



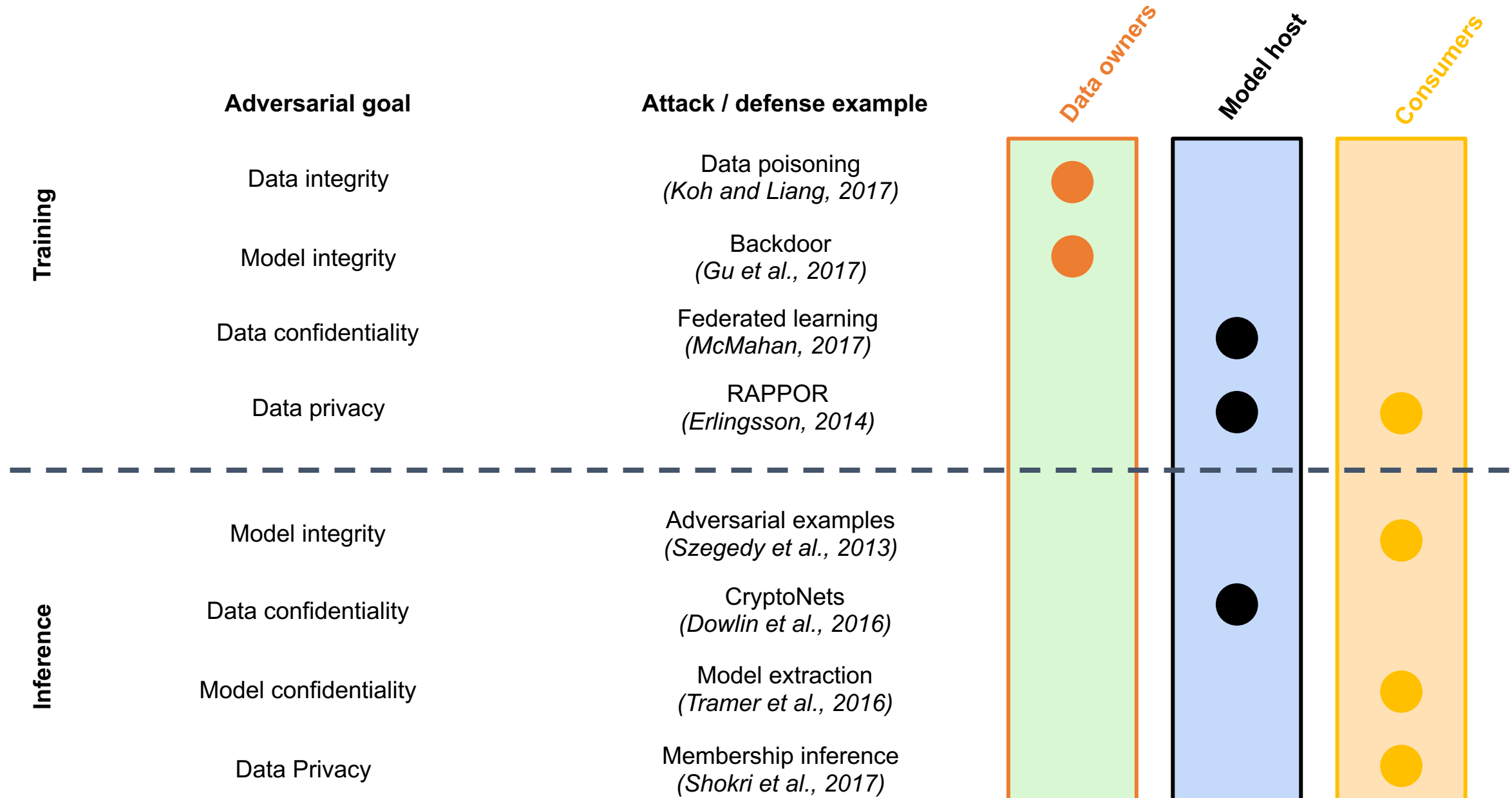
UNIVERSITY OF  
TORONTO



# Lecture 7: Differential Privacy

Oct 28

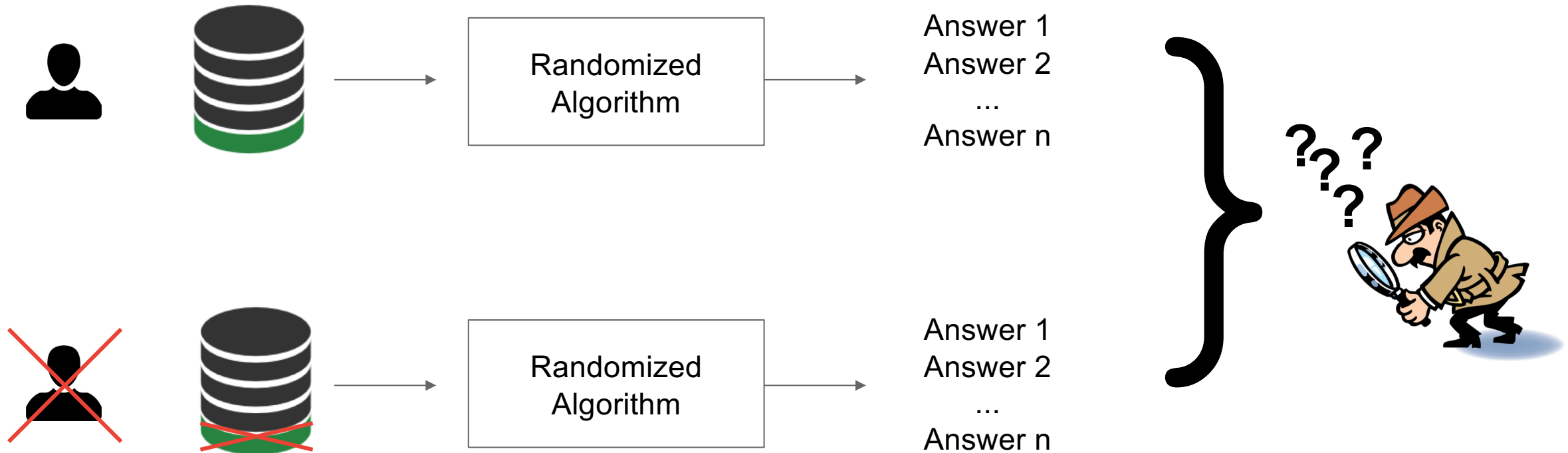
Prof. Nicolas Papernot



# Limitations of previous definitions: the case of k-anonymity

- Each record must be indistinguishable from  $k-1$  other records
  - Suppression -> replace features by wildcards
  - Generalization -> change age from number of years to bins
- Attacks:
  - Often use background knowledge
  - E.g., link attributes in private database and attributes from another database

# What is a differentially private algorithm?



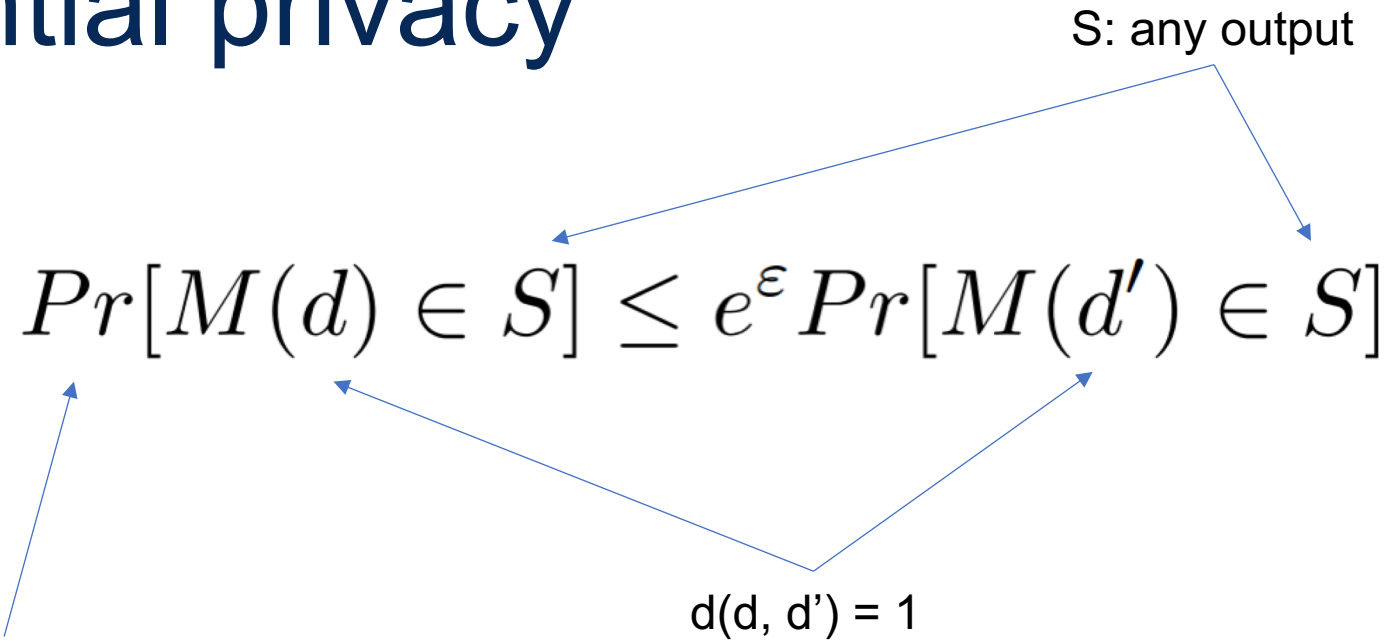


# Differential privacy

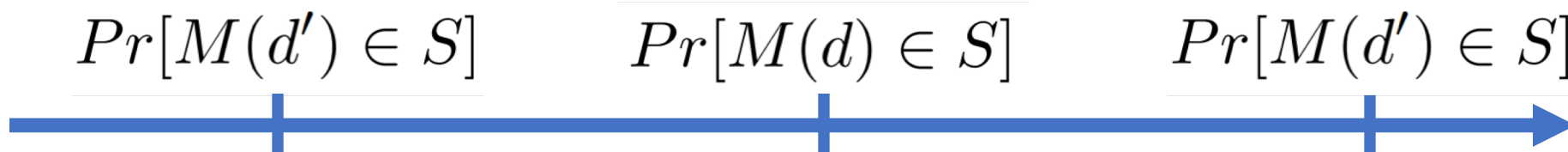
$$Pr[M(d) \in S] \leq e^\epsilon Pr[M(d') \in S]$$

S: any output

$d(d, d') = 1$



Probability (algorithm M is randomized)



# Why DP improves upon previous definitions

- Made assumptions about adversaries:
  - Value of  $k$  in  $k$ -anonymity depends on capabilities of adversary
  - Instead DP guarantee does not depend on:
    - What adversary knows (capability)
    - What adversary wants (goal)
- Precise metric for privacy leakage (bound on epsilon)

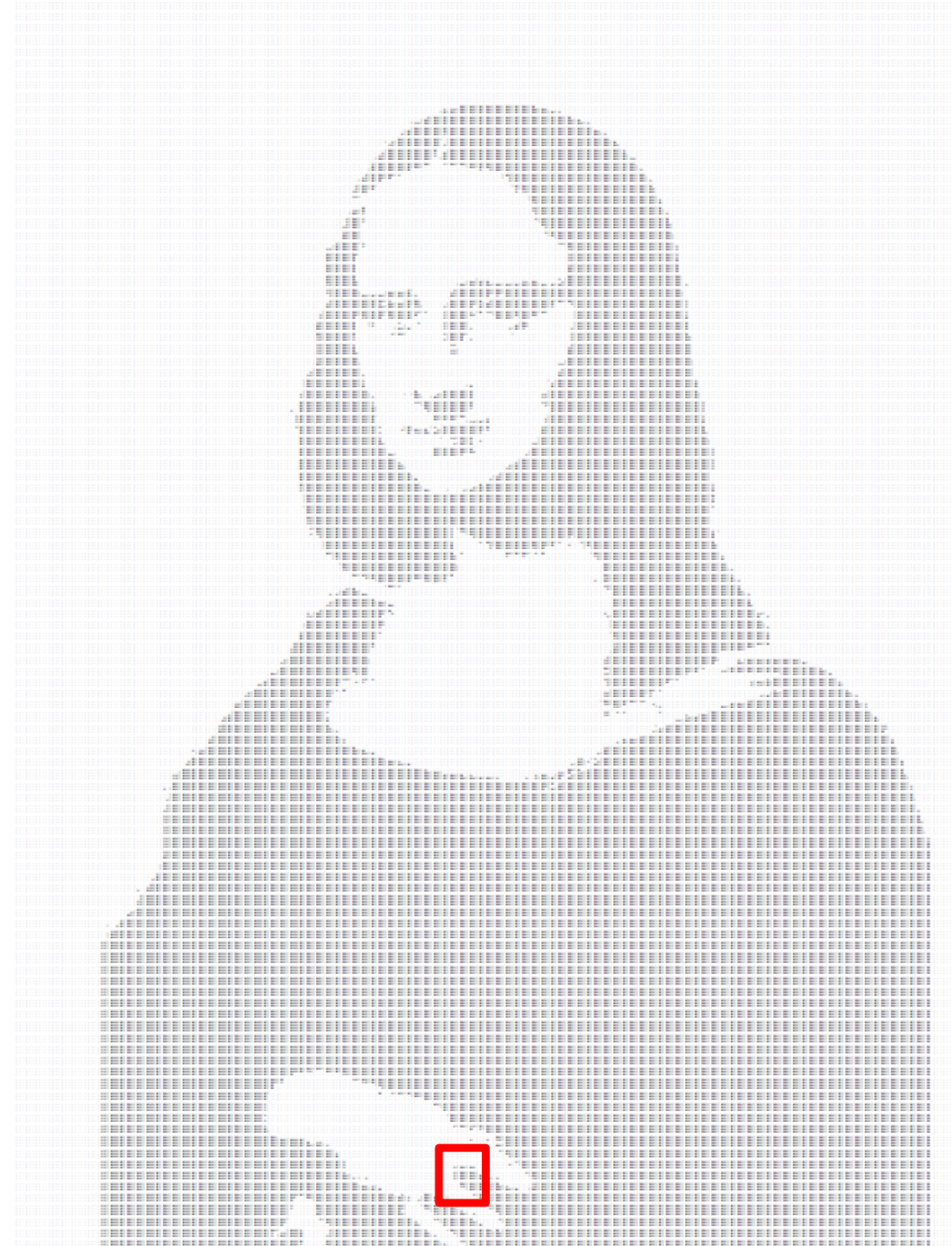
$$\Pr[M(d) \in S] \leq e^\epsilon \Pr[M(d') \in S]$$

- Robust to composition
  - Algorithm M1 has  $\epsilon$  DP
  - Algorithm M2 has  $\epsilon$  DP
  - Algorithms M1 and M2 have  $2\epsilon$  DP
- Group guarantees

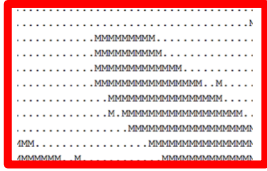
# What does that mean for a user?

- Pessimistic perspective: privacy is already lost
- DP moves forward by estimating cost of participating in a dataset
  - > *differential* privacy

# A Metaphor For Private Learning



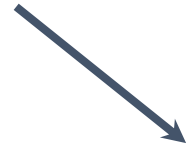
# An Individual's Training Data



.....M  
.....  
.....MMMMMMMMM.....  
.....MMMMMMMMMMM.....  
.....MMMMMMMMMMMMMMMMM.....  
.....MMMMMMMMMMMMMMMMMMMMM.....M.....  
.....MMMMMMMMMMMMMMMMMMMMMMMMM.....  
.....M.....MMMMMMMMMMMMMMMMMMMMMMMMM.....  
.....MMMMMMMMMMMMMMMMMMMMMMMMMMMMM.....  
1MM.....MMMMMMMMMMMMMMMMMMMMMMMMM.....  
1MMMMMM.....M.....MMMMMMMMMMMMMMMMMMMMMMMMM.....

# An Individual's Training Data

Each bit is flipped with  
probability  
50%



```

. . . . . M . . . . . MM . M . . . . . MMM . M . .
. . . . . . . . . . . . . . . . . . . . . . MM . . . MMMM . . .
. . . M . . MM . MM . . MMM . M . MM . M . . . M . . MM . .
. MM . . . . . MMM . . . . . MMMMMMMMMM . . . M . . . MM
. . M . . . . M . . . . . . . . . . MM . . MMMMMMMM . . . M . . .
M . . . . . M . . MM . MMMMMMMMMMMMMMMMMMMM . . . . M
. . . . . M . . . . . M . M . M . MMMMMM . . . MMMMM . . .
. . . M . . . . . M . MM . M . MM . . M . . M . . MM . MMMMM
M . . . M . M . . . . . M . M . . M . . MMM . MMMMM . MMMM
. MMM . M . . . . M . M . M . . . . . . . . . . MMMMMMMMMM . M

```



# Big Picture Remains!

# Are you a communist?

Algorithm:

1. Flip a first coin
2. If:
  - a. First coin was heads -> return correct answer
  - b. First coin was tails, flip second coin:
    - a. report true if heads
    - b. report false if tails

Plausible deniability

Is it still useful? What did you learn?



# Result of survey

- If person is communist:
  - With probability \_\_\_ they will respond correctly True
  - With probability \_\_\_ they will respond with the second coin flip
    - With probability \_\_\_ the second coin flip will return True
    - With probability \_\_\_ the second coin flip will return False
- Probability to say True \_\_\_
- Probability to say False \_\_\_
- Repeat exercise for a non-communist

# How private is our survey?

- Eps is such that  $0.75 = e^{\text{eps}} * 0.25$
- $\text{Eps} = \ln(3) \approx 1.1$
- If we changed probability of first coin flip to 75% saying true:
  - Eps is now such that  $0.75 + 0.25*0.5 = 0.875 = e^{\text{eps}} * 0.125$
  - $\text{Eps} = \ln(7) \approx 1.95$

# How to implement the survey in practice?

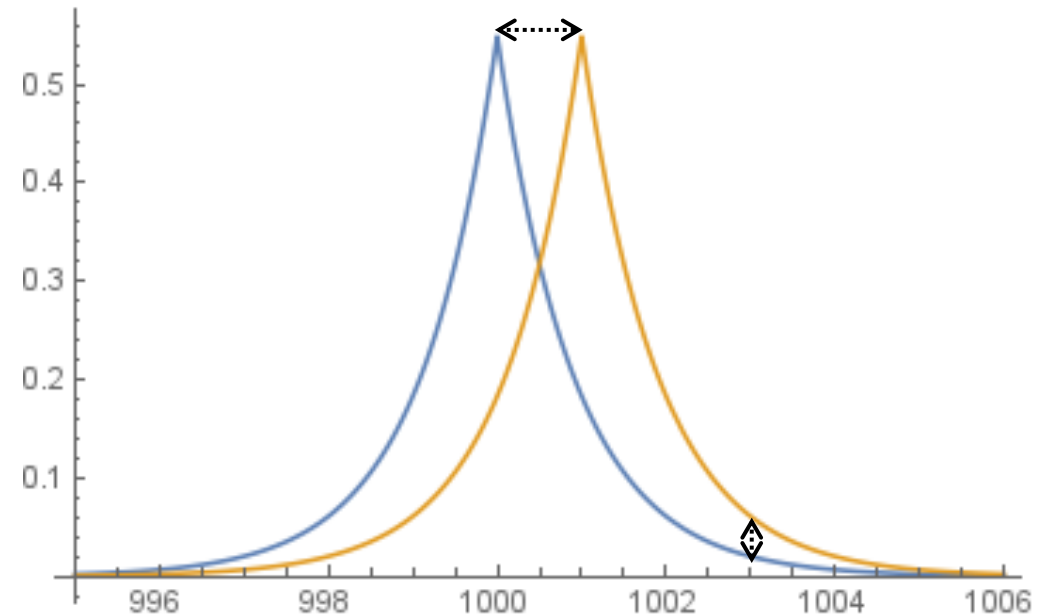
- Assume 10,000 participants
  - 3,000 say they are communist
  - 7,000 say they are not communist
- 50% answers are random so we remove 5,000/2 from each answer pool:
  - 500 are communist
  - 5,500 are not communist

# Another example: a privacy-preserving count query

Query: how many users have green eyes?

Adversarial knowledge: all eye colors besides one person's

Real answer $K=1000$	Real answer $K=1001$
Respond $1000 + \text{Laplace}(1/\epsilon)$	Respond $1001 + \text{Laplace}(1/\epsilon)$
Output 1003	



Probability of  $K=1001$  is  $e^\epsilon$  more likely than  $K=1000$

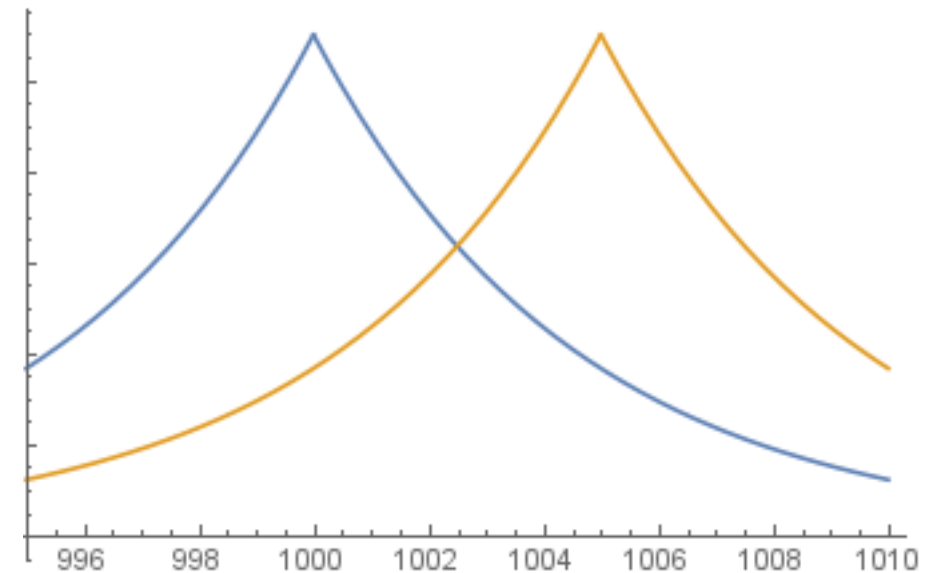
# Another example: a privacy-preserving count query

Query: average rating (between 0 and 5) submitted by users

Average is same than sum / number of users

Adversarial knowledge: all ratings besides one person's sum up to 1000

Real answer $K=1000$ (user votes 0)	Real answer $K=1005$ (user votes 5)
Respond $1000 + \text{Laplace}(5/\epsilon)$	Respond $1005 + \text{Laplace}(5/\epsilon)$



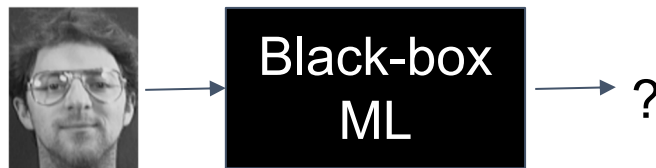
# One final consideration

- What if a user can contribute an outlier value?
  - Compute average of salaries where one individual has a very large salary
- Can pre-process data to remove outliers:
  - Good for privacy + accuracy when computing an average
  - Omission of data points creates new privacy issues
- Can relax definition of differential privacy:

$$Pr[M(d) \in S] \leq e^\epsilon Pr[M(d') \in S] + \delta$$

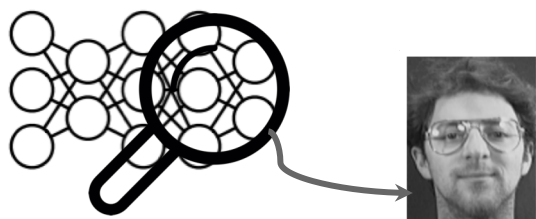
For more details: <https://www.youtube.com/watch?v=oQzaA5KG3pM> (watch first 5 minutes)

# Types of adversaries and our threat model



Model querying (**black-box adversary**)

Shokri et al. (2016) *Membership Inference Attacks*



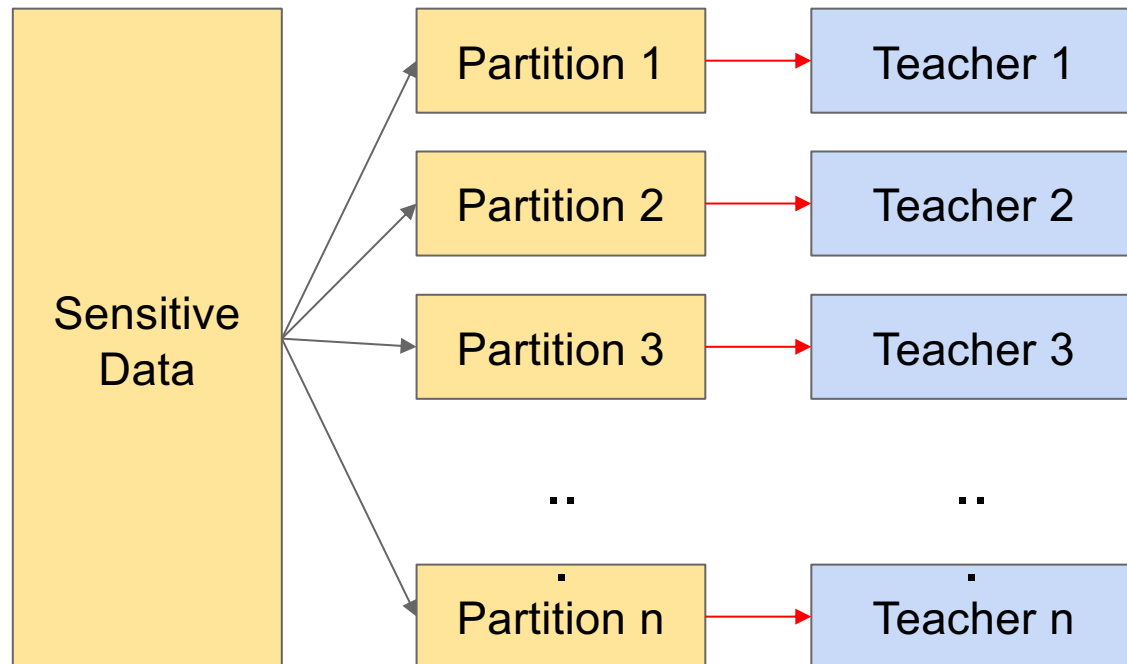
Model inspection (**white-box adversary**)

Zhang et al. (2017) *Understanding DL requires rethinking generalization*

## In our work, the threat model assumes:

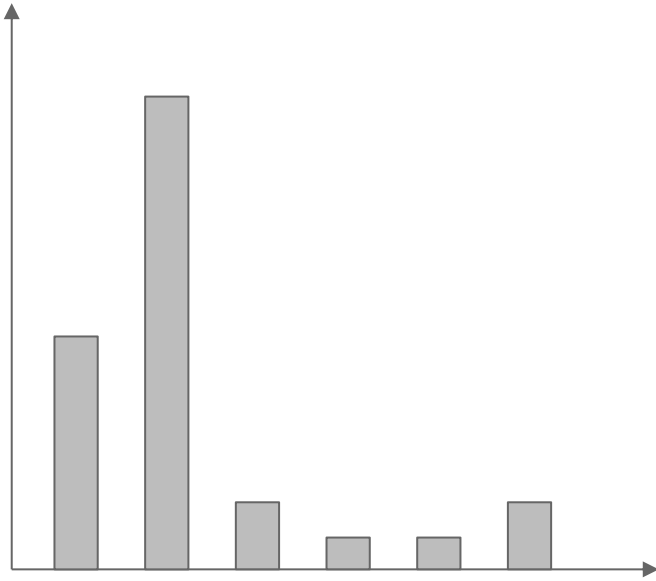
- Adversary can make a potentially unbounded number of queries
- Adversary has access to model internals

# Private Aggregation of Teacher Ensembles (PATE)



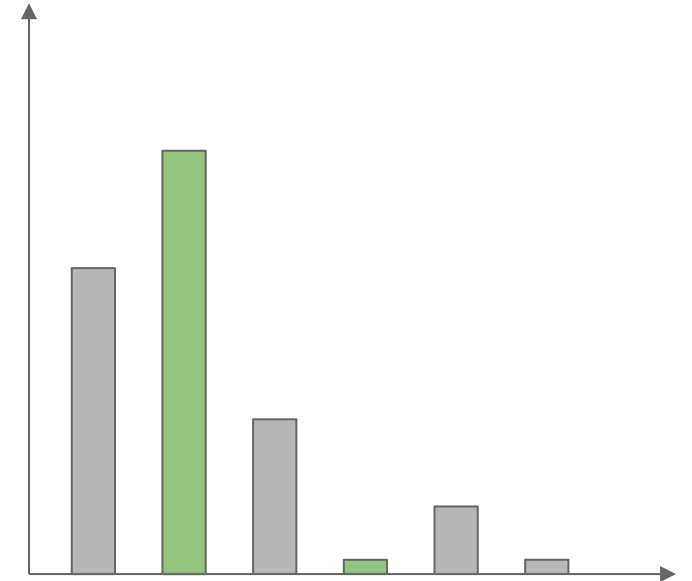


# Aggregation



Count votes

$$n_j(\vec{x}) = |\{i : i \in 1..n, f_i(\vec{x}) = j\}|$$

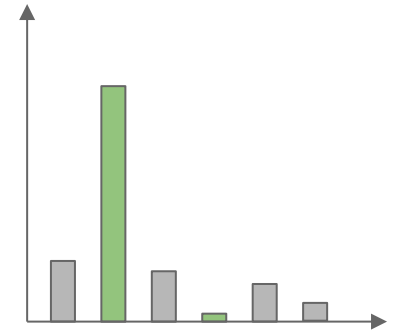


Take maximum

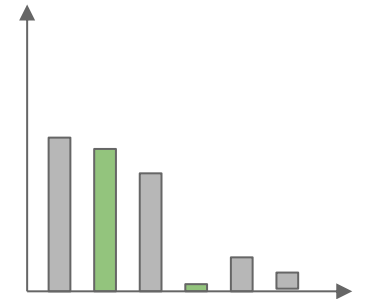
$$f(x) = \arg \max_j \{n_j(\vec{x})\}$$

# Intuitive privacy analysis

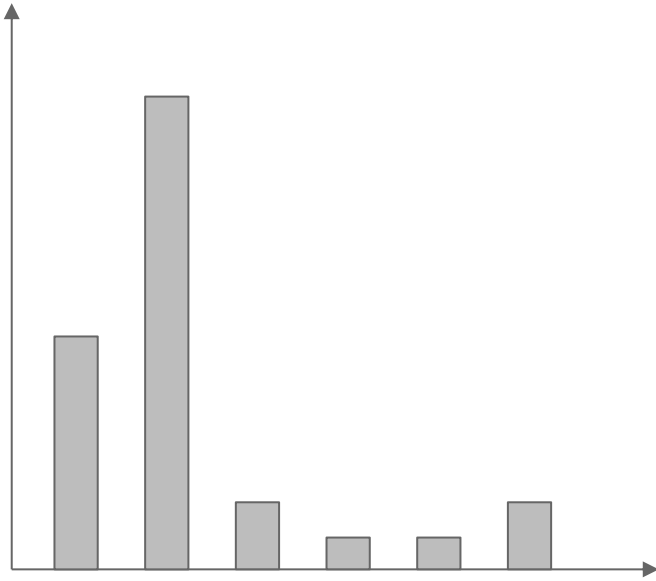
If most teachers agree on the label, it does not depend on specific partitions, so the privacy cost is small.



If two classes have close vote counts, the disagreement may reveal private information.



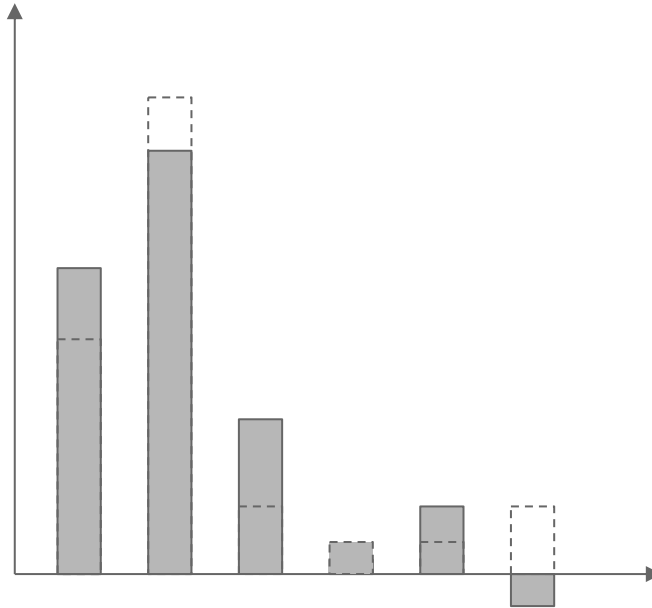
# Noisy aggregation



Count

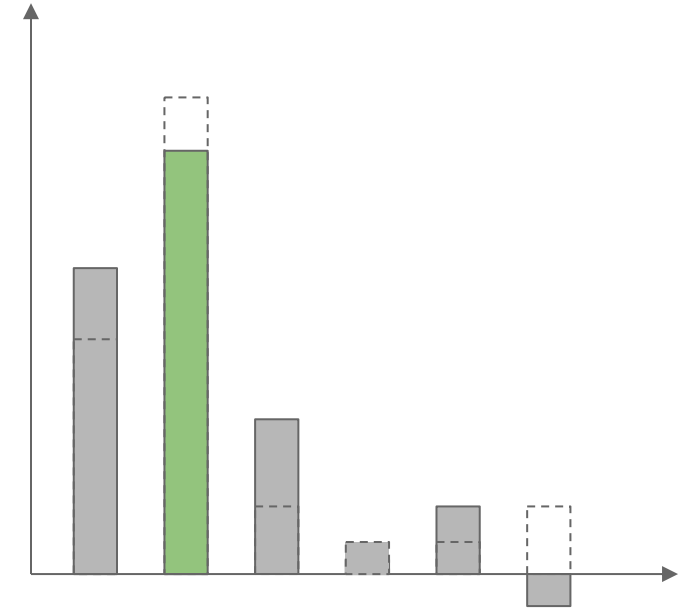
votes

$$n_j(\vec{x}) = |\{i : i \in 1..n, f_i(\vec{x}) = j\}|$$



Add Laplacian

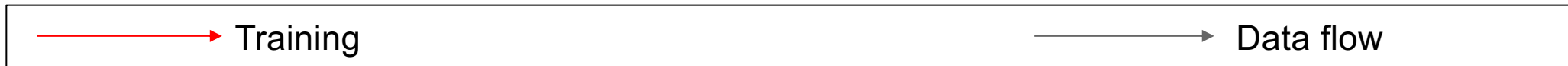
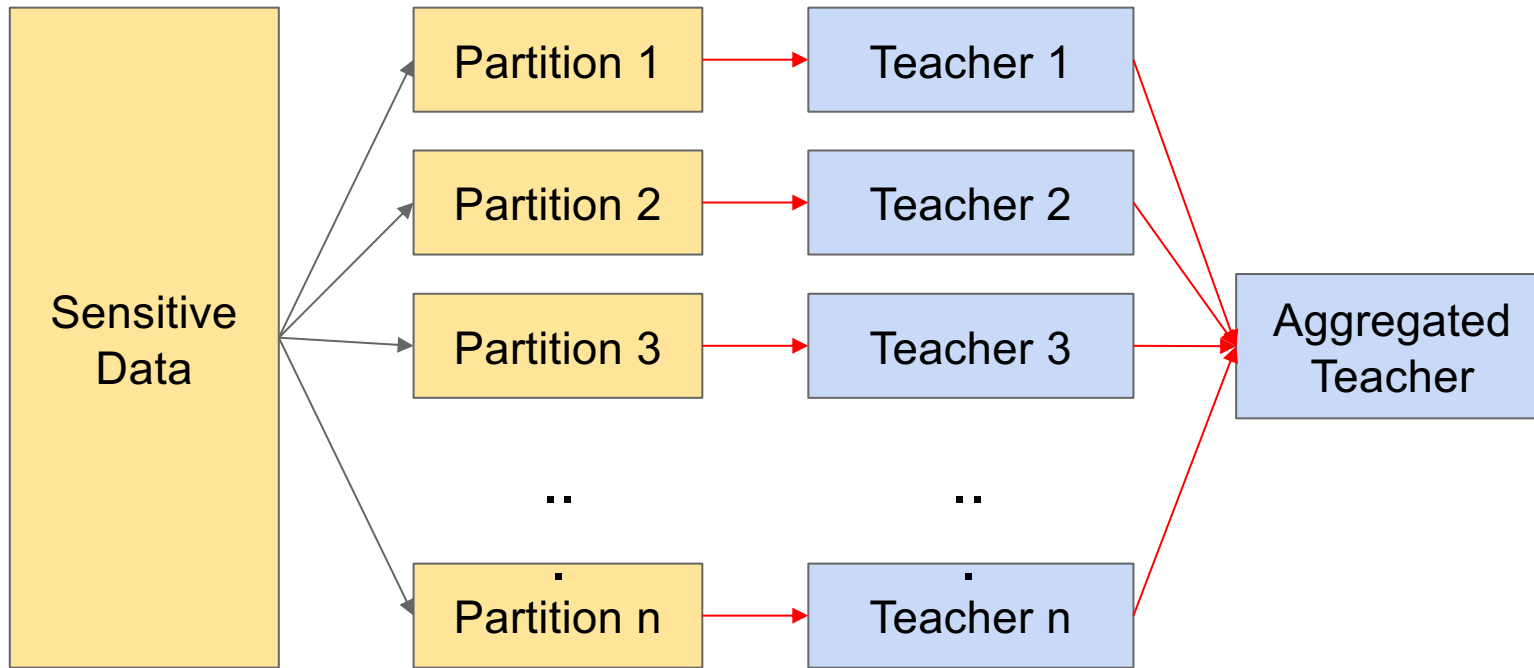
$$Lap\left(\frac{1}{\varepsilon}\right)$$



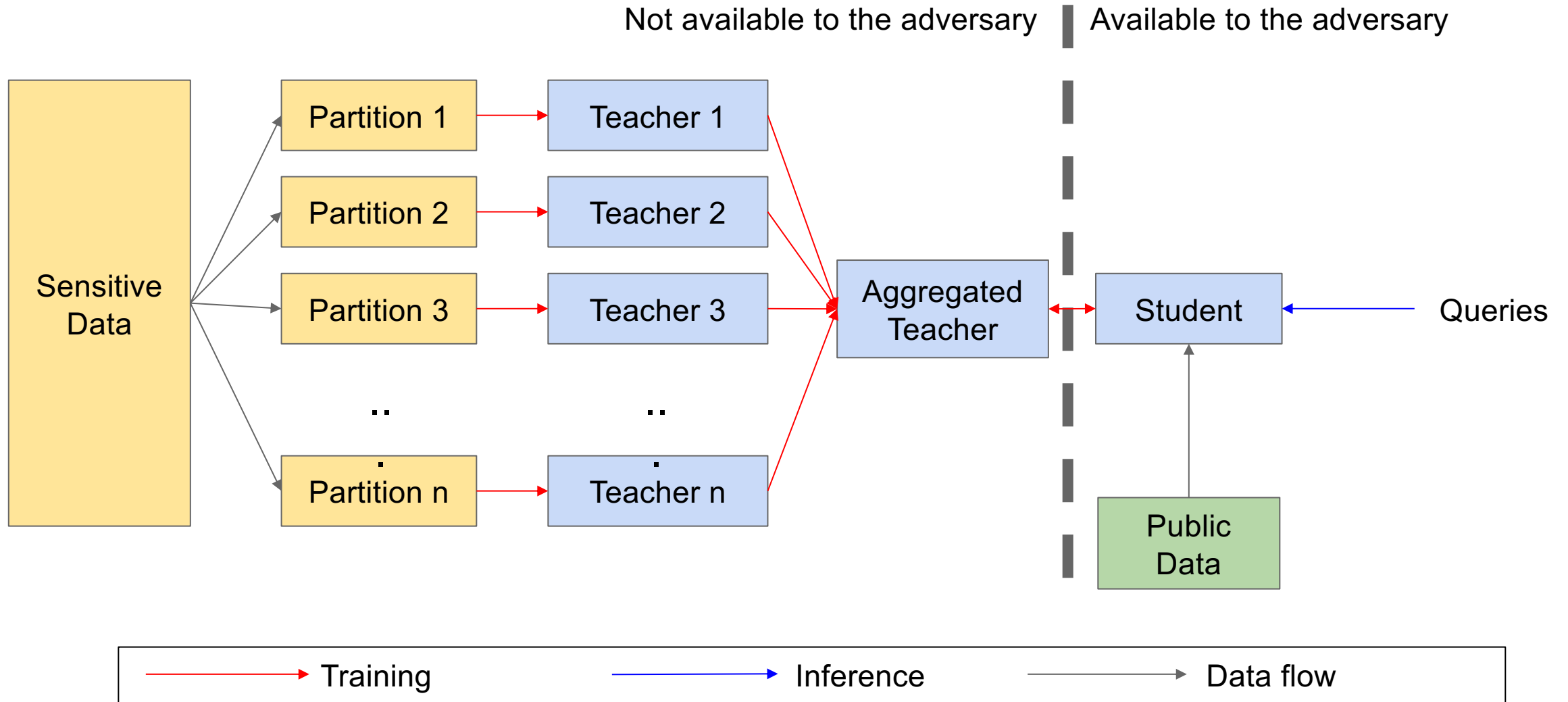
Take maximum

$$f(x) = \arg \max_j \left\{ n_j(\vec{x}) + Lap\left(\frac{1}{\varepsilon}\right) \right\}$$

# Teacher ensemble



# Student training

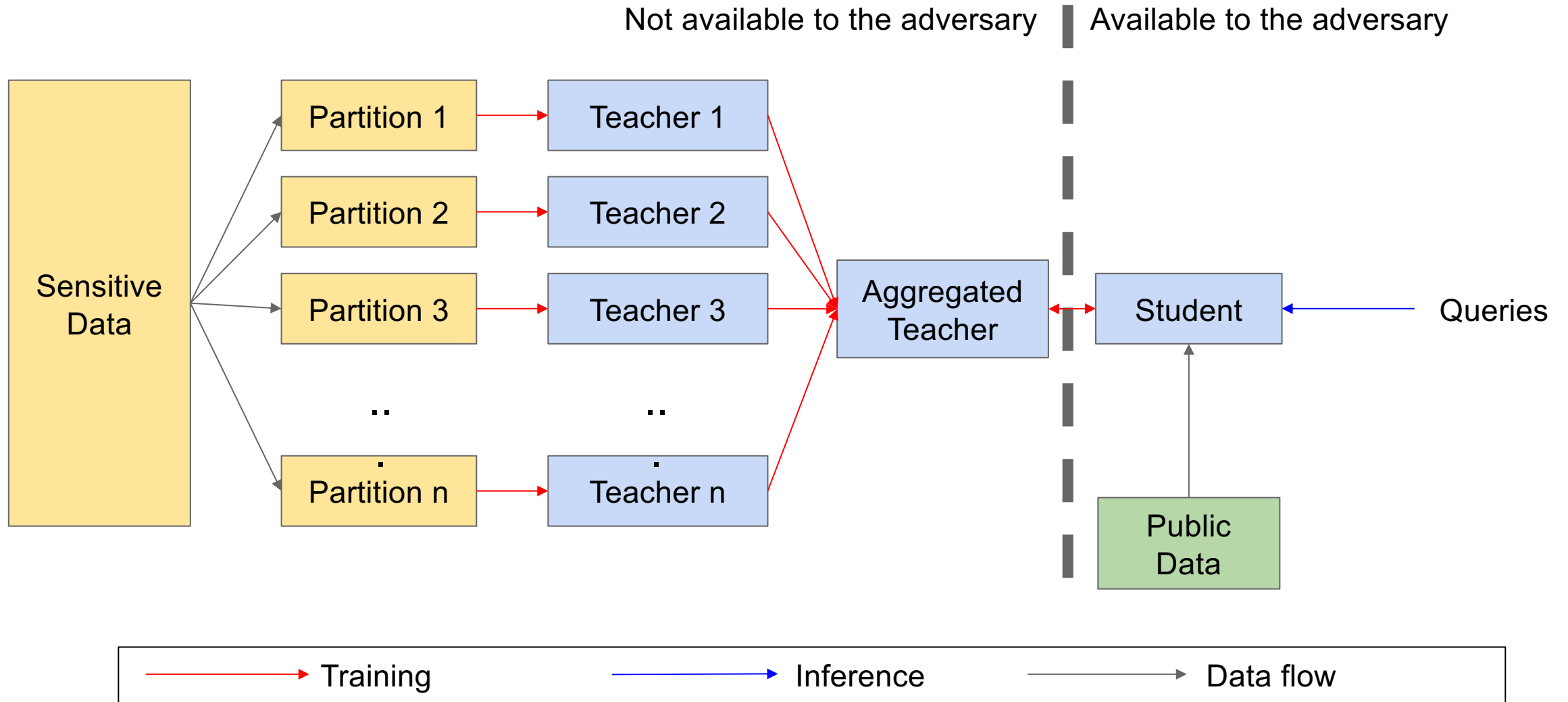


# Why train an additional “student” model?

The aggregated teacher violates our threat model:

- 1 Each prediction increases total privacy loss.**  
Privacy budgets create a tension between the accuracy and number of predictions.
- 2 Inspection of internals may reveal private data.**  
Privacy guarantees should hold in the face of white-box adversaries.

# Student training



# Deployment

Available to the adversary

Student

Queries

Inference



# Differential privacy analysis

## Differential privacy:

A randomized algorithm  $M$  satisfies  $(\epsilon, \delta)$  differential privacy if for all pairs of neighbouring datasets  $(d, d')$ , for all subsets  $S$  of outputs:

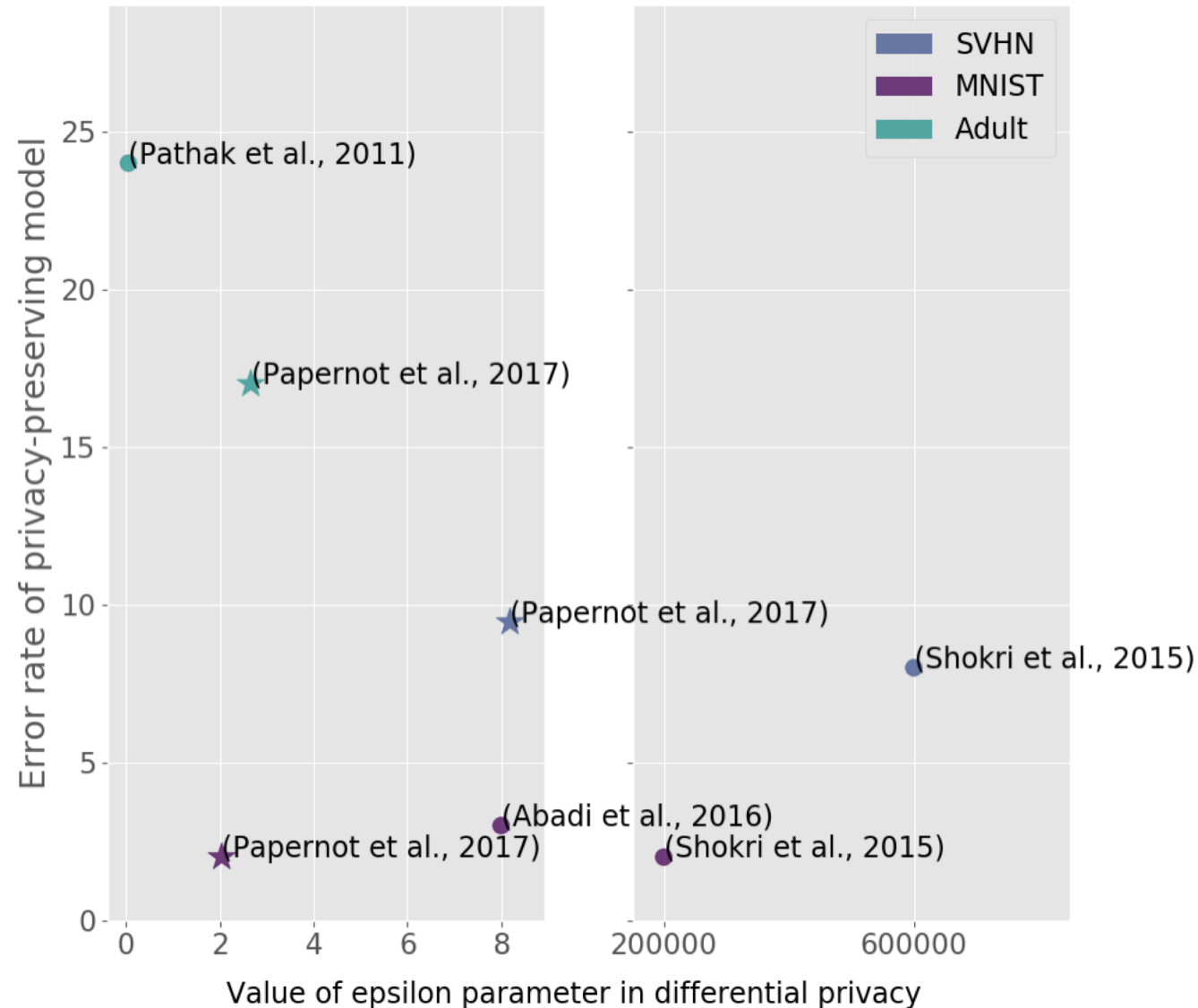
$$\Pr[M(d) \in S] \leq e^\epsilon \Pr[M(d') \in S] + \delta$$

Application of the **Moments Accountant** technique (Abadi et al, 2016)

Strong **quorum**  $\Rightarrow$  Small privacy cost

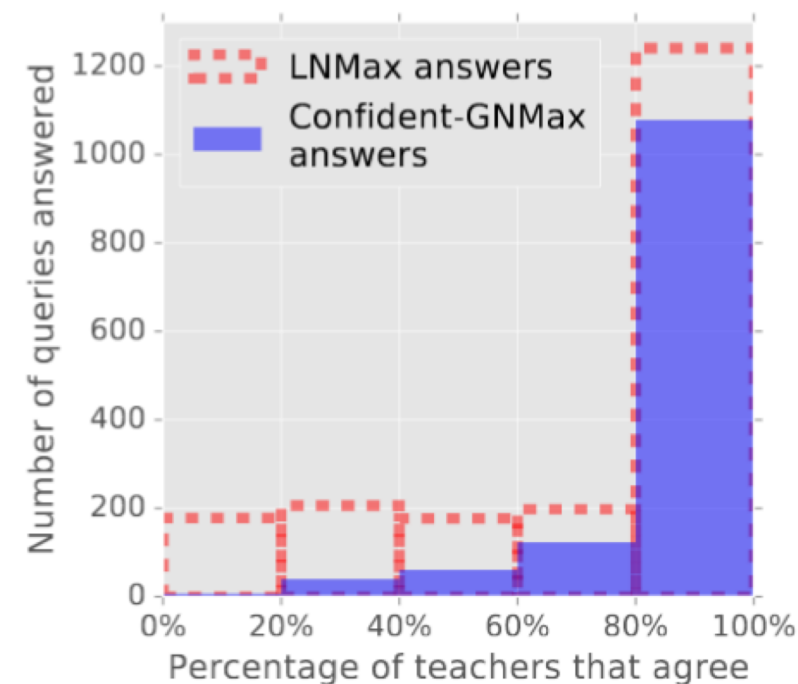
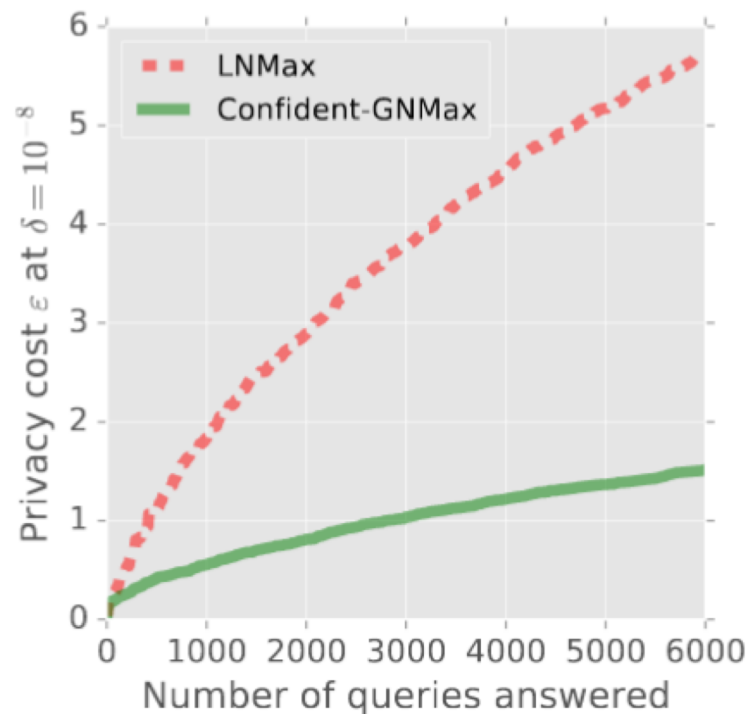
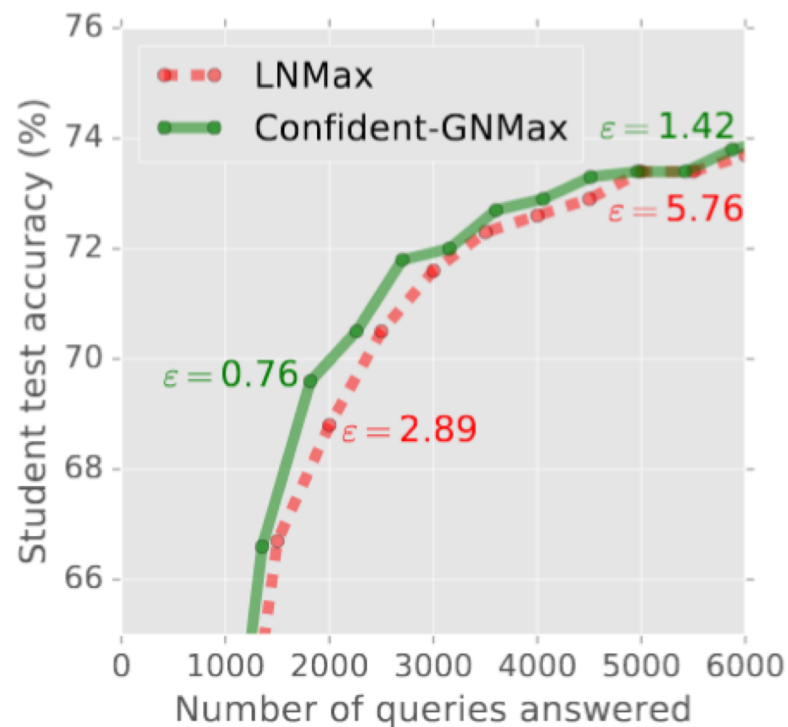
Bound is **data-dependent**: computed using the empirical quorum

# Trade-off between student accuracy and privacy

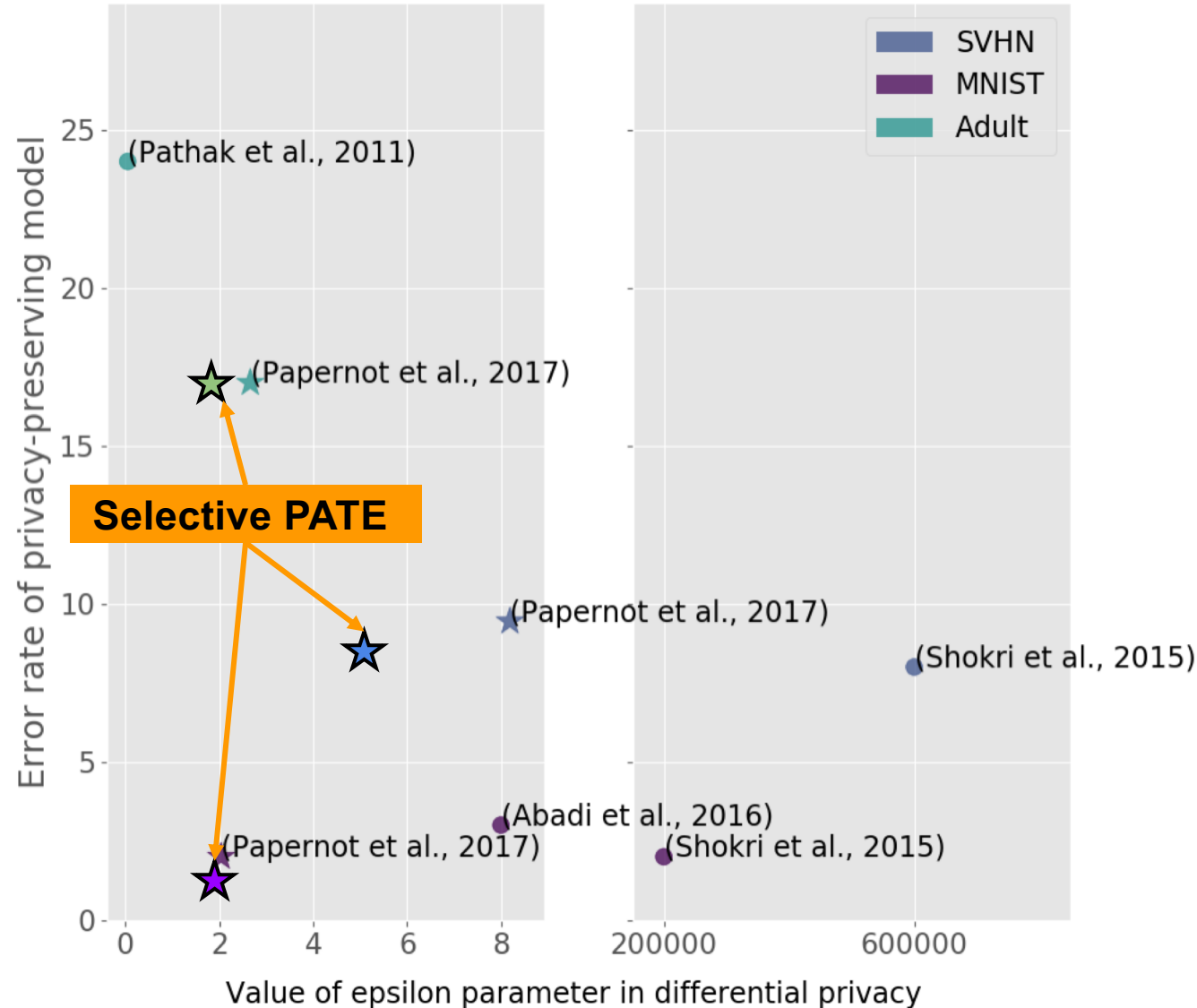


# Synergy between utility and privacy

1. Check privately for consensus
2. Run noisy argmax only when consensus is sufficient



# Trade-off between student accuracy and privacy



# How to train a model with SGD?

```
Initialize parameters  $\theta$ 
```

```
For  $t = 1..T$  do
```

```
    Sample batch  $B$  of training examples
```

```
    Compute average loss  $L$  on batch  $B$ 
```

```
    Compute average gradient of loss  $L$  wrt parameters  $\theta$ 
```

```
    Update parameters  $\theta$  by a multiple of gradient average
```

# How to train a model with differentially private SGD?

```
Initialize parameters  $\theta$ 
```

```
For  $t = 1..T$  do
```

```
    Sample batch  $B$  of training examples
```

```
    Compute per-example loss  $L$  on batch  $B$ 
```

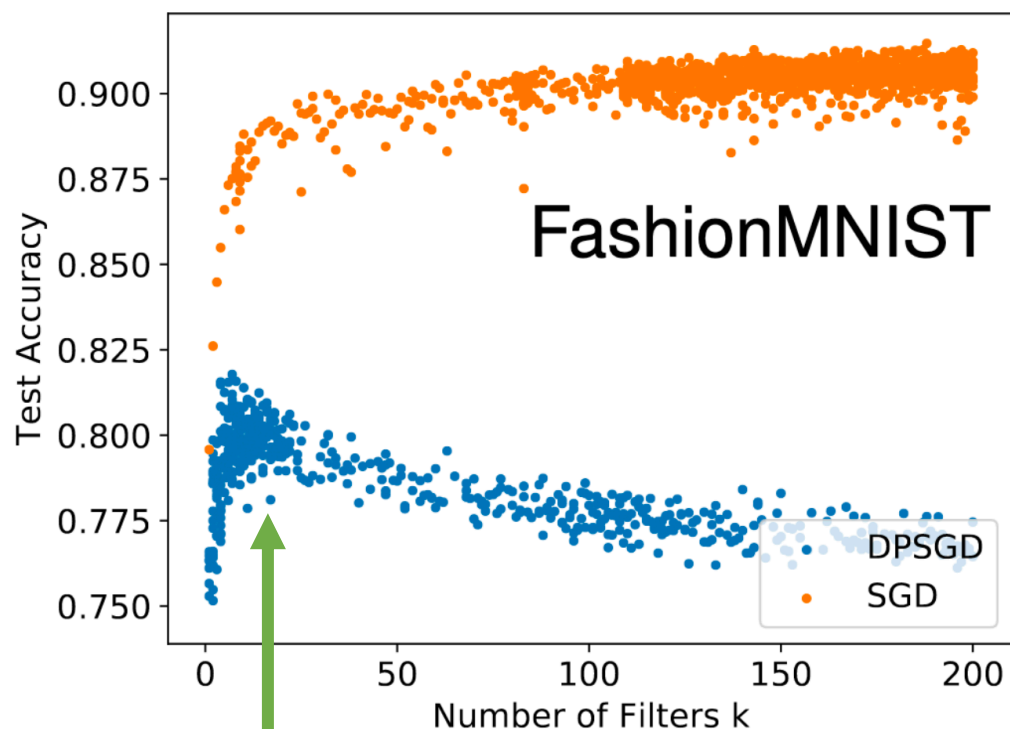
```
    Compute per-example gradients of loss  $L$  wrt parameters  $\theta$ 
```

```
    Ensure L2 norm of gradients  $< C$  by clipping
```

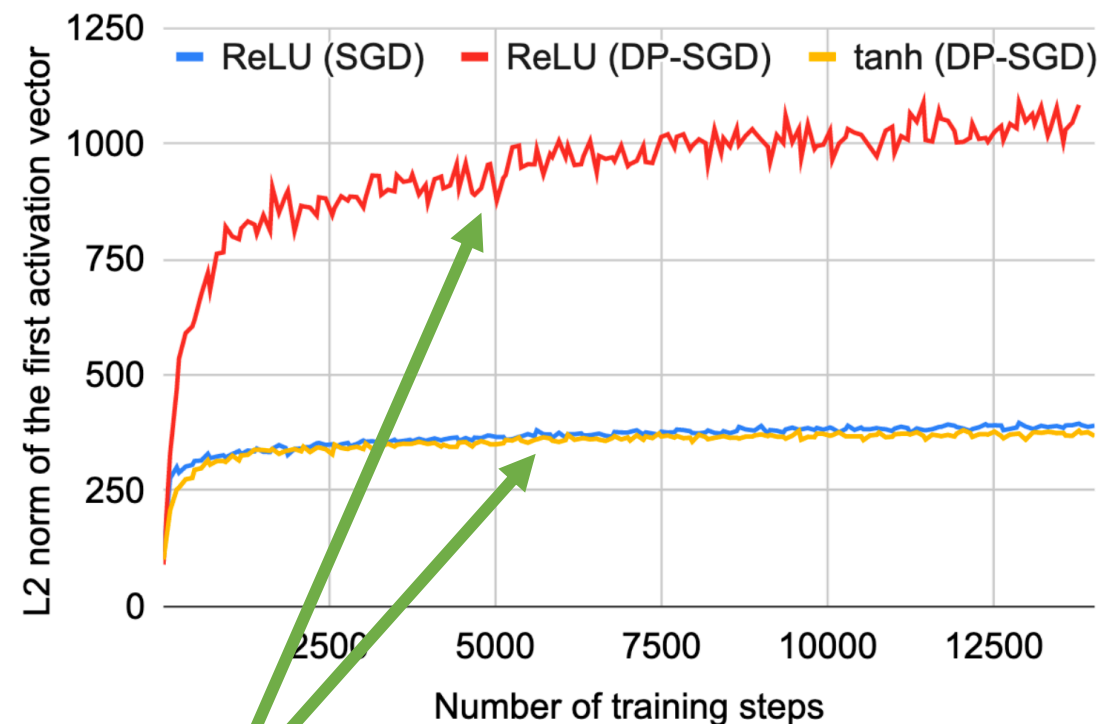
```
    Add Gaussian noise to average gradients (as a function of  $C$ )
```

```
    Update parameters  $\theta$  by a multiple of noisy gradient average
```

# Architectures for DP-SGD learning



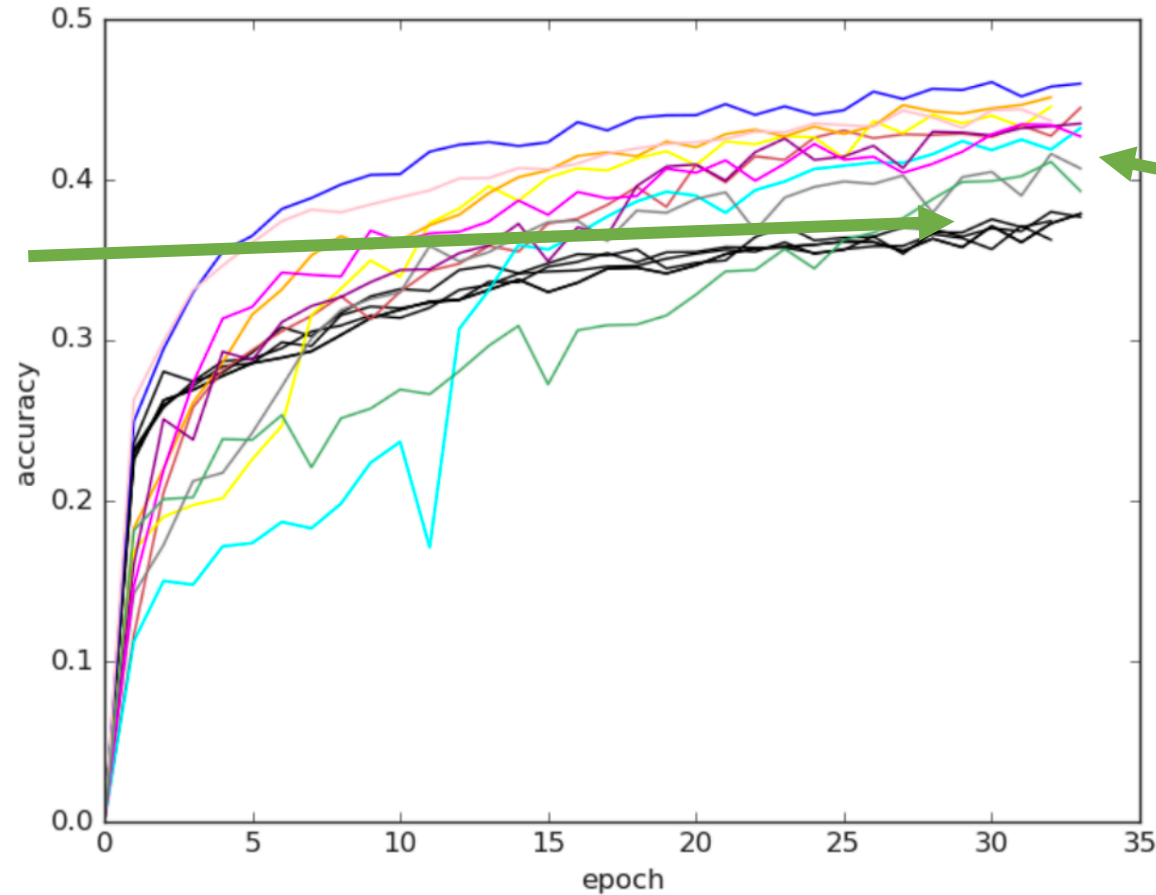
More capacity  
is not always  
helpful



Bounded activations help

# Initializations for DP-SGD learning

Low variance between different initializations when learning without differential privacy



High variance between initializations when learning with DP-SGD and initializing with Raghu et al.



# Hyperparameters for DP-SGD learning

Optimizer	Batch size	Epochs	Non-private		Differentially-private	
			Learning Rate	Test Acc.	Learning Rate	Test Acc.
SGD	256	40	$1.07 \cdot 10^{-1}$	90.3%	$3.32 \cdot 10^{-1}$	86.1%
	1024	7	$3.68 \cdot 10^{-1}$	86.3%	4.46	85.1%
Adam	256	40	$1.06 \cdot 10^{-3}$	90.5%	$1.32 \cdot 10^{-3}$	86.0%
	1024	7	$4.32 \cdot 10^{-3}$	88.7%	$7.08 \cdot 10^{-3}$	85.1%

Training with large batches for few epochs can be competitive in terms of wall-clock time

Best hyperparameter for non-private learning is not best hyperparameter for private learning

Adaptive optimizers are not necessarily helpful

# Architectures, initializations, hyperparameters for DP-SGD learning

Dataset	Technique	Acc.	$\epsilon$	$\delta$	Assumptions
MNIST	SGD w/ tanh (not private)	99.0%	$\infty$	0	-
MNIST	DP-SGD w/ ReLU	96.6%	2.93	$10^{-5}$	-
MNIST	<b>DP-SGD w/ tanh (ours)</b>	<b>98.1%</b>	<b>2.93</b>	$10^{-5}$	-
Fashion	SGD w/ ReLU (not private)	89.4%	$\infty$	0	-
Fashion	DP-SGD w/ ReLU	81.9%	2.7	$10^{-5}$	-
Fashion	<b>DP-SGD w/ tanh (ours)</b>	<b>86.1%</b>	<b>2.7</b>	$10^{-5}$	-
CIFAR10	Transfer + SGD (not private)	75%	$\infty$	0	-
CIFAR10	Transfer + DP-SGD (Abadi et al.)	67%	2	$10^{-5}$	Public Data
CIFAR10	<b>Transfer + DP-SGD (ours)</b>	<b>72%</b>	<b>2.1</b>	$10^{-5}$	Public Data

**Making the Shoe Fit: Architectures, Initializations, and Tuning for Learning with Privacy**

Papernot, Chien, Thakurta, Song, Erlingsson (in submission)

# Useful resources

- <https://desfontain.es/privacy/differential-privacy-awesomeness.html>
- <https://www.cis.upenn.edu/~aaroeth/Papers/privacybook.pdf>
- <https://github.com/tensorflow/privacy>