features.

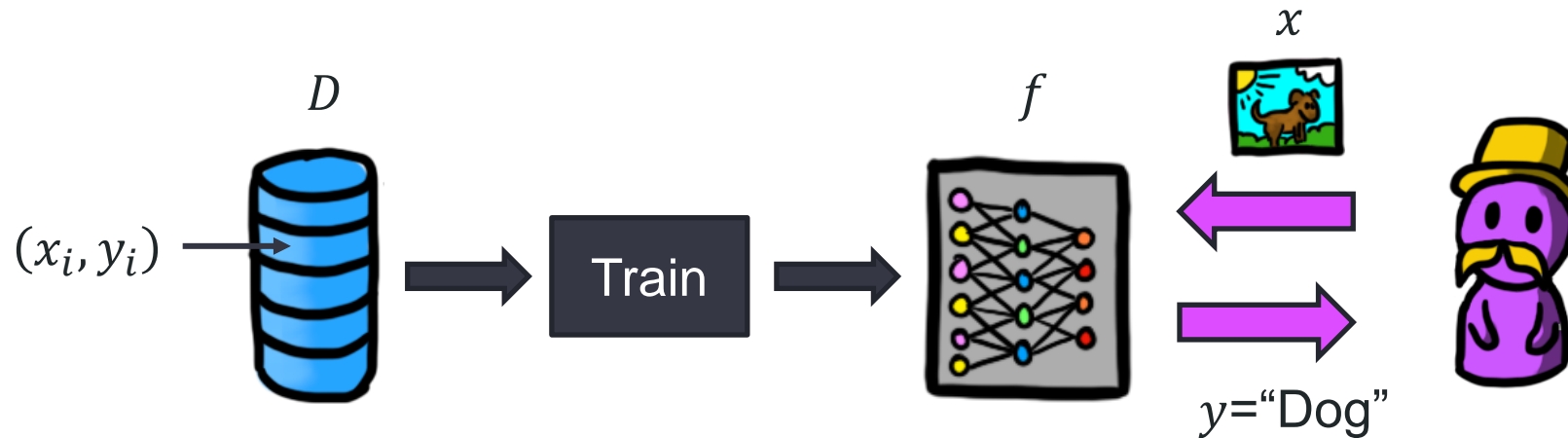
# Machine learning: quick primer

---

- For simplicity, we will focus on a *classification problem* with *supervised learning*.
  - Unsupervised or Reinforcement learning are other types
- We have a training set  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$  with  $n$  samples. Given a sample  $(x_i, y_i)$ ,  $x_i$  are the *features* and  $y_i$  is its *label*.
- We want to produce a function  $f: \mathcal{X} \rightarrow \mathcal{Y}$  that can *predict* a sample's label from its features.
- We will use the training set to train such a function. Ideally, it should correctly predict labels for *unseen* samples (e.g., samples in a testing set).
  - We will say that a model *generalizes* well if it has *high accuracy on unseen samples*
  - A model *overfits* if it works perfectly for samples in the training set but *does not generalize*.



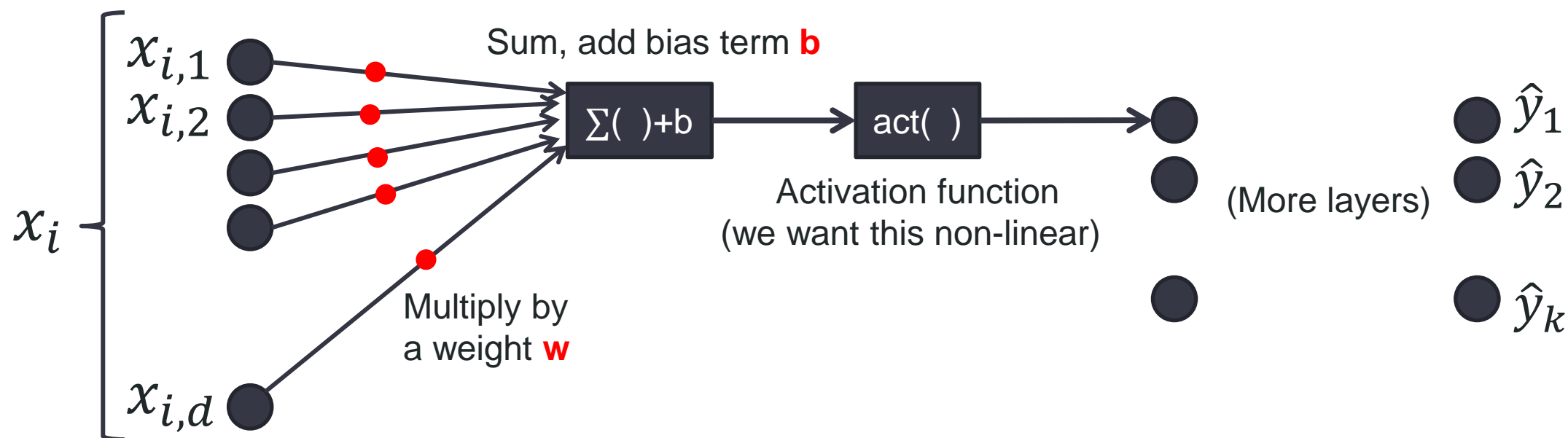
# Machine learning: quick primer



Usually, this gives **confidence scores** for each class:  $(\hat{y}_1, \hat{y}_2, \dots, \hat{y}_k)$   
For example: ["Dog", "Cat", "Mouse" ...]=[0.81, 0.10, 0.03, ...]

# Neural networks

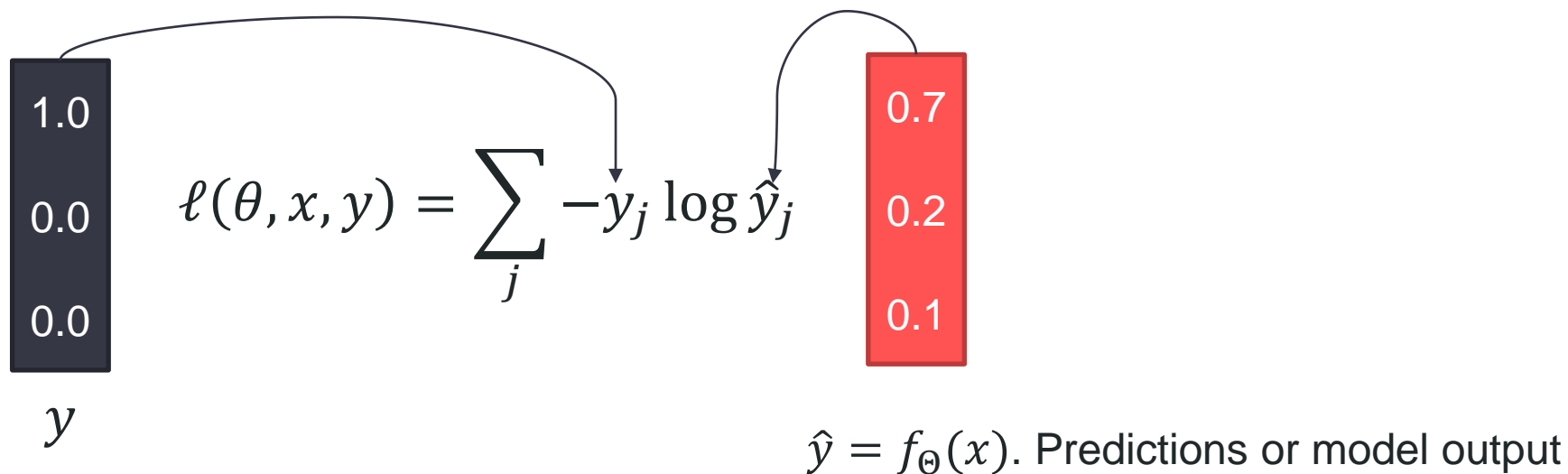
- There are many architectures for machine learning models (i.e., many structures for the function  $f$ ).
- One of the most popular are **neural networks**.



Training the model means tuning all  $w$ 's and  $b$ 's

# Loss Functions

- We define a *loss function* that we want to minimize:  $\ell(\theta, x, y)$ , where  $\theta$  are the parameters  $w$  and  $b$ .
  - For example, a typical loss function is  $\ell(\theta, x, y) = \sum_j -y_j \log \hat{y}_j$  where  $y_j$  is only 1 for the *true label* of the sample,  $j$ .

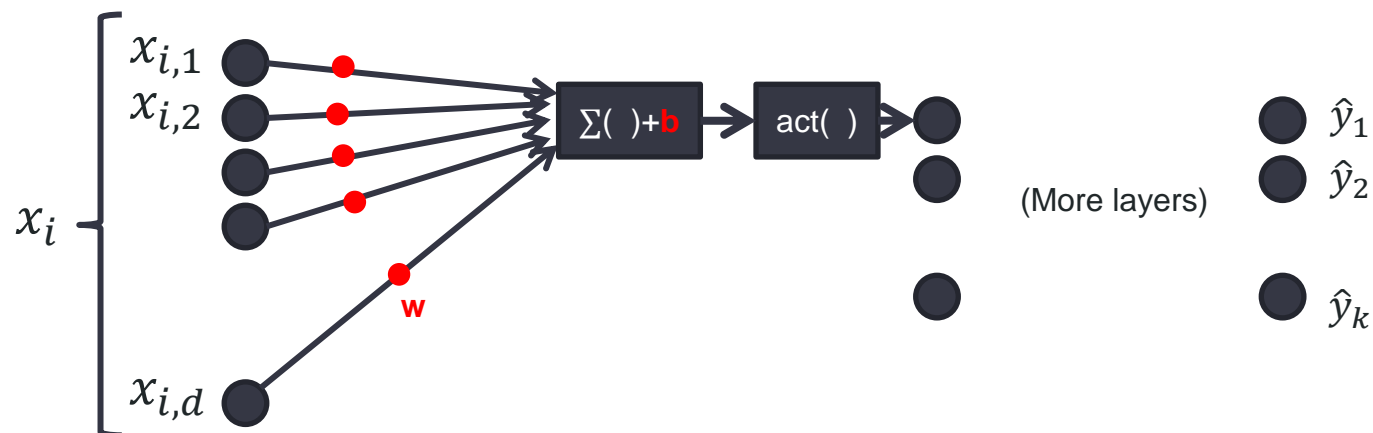


# Training neural networks

- Since we have the training set  $D$ , it makes sense to minimize the empirical loss in this training set:

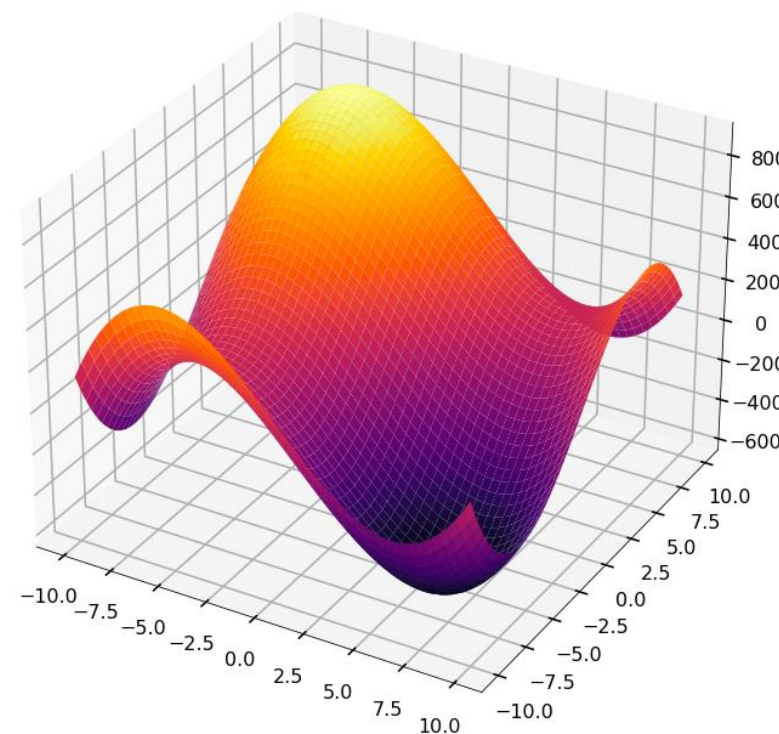
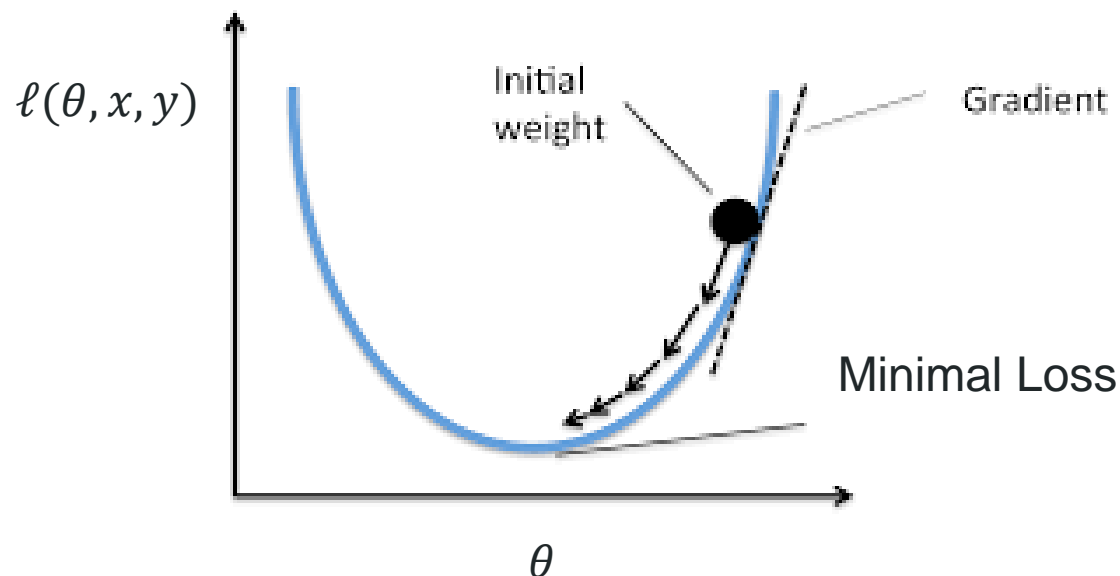
$$\mathcal{L}(\theta, D) = \frac{1}{N} \sum_i \ell(\theta, x_i, y_i)$$

- In practice, the minimization is done using **Stochastic Gradient Descent** (SGD).



# Gradient Descent

- The gradient of the loss  $\nabla \ell(\theta, x, y)$  evaluated at  $(x, y)$  is the derivative with respect to each parameter  $\theta_i$  (every **w** and **b**).
- It tells us the direction in which  $\theta$  should go to minimize the loss (for sample  $(x, y)$ ).



# Gradient Descent

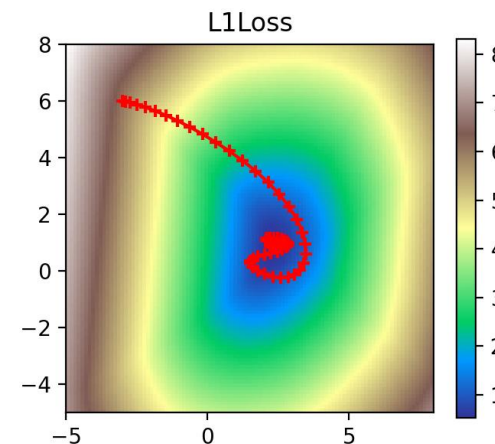
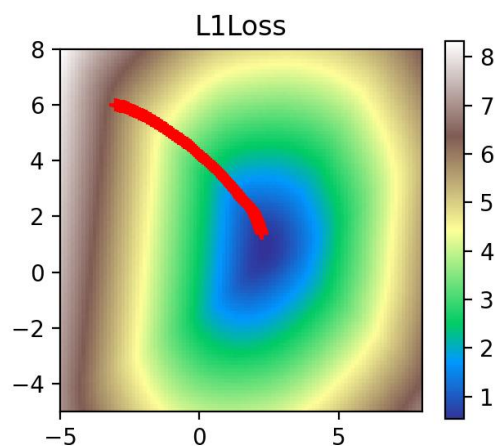
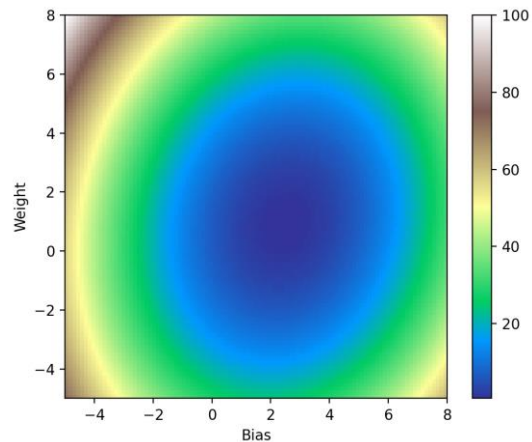
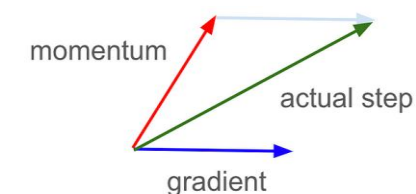
---

- We could minimize the loss by running several steps (**epochs**) of Gradient Descent:
  - For each step  $t \in [T]$ :
$$\theta_t = \theta_{t-1} - \eta \nabla \mathcal{L}(\theta_{t-1}, D)$$
  - $\eta$  is the *learning rate*
- This is expensive, so usually we do these iterations over a subset of the training sets (**batches**)
- Note  $\theta$  represents parameters,  $\eta$  and  $T$  are hyper-parameters

# Stochastic Gradient Descent – with Mini Batches

For each training step  $t \in [T]$ :

1. Take a batch  $B$  of  $L$  samples from  $D$
2. For each  $(x_i, y_i) \in B$ , compute the gradient  $g_i = \nabla \ell(\theta_{t-1}, x_i, y_i)$
3. Average the gradients  $g = \frac{1}{L} \sum_i g_i$
4. Descend  $\theta_t = \theta_{t-1} - \eta \cdot g$



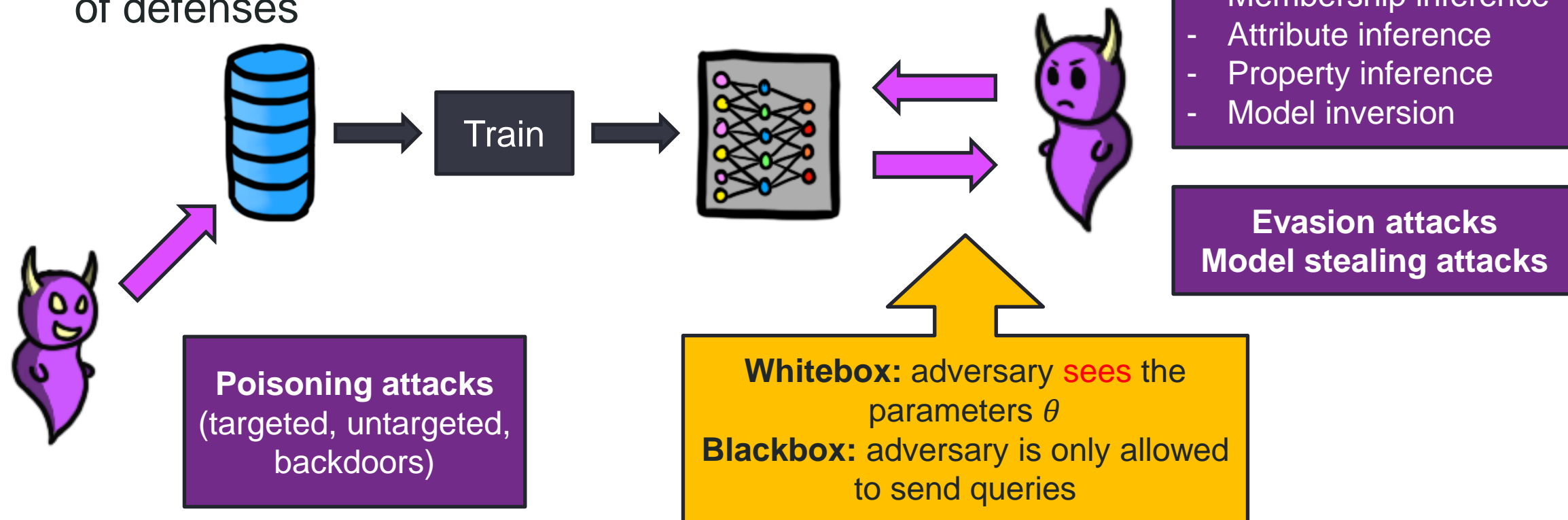
# Inference Attacks in ML

---



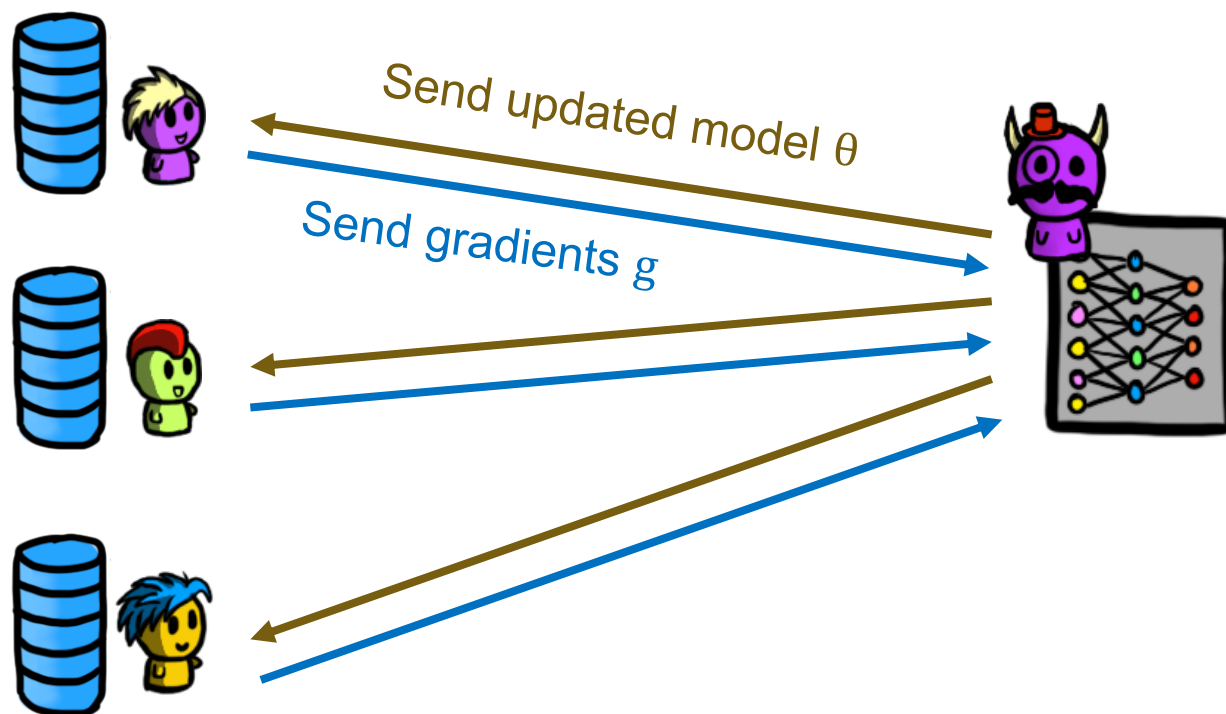
# Attacking ML models

- There are many types of attacks against ML
- Later we will see that there are also different types of defenses



# Attacking ML models in Federated Learning

- Federated Learning: a centralized server builds a model, a set of clients send updates (**gradients**) using their local datasets

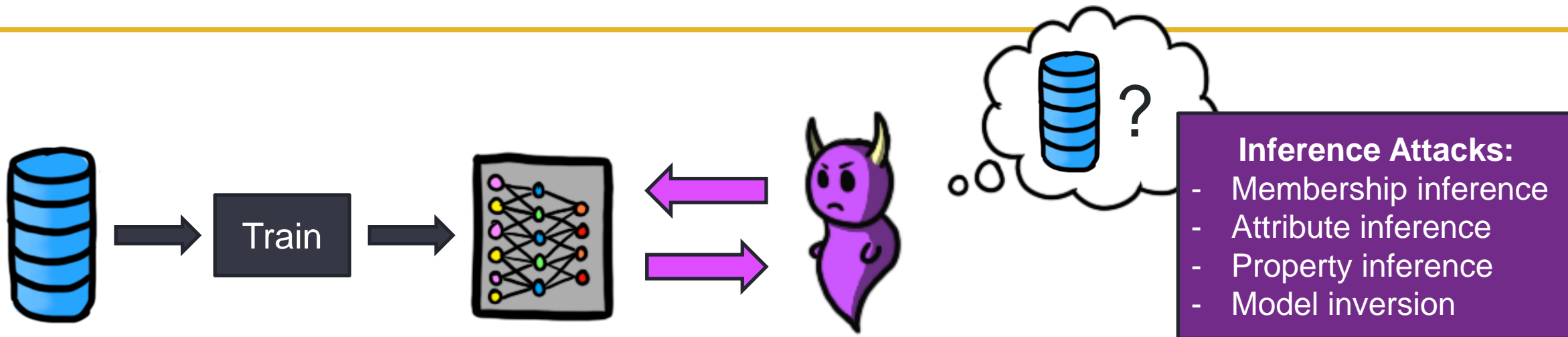


**Poisoning attacks**  
(targeted,  
untargeted,  
backdoors)

**Inference Attacks:**  
(adv **sees** all intermediate  
gradients, can potentially  
send **malicious  $\theta$** )

- Membership inference
- Attribute inference
- Property inference
- ...

# Inference attacks



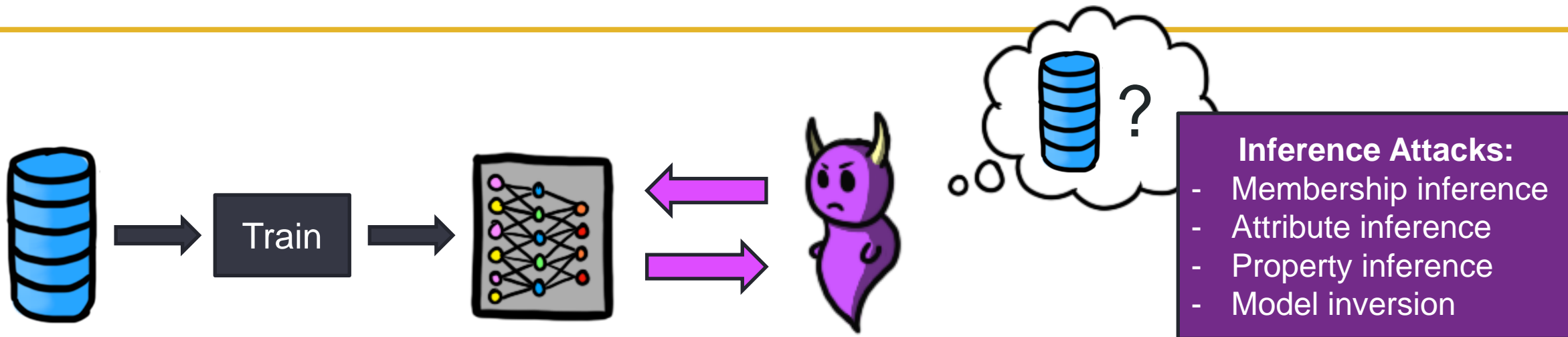
**Membership Inference:**  
Is a given sample **in** the training set?

**Attribute Inference:**  
Given a sample with some missing attributes, **can we guess** them?

**Property Inference:**  
Given a property about the *whole* training set, can we guess if it's **true or not?**

**Model inversion:**  
Given a label, can we find a representative element of this class? (**learn  $x$  from  $y$** )

# Inference attacks



## Membership Inference:

Is a given sample **in** the training set?

## Attribute Inference:

Given a sample with some missing attributes, **can we guess** them?

## Property Inference:

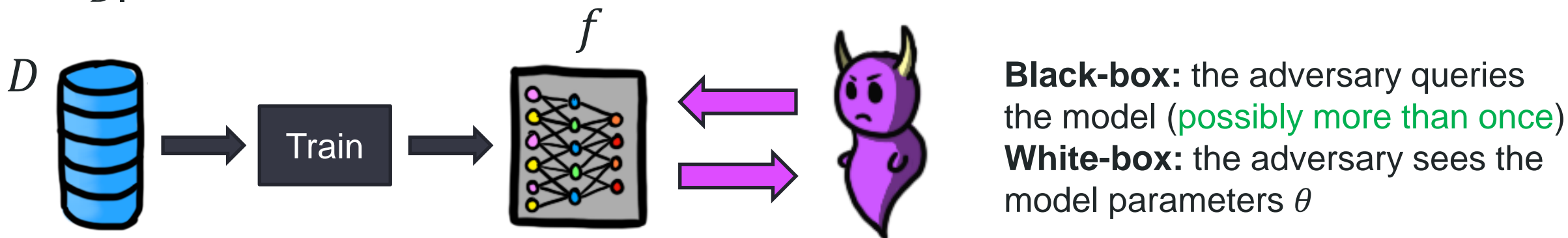
Given a property about the *whole* training set, can we guess if it's **true or not**?

## Model inversion:

Given a label, can we find a representative element of this class? (**learn  $x$  from  $y$** )

# Membership Inference Attacks (MIAs)

- Given a sample  $(x, y)$ , and a model  $f$  trained with dataset  $D$ , guess whether  $(x, y) \in D$ .



- With only black-box access, and a model that outputs confidence scores:
  - $f(x) = [\hat{v}_1, \hat{v}_2, \dots, \hat{v}_i, \dots]$  where  $\hat{v}_i$  are confidence scores for label  $i$

**Q:** If you were the adversary, with a target sample  $(x, y)$  and black-box access to the model  $f$ , how would you guess if the target sample is a member?

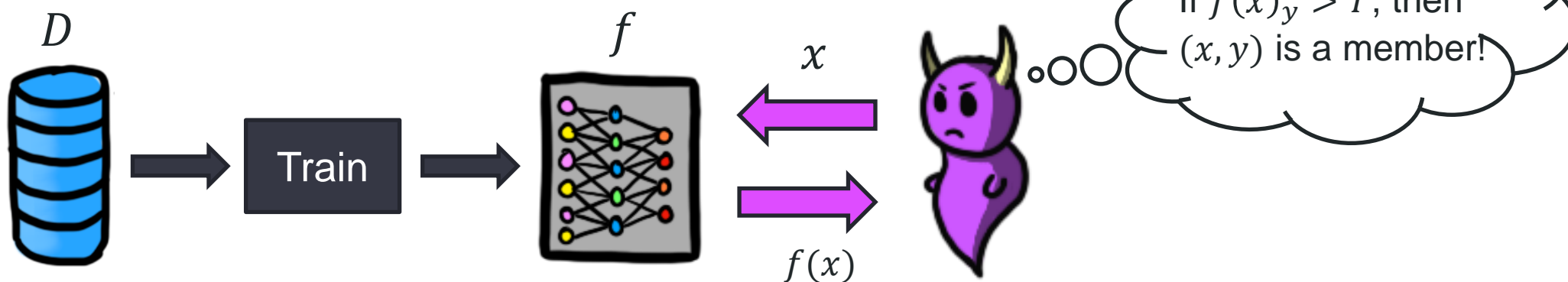
# Threshold Attacks

- *Idea:* the model will be **more confident** on samples it has seen during training.

## Threshold attack

- This attack queries the model on **sample  $x$**  and then measures the confidence score assigned to its **true label  $y$** .
- If the confidence score is above some **threshold**, then the attack decides the sample is a *member*.

**Q:** how can the attacker compute this threshold?



Yeom et al. "Privacy risk in machine learning: Analyzing the connection to overfitting." *CSF*, 2018.

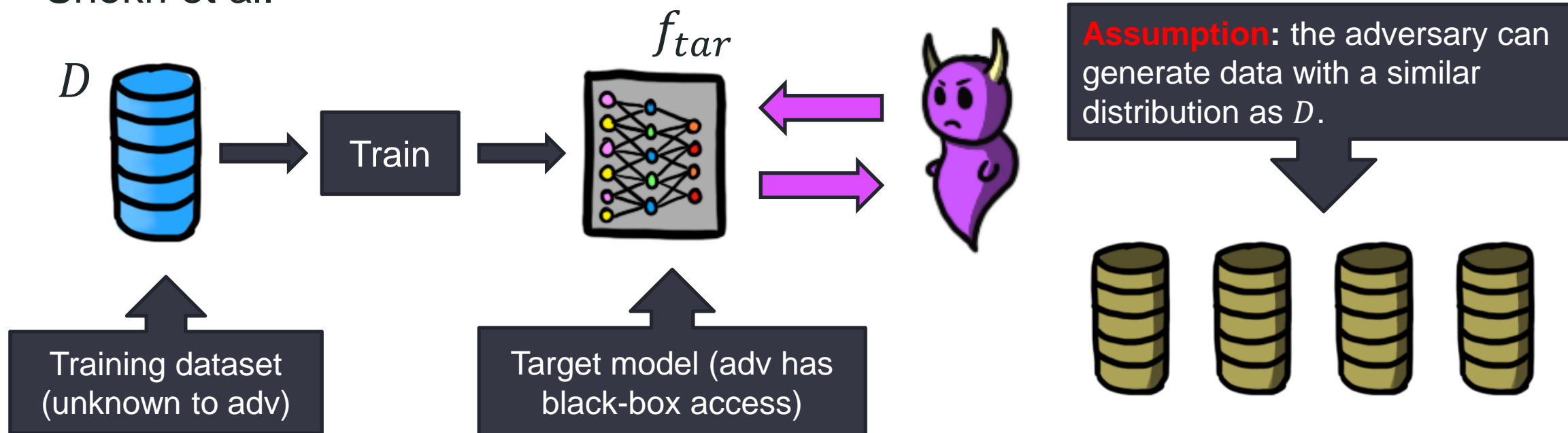
# Neural Network-based Attacks

- Other MIAs use Machine Learning against Machine Learning.



# Neural Network-based Attacks

- Other MIAs use Machine Learning against Machine Learning.
- The first NN-based attack (which was also the first MIA) was proposed by Shokri et al.

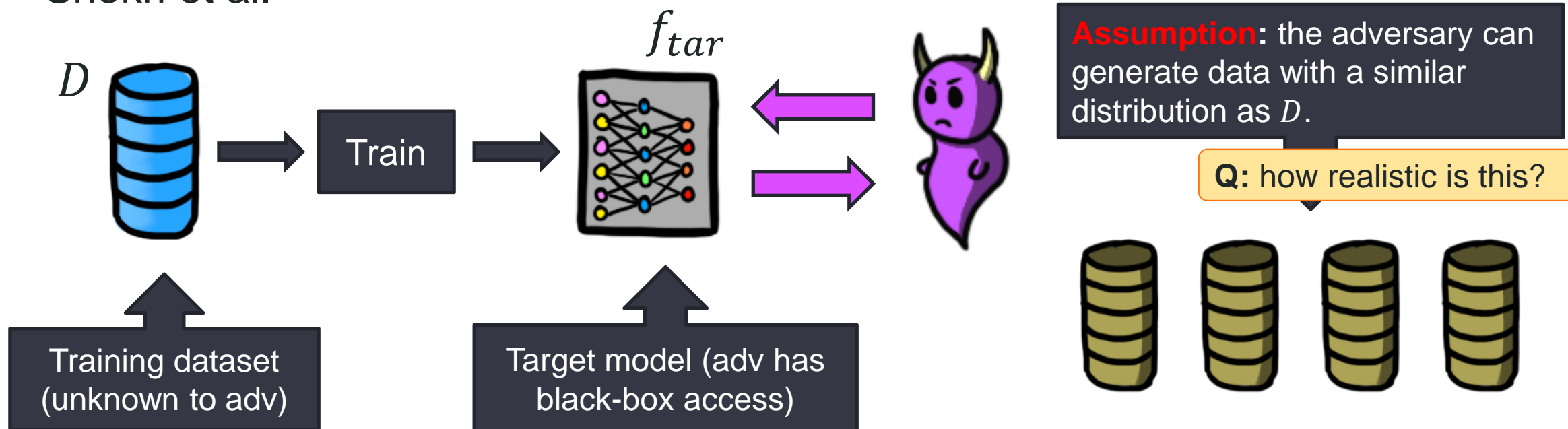


Shokri, Reza, et al. "Membership inference attacks against machine learning models." *IEEE symposium on security and privacy (SP)*, 2017.



# Neural Network-based Attacks

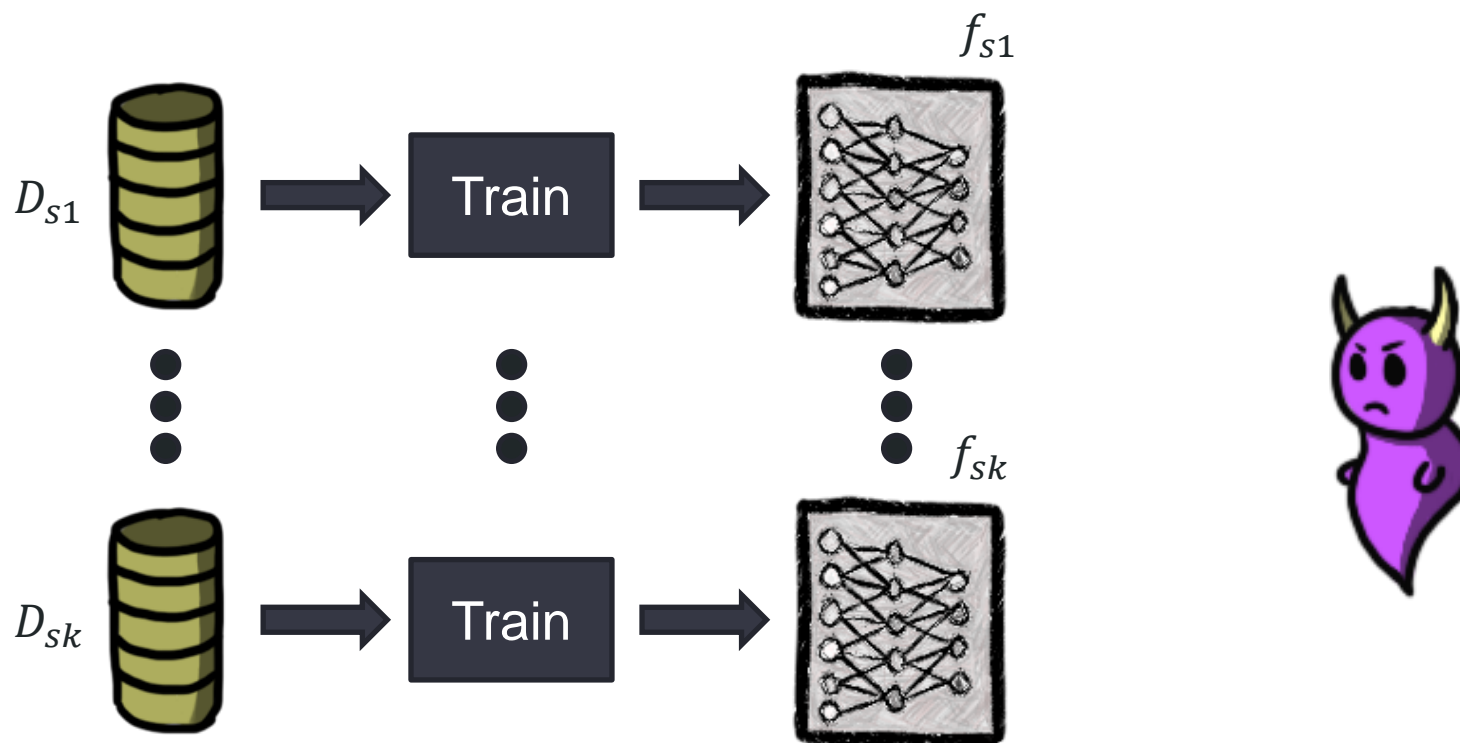
- Other MIAs use Machine Learning against Machine Learning.
- The first NN-based attack (which was also the first MIA) was proposed by Shokri et al.



Shokri, Reza, et al. "Membership inference attacks against machine learning models." *IEEE symposium on security and privacy (SP)*, 2017.

# Shokri et al.'s attack

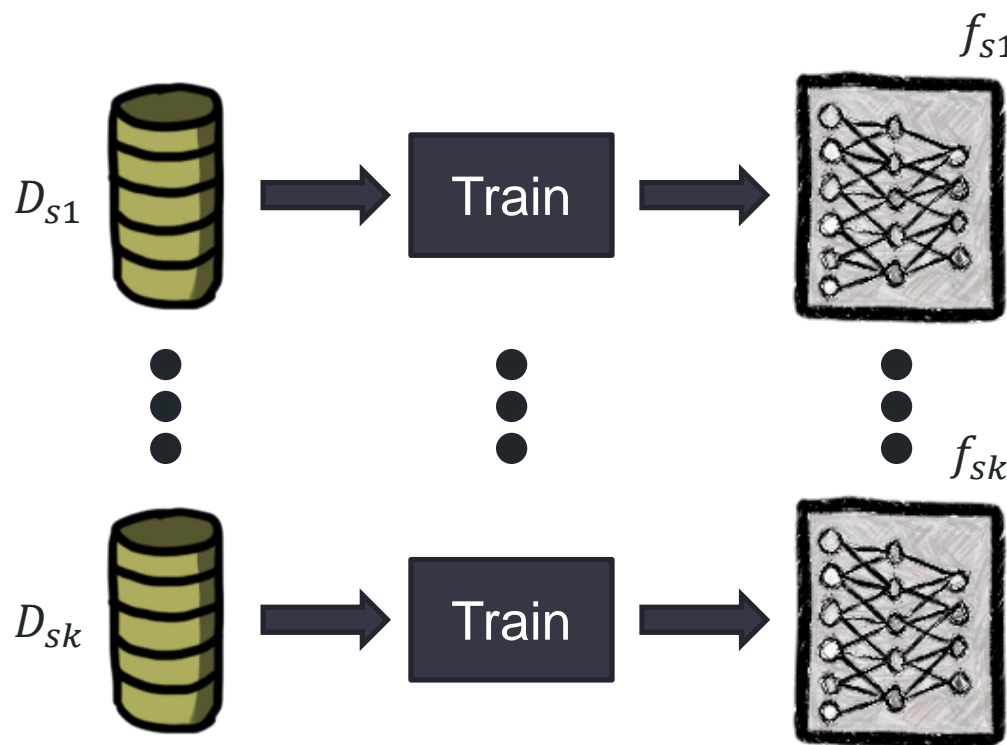
1. Generate **shadow** training datasets  $D_{s1}, D_{s2}, \dots, D_{sk}$  (based on  $D'$  with distribution similar to  $D$ ).
2. Train  $k$  **shadow models**  $f_{s1}, \dots, f_{sk}$  (same classification task as the target model "AWS MLaaS").



Shokri, Reza, et al. "Membership inference attacks against machine learning models." *IEEE symposium on security and privacy (SP)*, 2017.

# Shokri et al.'s attack

1. Generate **shadow** training datasets  $D_{s1}, D_{s2}, \dots, D_{sk}$  (based on  $D'$  with distribution similar to  $D$ ).
2. Train  $k$  **shadow models**  $f_{s1}, \dots, f_{sk}$  (same classification task as the target model "AWS MLaaS").

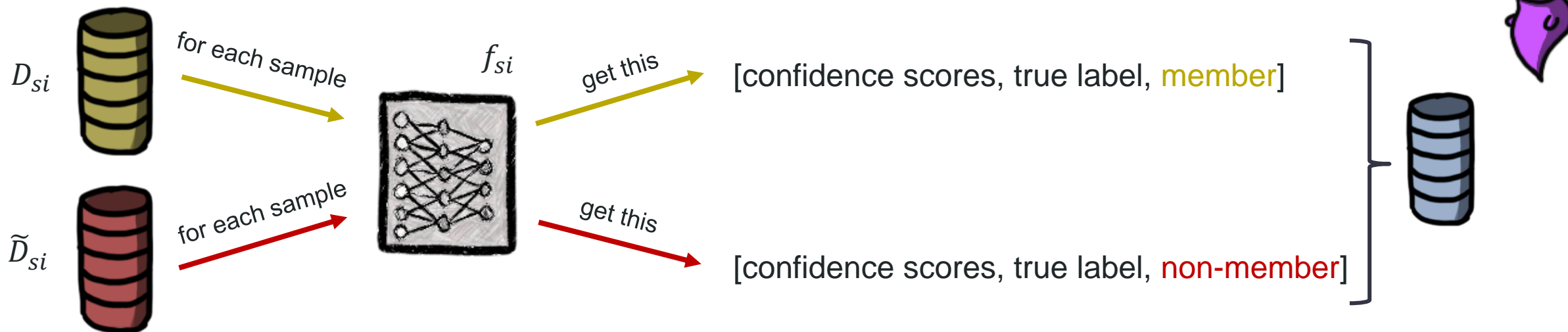


Works even with different models! (but better if you know the actual one)

Shokri, Reza, et al. "Membership inference attacks against machine learning models." *IEEE symposium on security and privacy (SP)*, 2017.

# Shokri et al.'s attack

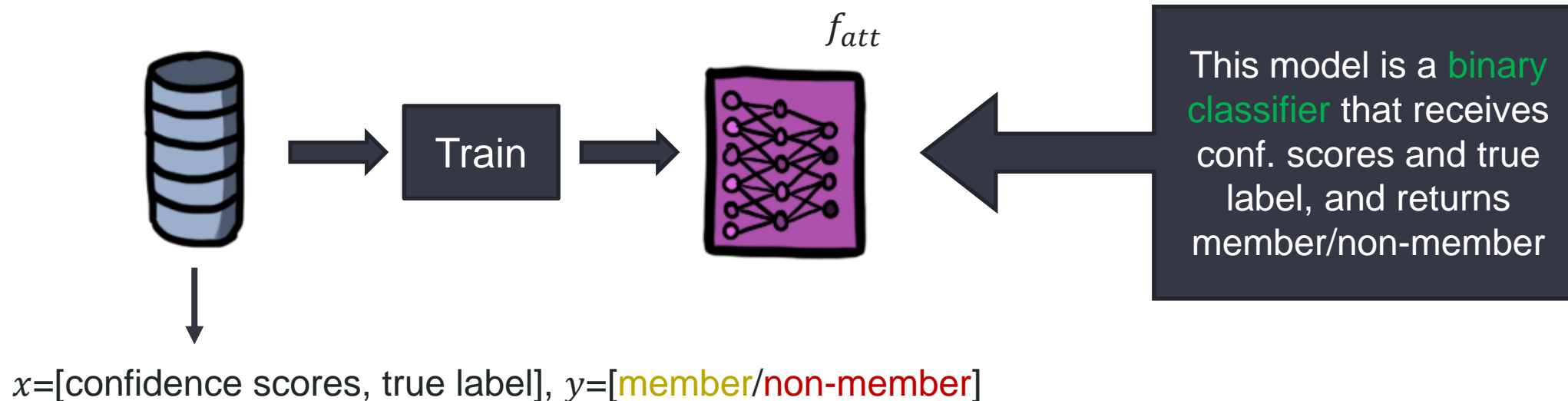
3. Generate **shadow** test data  $\tilde{D}_{s1}, \tilde{D}_{s2}, \dots, \tilde{D}_{sk}$ .
4. For each shadow model  $i \in [k]$ : get the confidence scores for each sample in  $D_{si}$  and  $\tilde{D}_{si}$ . Create a dataset with **[confidence scores, true label, membership]** for each sample.



Shokri, Reza, et al. "Membership inference attacks against machine learning models." *IEEE symposium on security and privacy (SP)*, 2017.

# Shokri et al.'s attack

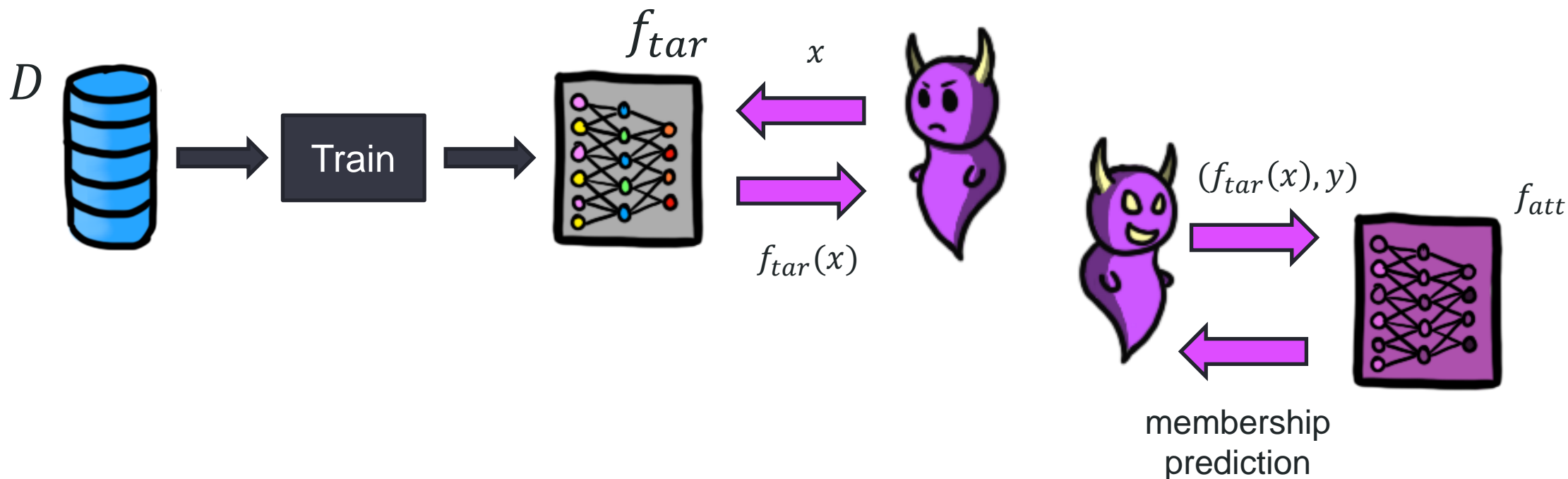
5. With the new dataset, that contains [confidence scores, true label, membership status] computed with all the shadow models, train a new **attack model**  $f_{att}$  to predict the **membership status** from [confidence scores, true label]



Shokri, Reza, et al. "Membership inference attacks against machine learning models." *IEEE symposium on security and privacy (SP)*, 2017.

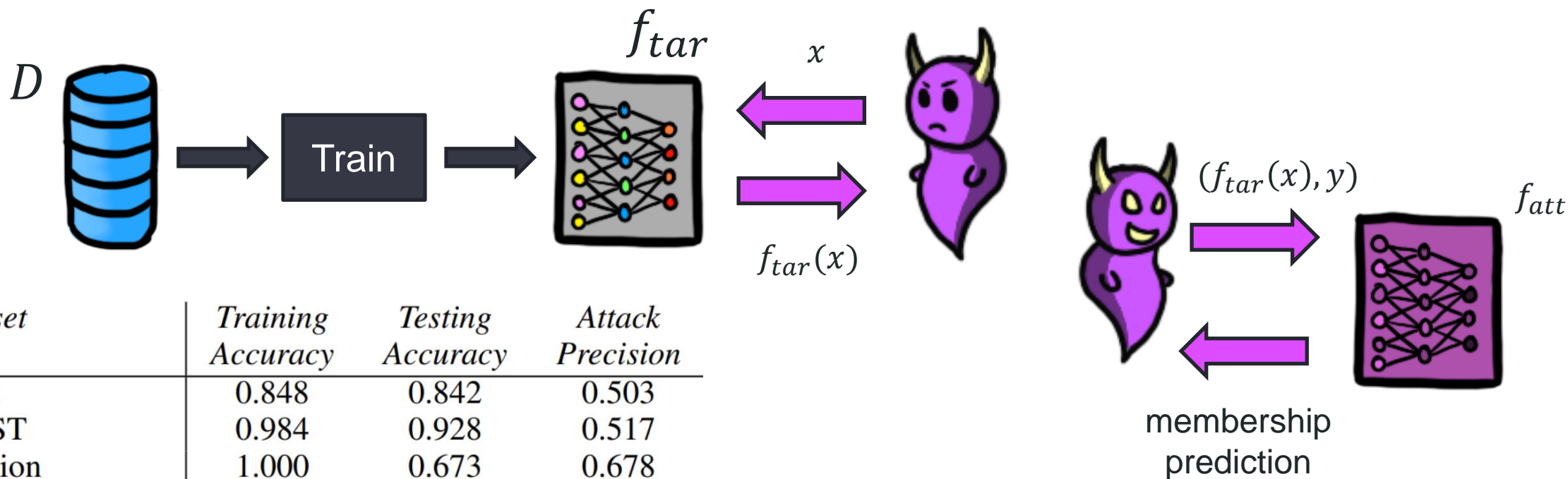
# Shokri et al.'s attack

6. Get the confidence scores of the target sample in the target model  $f_{tar}$ .
7. Evaluate those **[confidence scores, true label]** samples in the attack model  $f_{att}$ .



Shokri, Reza, et al. "Membership inference attacks against machine learning models." *IEEE symposium on security and privacy (SP)*, 2017.

# Shokri et al.'s attack

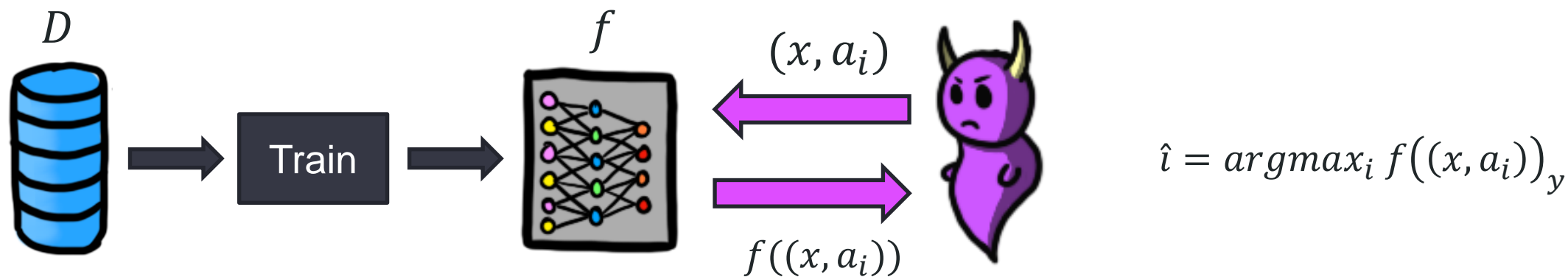


Dataset	Training Accuracy	Testing Accuracy	Attack Precision
Adult	0.848	0.842	0.503
MNIST	0.984	0.928	0.517
Location	1.000	0.673	0.678
Purchase (2)	0.999	0.984	0.505
Purchase (10)	0.999	0.866	0.550
Purchase (20)	1.000	0.781	0.590
Purchase (50)	1.000	0.693	0.860
Purchase (100)	0.999	0.659	0.935
TX hospital stays	0.668	0.517	0.657

The higher the discrepancy between training and testing accuracy  
 → The more likely membership inference attack can happen.

# Attribute Inference Attacks

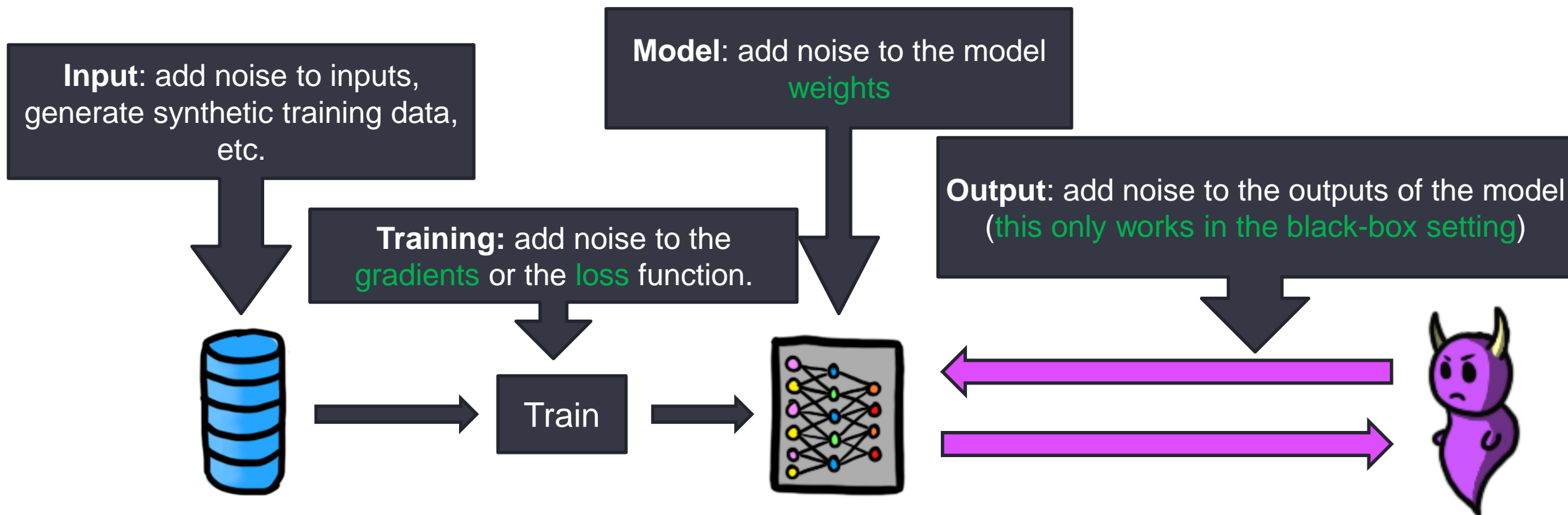
- Each sample is  $z = (x, a, y)$ , where  $x$  is the features,  $a$  is a **privacy-sensitive attribute**, and  $y$  is the label.
- The adversary has a sample  $z = (x, ?, y)$ , and wants to learn the attribute.
- Assume the space of all attributes is  $\mathcal{A} = \{a_1, a_2, \dots, a_m\}$
- Simple attack: query for all possible samples  $(x, a_1), \dots, (x, a_m)$ .  
→ The true attribute is probably the one that yields a **highest confidence** score for the true class  $y$ .





# Defending against inference attacks

- Where do we defend?



# Defenses against inference attacks

---

# Differentially Private Stochastic Gradient Descent (DP-SGD)

- Adds privacy during the **training** step, modifying SGD.
- **Recall Differential Privacy:** we want to limit the effect that **a single training set sample** has on the output (the “**output**” of the training algorithm is the “**model!**”)

## SGD

For each training step  $t \in [T]$ :

1. Take a batch  $B$  of  $L$  samples from  $D$ .
2. For each  $(x_i, y_i) \in B$ , compute the gradient:

$$g_i = \nabla \ell(\theta_{t-1}, x_i, y_i)$$

3. Average the gradients  $g = \frac{1}{L} \sum_i g_i$ .
4. Descend  $\theta_t = \theta_{t-1} - \eta \cdot g$ .

## “Private” SGD

For each training step  $t \in [T]$ :

1. Take a batch  $B$  of  $L$  samples from  $D$ .
2. For each  $(x_i, y_i) \in B$ , compute the gradient:

$$g_i = \nabla \ell(\theta_{t-1}, x_i, y_i)$$

3. Average the gradients **and add noise**  $g = \frac{1}{L} (\sum_i g_i + \mathcal{N}(0, \sigma^2))$ .
4. Descend  $\theta_t = \theta_{t-1} - \eta \cdot g$ .

**Q:** Is it enough to add noise to the gradients?

# Differentially Private Stochastic Gradient Descent (DP-SGD)

- The gradient could potentially be **unbounded** → Here, unbounded sensitivity is bad for DP (Algorithm is highly sensitive to individual data points)
- We **clip** the gradients to ensure their  $\ell_2$  norm is at most  $C$ .
  - $C$  is the *clipping threshold* (1 is usually a good value)
  - $C$  is independent of the data

## SGD

For each training step  $t \in [T]$ :

1. Take a batch  $B$  of  $L$  samples from  $D$ .
2. For each  $(x_i, y_i) \in B$ , compute the gradient:

$$g_i = \nabla \ell(\theta_{t-1}, x_i, y_i)$$

3. Average the gradients  $g = \frac{1}{L} \sum_i g_i$ .
4. Descend  $\theta_t = \theta_{t-1} - \eta \cdot g$ .

## DP-SGD

For each training step  $t \in [T]$ :

1. Take a batch  $B$  of  $L$  samples from  $D$ .
2. For each  $(x_i, y_i) \in B$ , compute the gradient:

$$g_i = \nabla \ell(\theta_{t-1}, x_i, y_i)$$

3. Clip the gradients:  $g_i = g_i / \max\left(1, \frac{\|g_i\|_2}{C}\right)$

4. Sum the gradients  $g = \sum_i g_i$ .

5. Add noise:  $g = g + \mathcal{N}(0, \sigma^2 C^2)$

6. Descend  $\theta_t = \theta_{t-1} - \eta \cdot \frac{1}{L} g$ .

# Differentially Private Stochastic Gradient Descent (DP-SGD)

- The gradient could potentially be **unbounded** → bad for DP
- We **clip** the gradients to ensure  $\|g_i\|_2 \leq C$  norm:
  - $C$  is the *clipping threshold* (1 is used)
  - $C$  is independent of the data

**NOTE:**  $C$  should be **large enough** to avoid significantly impeding the learning process, but **small enough** to ensure that the gradients are not too sensitive to individual data points.

For each training step  $t$ :

1. Take a batch  $B$  of  $n$  samples from  $D$ .
2. For each  $(x_i, y_i) \in B$ , compute the gradient:  
 $g_i = \nabla \ell(\theta_{t-1}, x_i, y_i)$

$$g_i = \nabla \ell(\theta_{t-1}, x_i, y_i)$$

3. Average the gradients  $g = \frac{1}{L} \sum_{i=1}^L g_i$ .
4. Descend  $\theta_t = \theta_{t-1} - \eta \cdot g$ .

Of  $L$  samples from  $D$ .

$(x_i, y_i) \in B$ , compute the gradient:  
 $g_i = \nabla \ell(\theta_{t-1}, x_i, y_i)$

Clip the gradients:  $g_i = g_i / \max\left(1, \frac{\|g_i\|_2}{C}\right)$

Sum the gradients  $g = \sum_i g_i$ .

Add noise:  $g = g + \mathcal{N}(0, \sigma^2 C^2)$

Descend  $\theta_t = \theta_{t-1} - \eta \cdot \frac{1}{L} g$ .

# DP-SGD: keeping track of $\epsilon, \delta$

- Note that a single sample will participate in multiple training steps  $\rightarrow$  there will be some **sequential composition** involved.
- We need to **keep track of  $\epsilon, \delta$** .
- For a **fixed amount of noise  $\sigma$** , if we do not keep track of  $\epsilon, \delta$ , we can end up with a very large  $\epsilon$ , which is bad.
  - The actual *true*  $\epsilon$  will be **smaller** than the  $\epsilon$  we can compute theoretically. – e.g., due to batching, one sample may not appear in a given training step.
  - We can only guarantee an  $\epsilon$  we can prove w/ theory.

## DP-SGD

For each training step  $t \in [T]$ :

1. Take a batch  $B$  of  $L$  samples from  $D$ .
2. For each  $(x_i, y_i) \in B$ , compute the gradient:

$$g_i = \nabla \ell(\theta_{t-1}, x_i, y_i)$$

3. **Clip the gradients:**  $g_i = g_i / \max\left(1, \frac{\|g_i\|_2}{c}\right)$
4. Sum the gradients  $g = \sum_i g_i$ .
5. **Add noise:**  $g = g + \mathcal{N}(0, \sigma^2 C^2)$
6. Descend  $\theta_t = \theta_{t-1} - \eta \cdot \frac{1}{L} g$ .

# DP-SGD: keeping track of $\epsilon, \delta$

- First, we choose a  $\delta$ . Recall that this should be smaller than  $\delta < \frac{1}{N}$ .
  - The reason is the following: a training algorithm that simply publishes a random training set record would provide  $(\epsilon = 0, \delta = 1/N)$ -DP. However, we know this is not private enough.

## DP-SGD

For each training step  $t \in [T]$ :

1. Take a batch  $B$  of  $L$  samples from  $D$ .
2. For each  $(x_i, y_i) \in B$ , compute the gradient:

$$g_i = \nabla \ell(\theta_{t-1}, x_i, y_i)$$

3. **Clip the gradients:**  $g_i = g_i / \max\left(1, \frac{\|g_i\|_2}{c}\right)$
4. Sum the gradients  $g = \sum_i g_i$ .
5. **Add noise:**  $g = g + \mathcal{N}(0, \sigma^2 C^2)$
6. Descend  $\theta_t = \theta_{t-1} - \eta \cdot \frac{1}{L} g$ .

# DP-SGD: keeping track of $\epsilon, \delta$

**Q:** Given  $\delta, \sigma, C, T$ , and assuming each sample in  $D$  is used *once per training step*, what is the total  $\epsilon$  we get?

- Use naive composition

$\Delta_2^2 = C^2$  : second-order sensitivity

Typically refers to the maximum **L<sub>2</sub> norm** of the gradients, and it is directly related to the **gradient clipping bound**

## DP-SGD

For each training step  $t \in [T]$ :

1. Take a batch  $B$  of  $L$  samples from  $D$ .
2. For each  $(x_i, y_i) \in B$ , compute the gradient:

$$g_i = \nabla \ell(\theta_{t-1}, x_i, y_i)$$

3. **Clip the gradients:**  $g_i = g_i / \max\left(1, \frac{\|g_i\|_2}{c}\right)$
4. Sum the gradients  $g = \sum_i g_i$ .
5. **Add noise:**  $g = g + \mathcal{N}(0, \sigma^2 C^2)$
6. Descend  $\theta_t = \theta_{t-1} - \eta \cdot \frac{1}{L} g$ .

$f(D) + Y$  is  $(\epsilon, \delta)$ -DP if

$$Y \sim N(0, \sigma^2)$$

$$\sigma^2 = 2 \ln\left(\frac{1.25}{\delta}\right) \Delta_2^2 / \epsilon^2$$



# DP-SGD: keeping track of $\epsilon, \delta$

**Q:** Given  $\delta, \sigma, C, T$ , and assuming each sample in  $D$  is used *once per training step*, what is the total  $\epsilon$  we get?

- Use naive composition

**A:**  $\sigma^2 = 2 \ln \left( \frac{1.25}{\delta} \right) \Delta_2^2 / \epsilon^2 \rightarrow \epsilon_{step} = \sqrt{2 \ln \left( \frac{1.25}{\delta} \right) C^2 / \sigma^2}$

for each step. Then naïve composition gives

$$\epsilon_{total} = T \cdot C / \sigma \sqrt{2 \ln \left( \frac{1.25}{\delta} \right)}$$

\*Note: this question is very over simplified

## DP-SGD

For each training step  $t \in [T]$ :

1. Take a batch  $B$  of  $L$  samples from  $D$ .
2. For each  $(x_i, y_i) \in B$ , compute the gradient:

$$g_i = \nabla \ell(\theta_{t-1}, x_i, y_i)$$

3. Clip the gradients:  $g_i = g_i / \max \left( 1, \frac{\|g_i\|_2}{c} \right)$
4. Sum the gradients  $g = \sum_i g_i$ .
5. Add noise:  $g = g + \mathcal{N}(0, \sigma^2 C^2)$
6. Descend  $\theta_t = \theta_{t-1} - \eta \cdot \frac{1}{L} g$ .

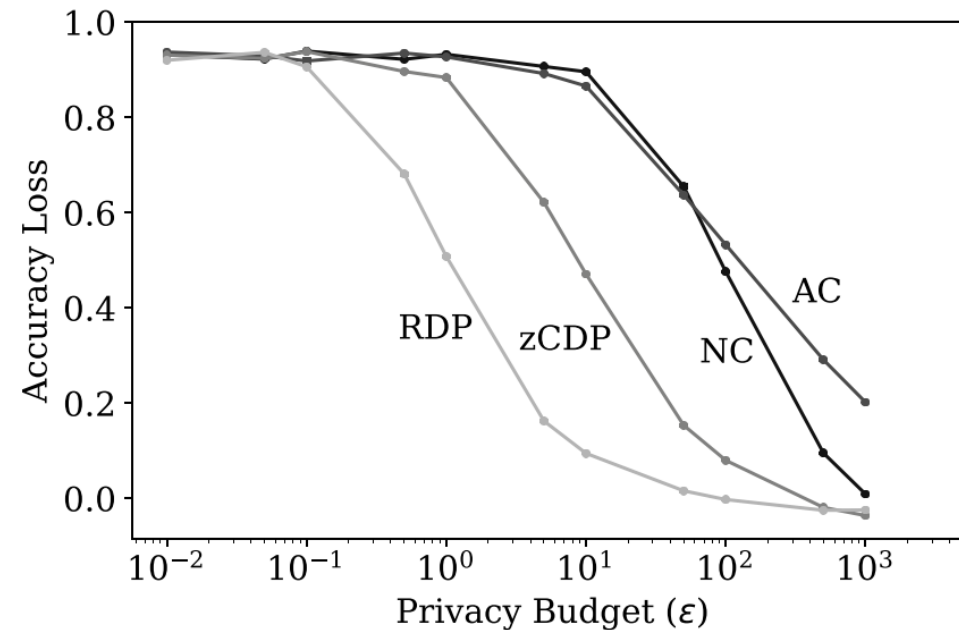
$f(D) + Y$  is  $(\epsilon, \delta)$ -DP if

$$Y \sim N(0, \sigma^2)$$

$$\sigma^2 = 2 \ln \left( \frac{1.25}{\delta} \right) \Delta_2^2 / \epsilon^2$$

# DP-SGD: keeping track of $\epsilon, \delta$

- Renyi Differential Privacy (RDP) provides a tighter  $\epsilon, \delta$  bound.
  - Better suited to Gaussian Noise
  - Keeps track of more information
- This means that, for a given  $\sigma, C$ , and  $\delta$ , RDP tells us our actual  $\epsilon$  is smaller than what Advanced Composition (AC) tells us.
- **In other words**, for a target privacy budget  $\epsilon$ , using RDP we need to add less noise than using AC.
  - E.g., again, because a sample may be excluded from a given training step
- Note that, even with RDP, we need  $\epsilon > 100$  if we **do not want any accuracy loss**

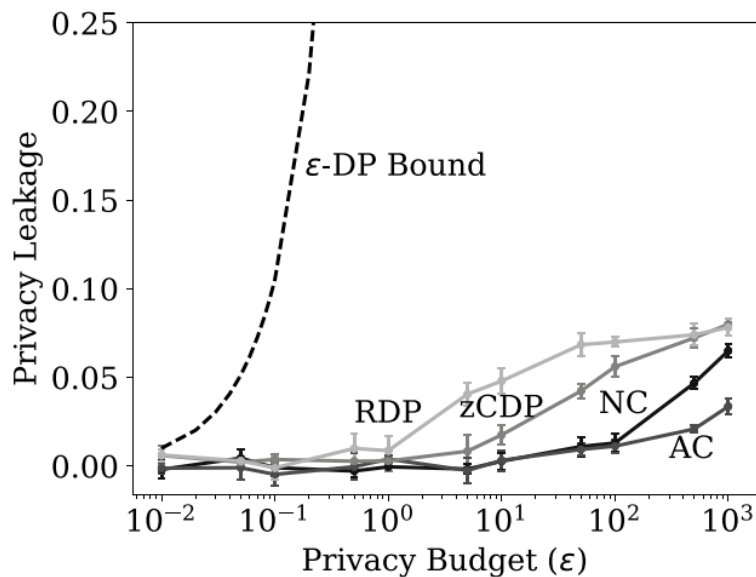


zCDP: Zero-Concentrated Differential Privacy  
NC: Noise-Contrastive privacy  
AC: Approximate Composition

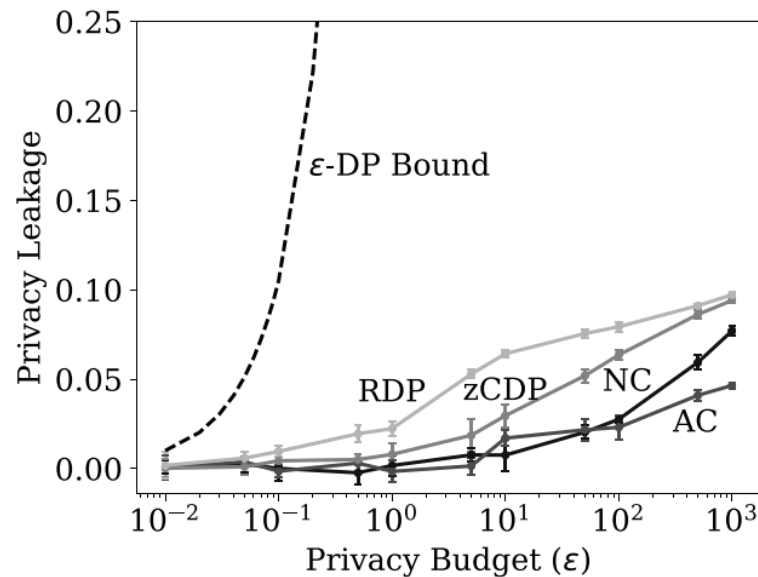
Jayaraman, Bargav, and David Evans. "Evaluating differentially private machine learning in practice." *USENIX Security Symposium*. 2019.

# DP-SGD: theoretical vs empirical privacy

- Both attacks we've seen perform similarly
- It seems that  $\epsilon = 100$  or even  $\epsilon = 1000$  still provides good empirical privacy
- The theoretical bound “worst-case” on the privacy leakage provided by DP is very loose “overestimates the actual risk”.



(a) Shokri et al. membership inference



(b) Yeom et al. membership inference

Jayaraman, Bargav, and David Evans. "Evaluating differentially private machine learning in practice." *USENIX Security Symposium*. 2019.

# Issues of DP-SGD

---

- We saw that, for strong theoretical privacy (e.g.,  $\epsilon < 1$ ), the models usually **lose all utility**.
- For very weak theoretical privacy (e.g.,  $\epsilon = 100$ ), some models achieve reasonable utility.
- However, DP-SGD with  $\epsilon = 100$  seems to provide enough protection against existing attacks.

**Q:** Is it OK to use  $\epsilon = 100$ ?

# Issues of DP-SGD

---

- We saw that, for strong theoretical privacy (e.g.,  $\epsilon < 1$ ), the models usually **lose all utility**.
- For very weak theoretical privacy (e.g.,  $\epsilon = 100$ ), some models achieve reasonable utility.
- However, DP-SGD with  $\epsilon = 100$  seems to provide enough protection against existing attacks.

**Q:** Is it OK to use  $\epsilon = 100$ ?

**A:** It might be OK to use DP-SGD tuned to  $\epsilon = 100$ , but at that point we might as well use defenses that do not provide DP, since the DP guarantee is already meaningless at that point.

# Private Aggregation of Teacher Ensembles (PATE)

1. Train teacher models with disjoint subsets of the training data
2. Use the teachers to label some (incomplete) public data
3. Use the labeled public data to train a student model

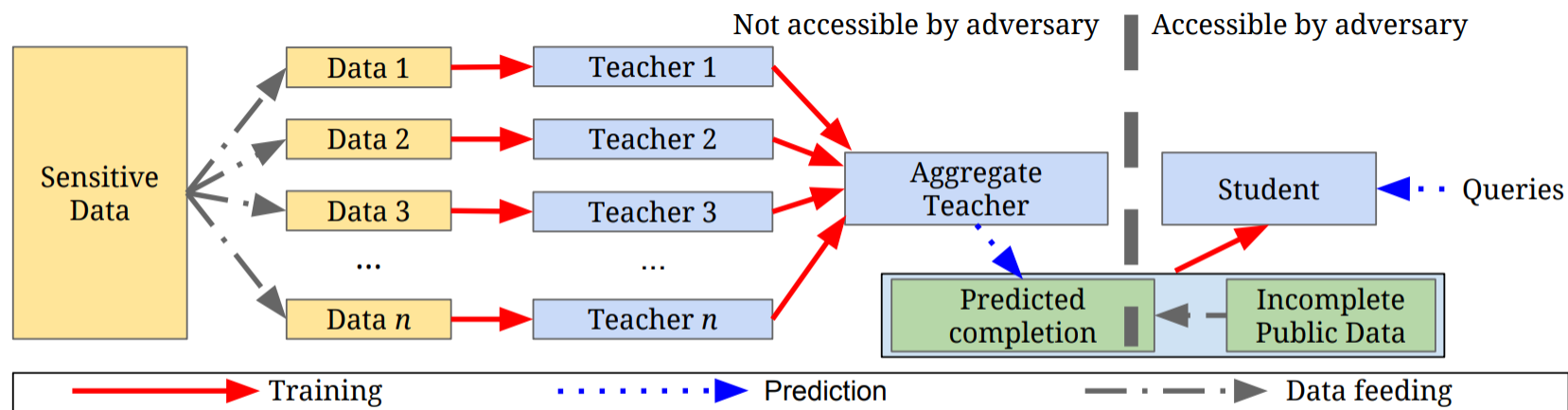


Figure 1: Overview of the approach: (1) an ensemble of teachers is trained on disjoint subsets of the sensitive data, (2) a student model is trained on public data labeled using the ensemble.

# Private Aggregation of Teacher Ensembles (PATE)

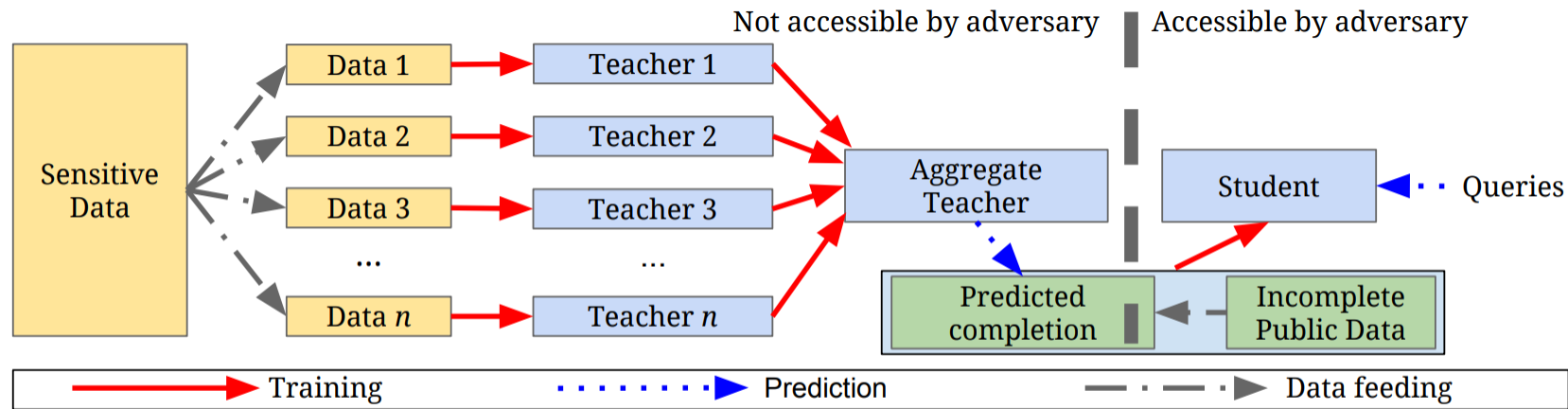


Figure 1: Overview of the approach: (1) an ensemble of teachers is trained on disjoint subsets of the sensitive data, (2) a student model is trained on public data labeled using the ensemble.

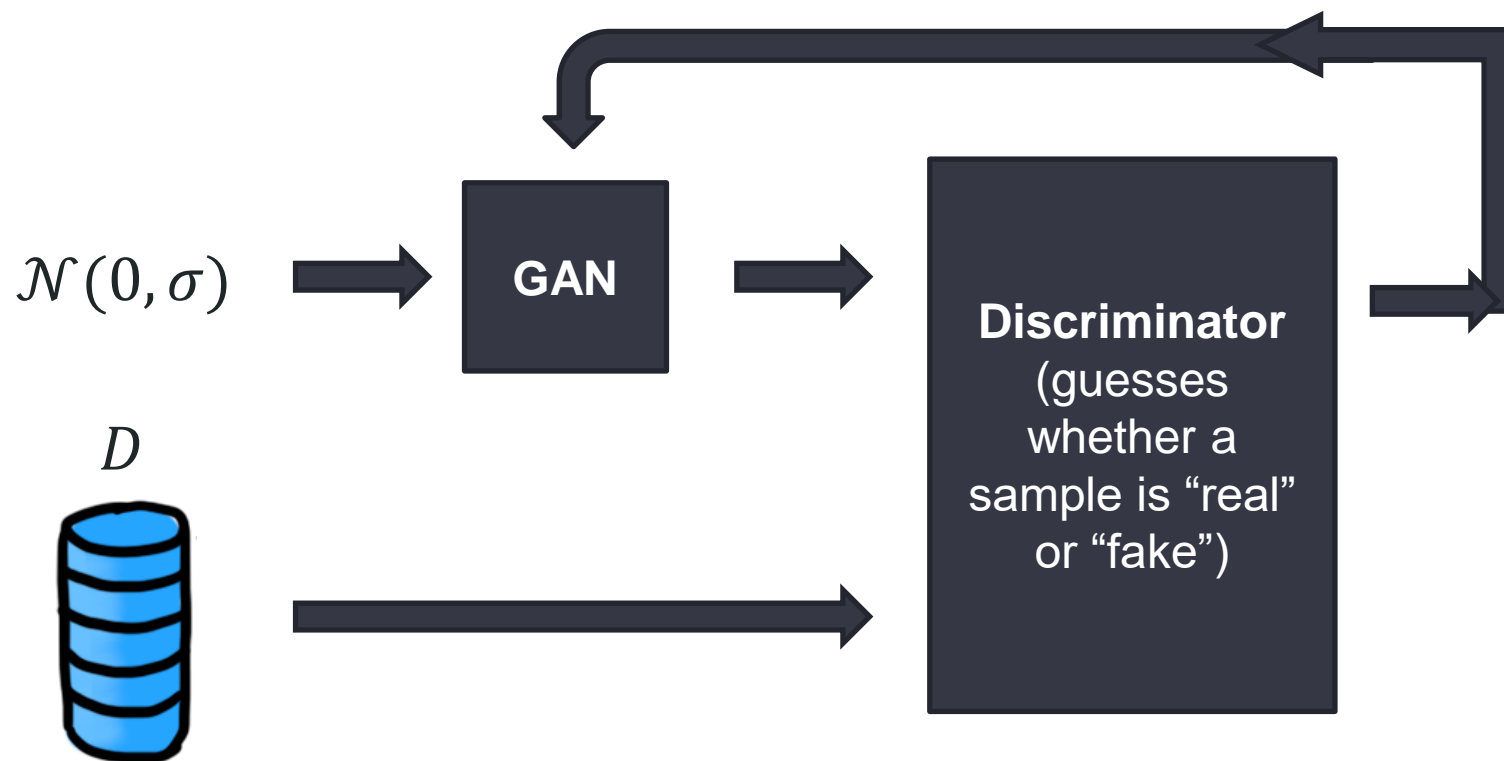
- For a sample from the **incomplete public data**  $\vec{x}$ , let  $n_j(\vec{x})$  be the number of teachers that **voted** for label  $j$ .
- Instead of labeling by taking  $\operatorname{argmax}_j \{n_j(\vec{x})\}$ , we can add Laplacian noise to provide DP:

$$\operatorname{argmax}_j \left\{ n_j(\vec{x}) + \operatorname{Lap} \left( \frac{1}{\gamma} \right) \right\}$$

Papernot, Nicolas, et al. "Semi-supervised knowledge transfer for deep learning from private training data." *ICLR 2017*

# Synthetic Data Generation

- For example, by using a GAN to generate real-looking synthetic samples:

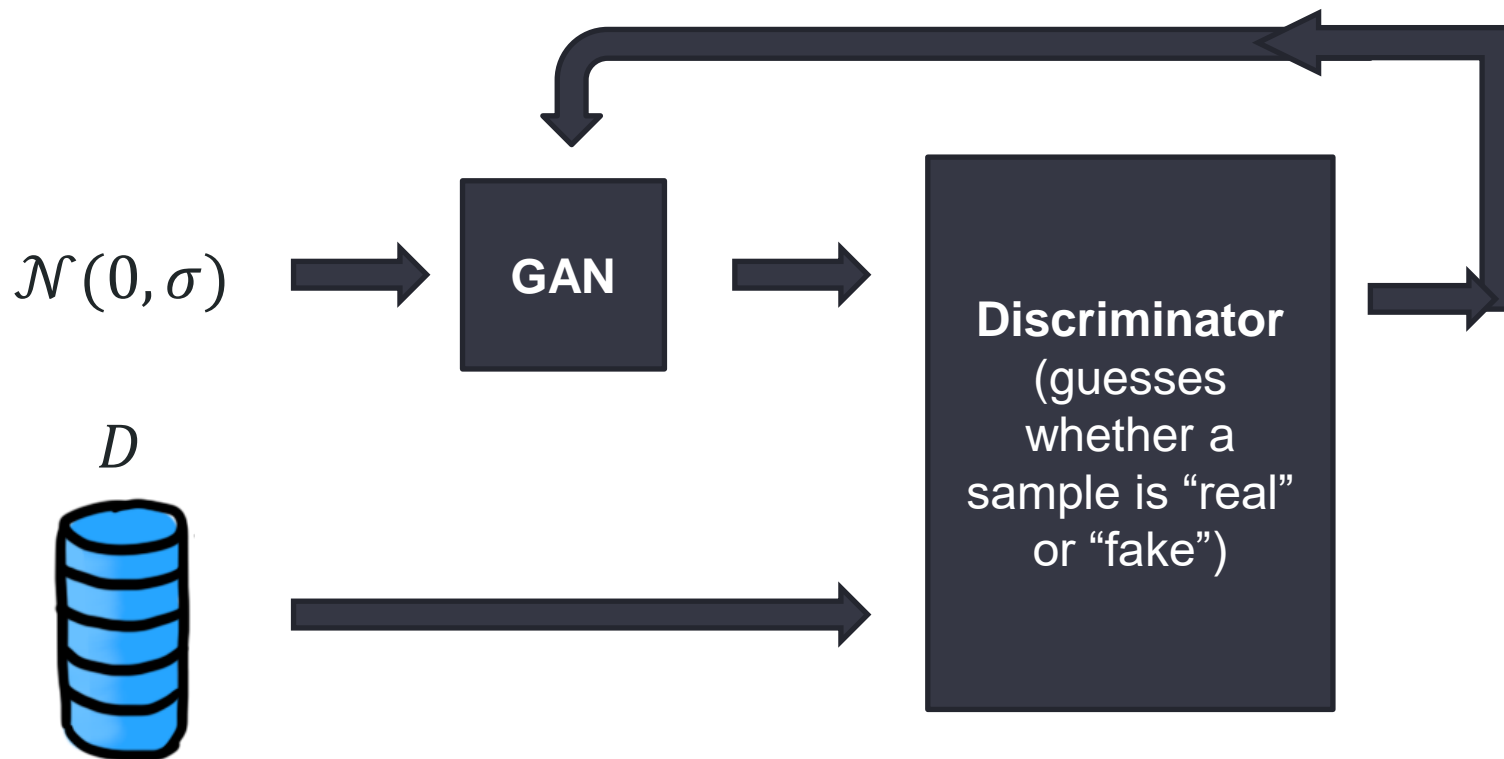


If we train the GAN using privacy-preserving training algorithms (e.g., DP-SGD on the discriminator), we can use it to generate a privacy-preserving synthetic dataset!



# Synthetic Data Generation

- For example, by using a GAN to generate real-looking synthetic samples:

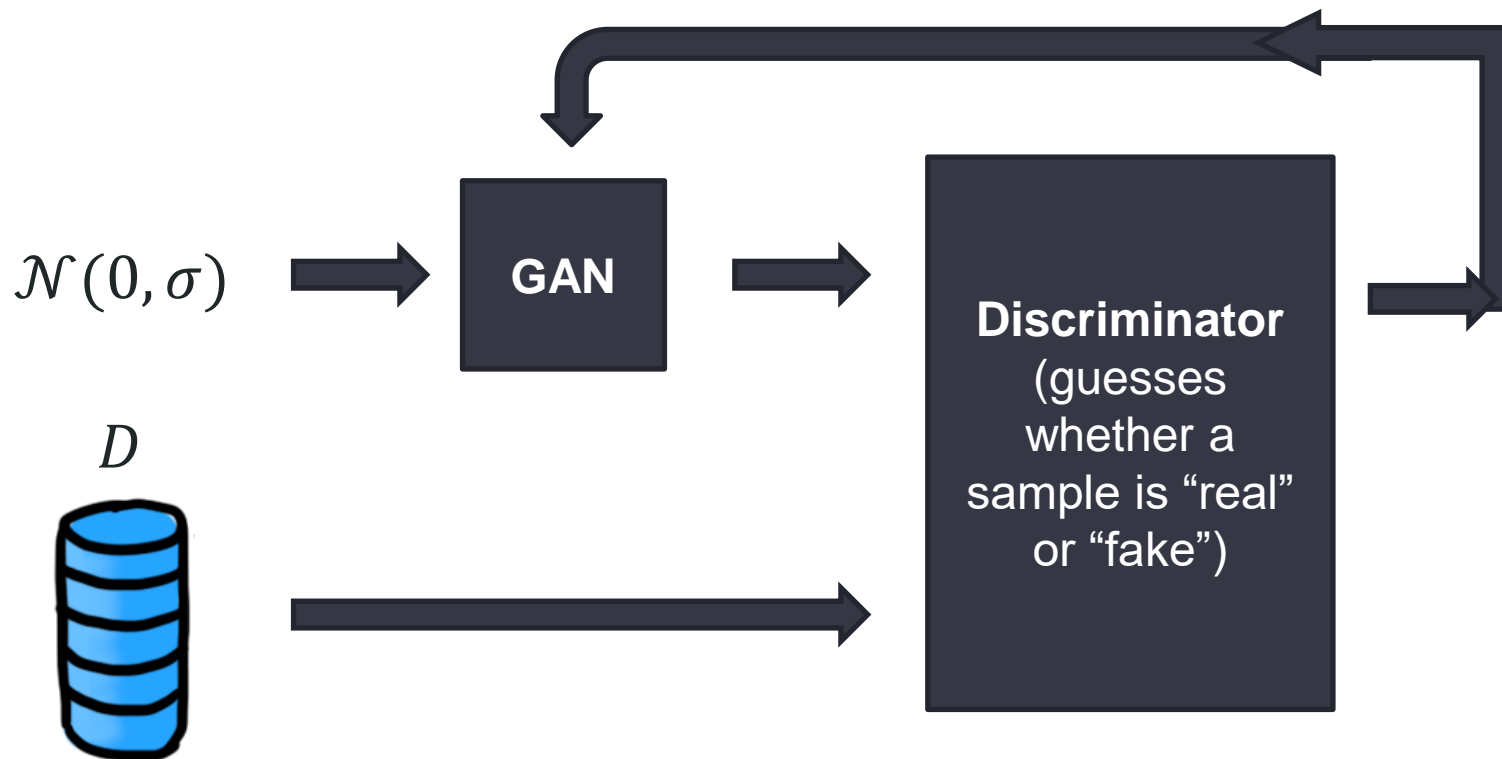


If we train the GAN using privacy-preserving training algorithms (e.g., DP-SGD on the discriminator), we can use it to generate a privacy-preserving synthetic dataset!

**Q:** What can we do with the resulting dataset?

# Synthetic Data Generation

- For example, by using a GAN to generate real-looking synthetic samples:



If we train the GAN using privacy-preserving training algorithms (e.g., DP-SGD on the discriminator), we can use it to generate a privacy-preserving synthetic dataset!

**Q:** What can we do with the resulting dataset?

**A:** Anything!

# Other defenses

---

- There are defenses that add **noise** to the confidence scores (MemGuard [Jia et al.]), but are **not very effective**.
- MIAs can work even if the model **just leaks the predicted label** (and not the confidence scores)
- Sometimes, **generalization** is a good defense by itself:
  - A well-generalized model will perform similarly in members (training set) and non-members (testing set)
  - Therefore, **it will be harder** for an adversary to decide whether a sample is a member or non-member if the model generalizes well.
  - Generalization is also **good for utility** (improves test accuracy), so it's a win-win.