

# DP-Finder: Finding Differential Privacy Violations by Sampling and Optimization



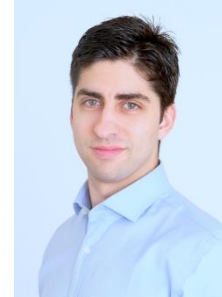
Benjamin Bichsel



Timon Gehr



Dana  
Drachsler-Cohen

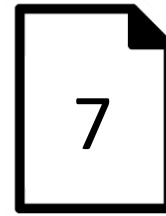


Petar Tsankov

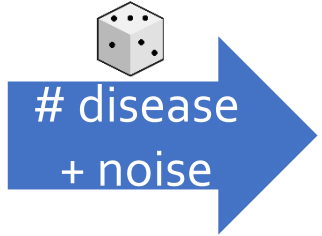


Martin Vechev

# Differential Privacy – Basic Setting

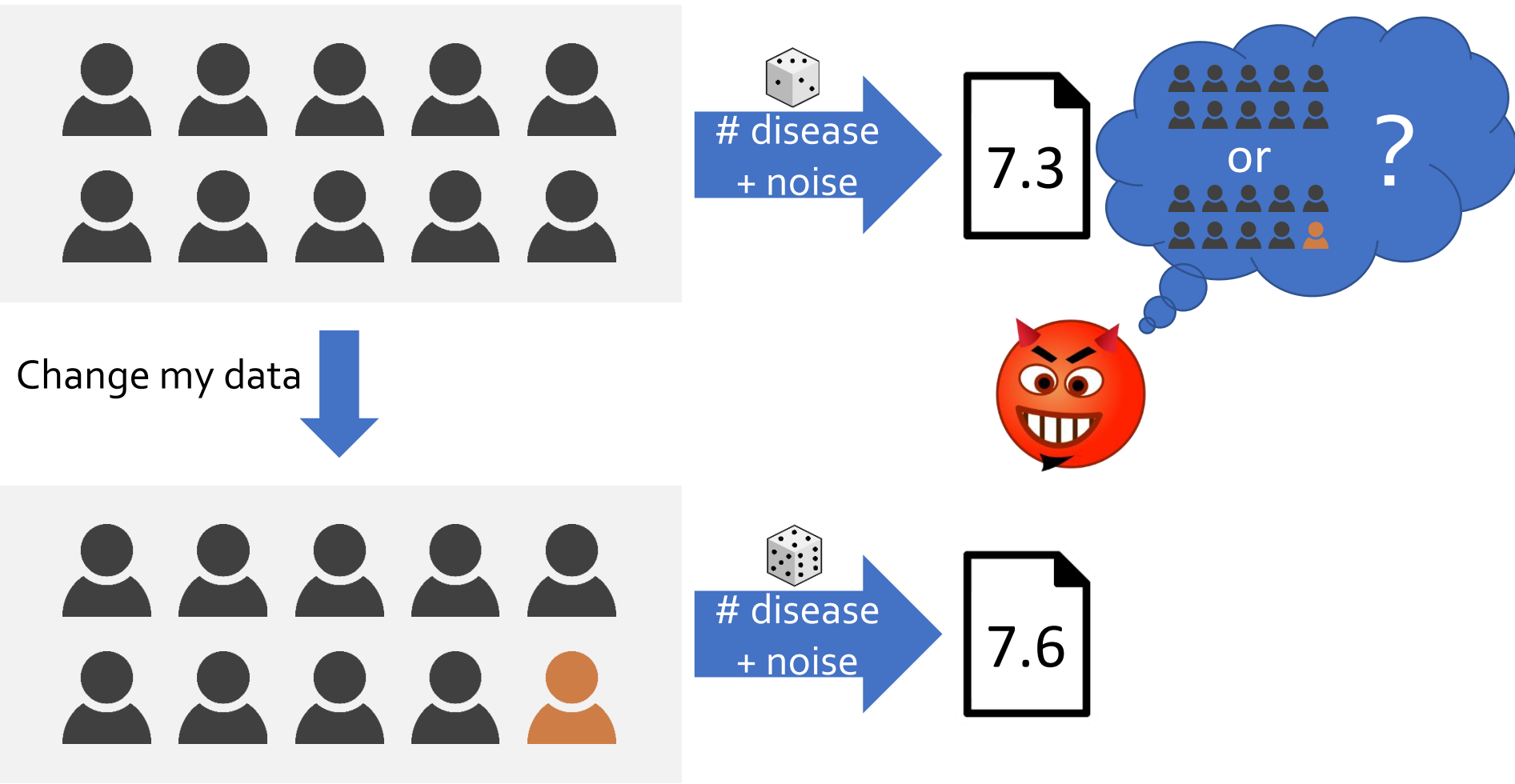


# Differential Privacy – Basic Setting

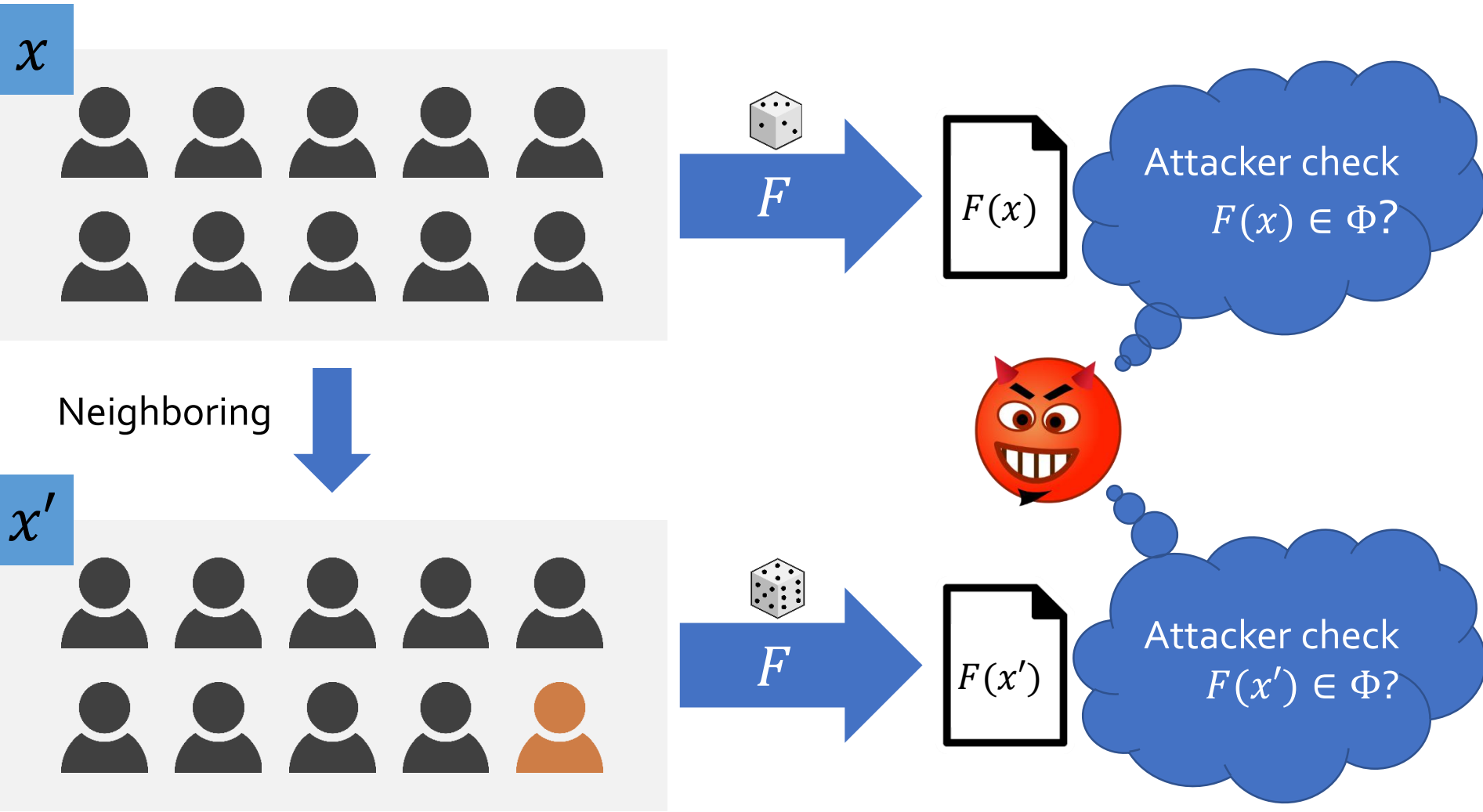


What about my privacy?

# Differential Privacy - Intuition



# Differential Privacy – More Abstractly



# Differential Privacy - Definition

$x$



Neighbouring



$x'$



$\epsilon$ -DP:

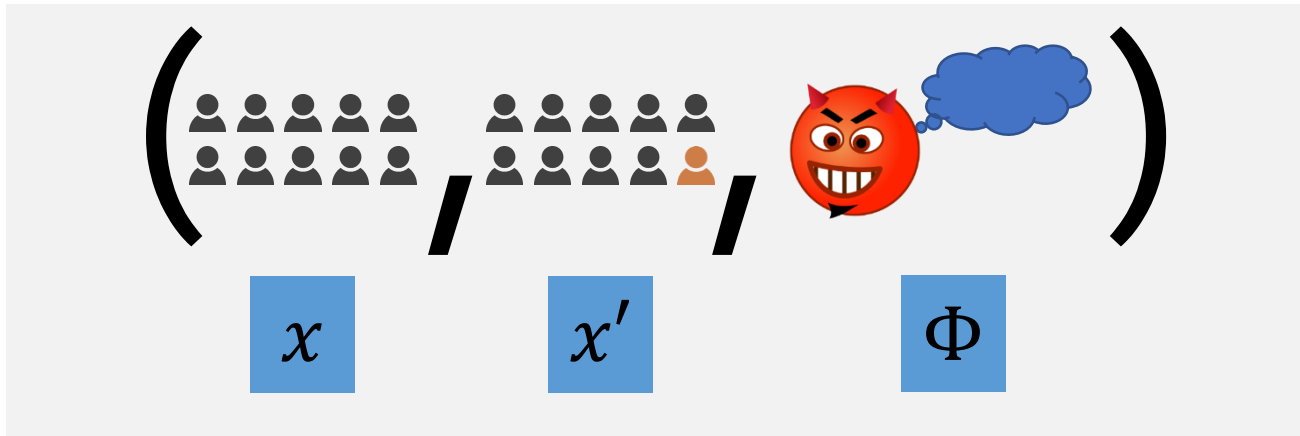
$$\frac{\Pr[F(x) \in \Phi]}{\Pr[F(x') \in \Phi]} \leq \exp(\epsilon) \approx 1 + \epsilon$$



Challenges induced by DP:

- Proving/checking  $\epsilon$ -DP is hard (buggy algorithms)
- Proof strategies not complete
- Proofs only provide upper bounds

# $\epsilon$ -DP Counterexamples



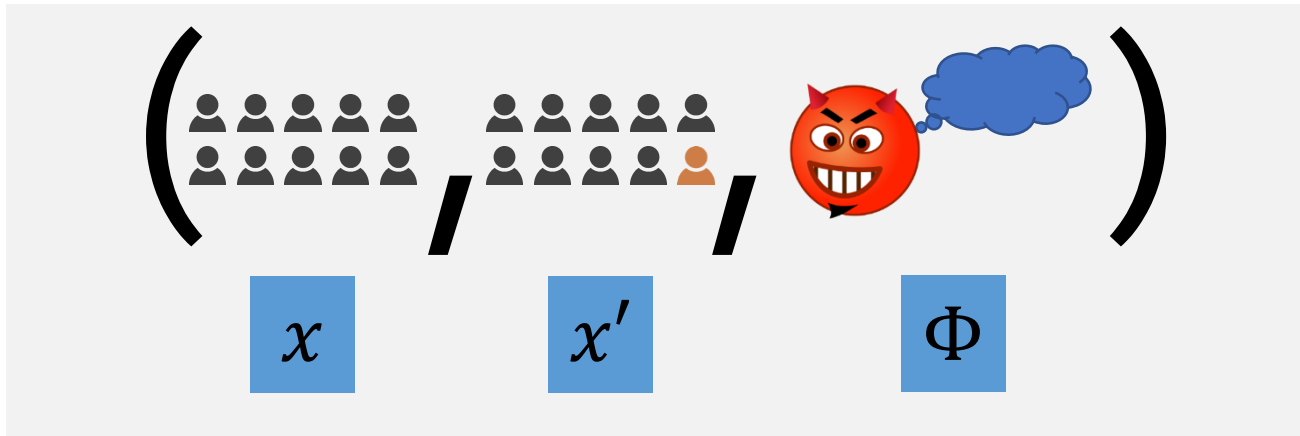
that violate  $\epsilon$ -DP:

$$\frac{\Pr[F(x) \in \Phi]}{\Pr[F(x') \in \Phi]} > \exp(\epsilon)$$

$\Leftrightarrow$

$$\log \frac{\Pr[F(x) \in \Phi]}{\Pr[F(x') \in \Phi]} > \epsilon$$

# $\epsilon$ -DP Counterexamples



that violate  $\epsilon$ -DP:

$$\frac{\Pr[F(x) \in \Phi]}{\Pr[F(x') \in \Phi]} > \exp(\epsilon)$$

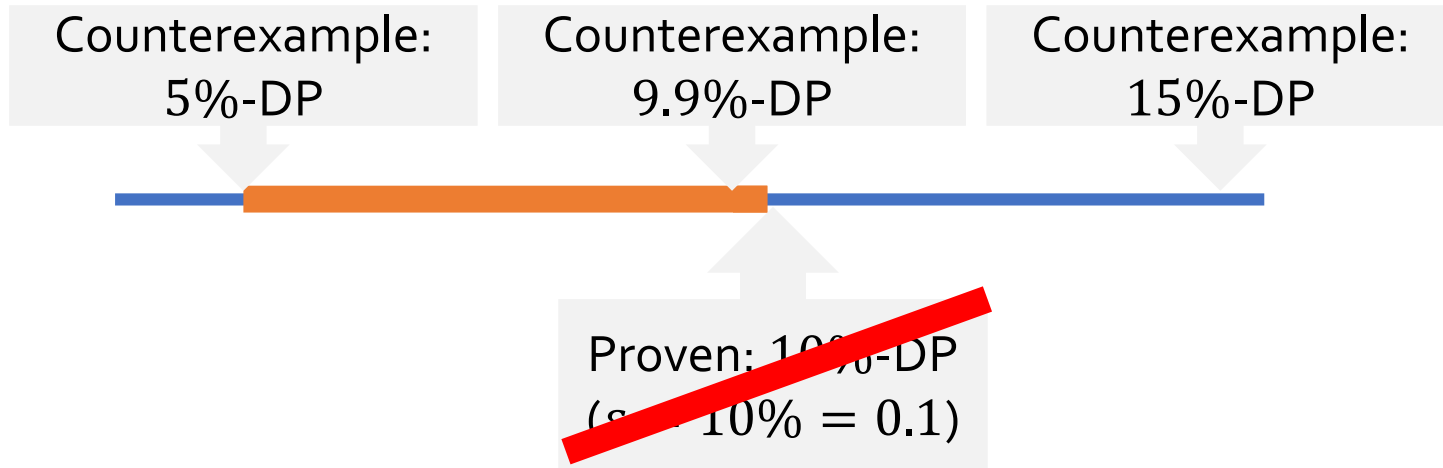
$\Leftrightarrow$

Maximize  
 $\epsilon(x, x', \Phi)$

$$\log \frac{\Pr[F(x) \in \Phi]}{\Pr[F(x') \in \Phi]} > \epsilon$$



# Bounds on "true" $\varepsilon$



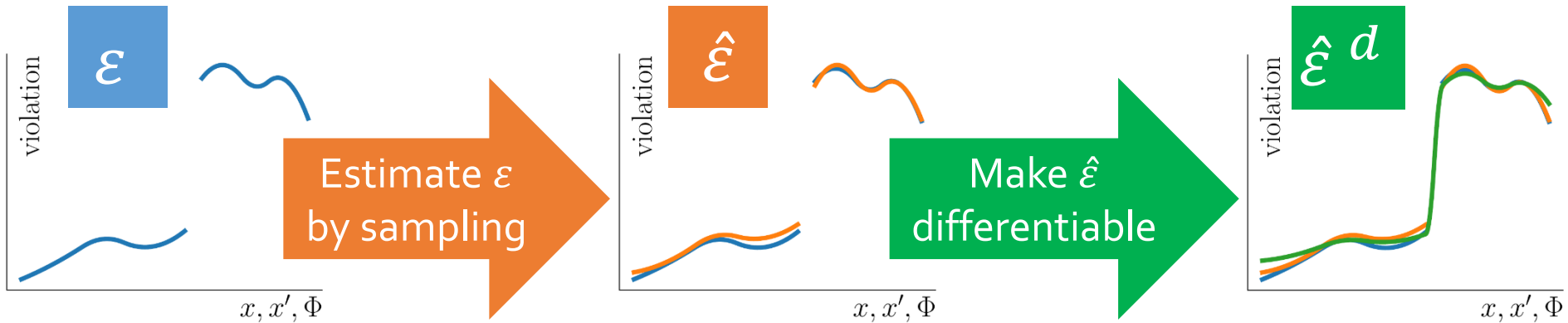
**Evaluation:** We get precise and large  $\varepsilon$ , close to known upper bounds

# $\varepsilon$ -DP Counterexamples

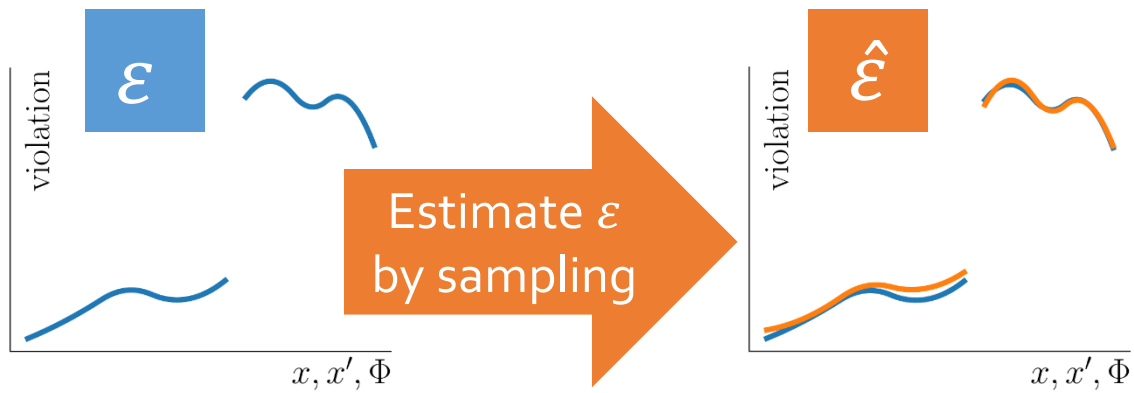
**Goal:** Maximize  $\varepsilon(x, x', \Phi)$

**Challenge 1:** Expensive to compute  $\varepsilon$  precisely

**Challenge 2:** Search space is sparse: Few  $x, x', \Phi$  lead to large  $\varepsilon(x, x', \Phi)$



# Step 1: Estimate $\varepsilon$



# Estimating $\varepsilon$

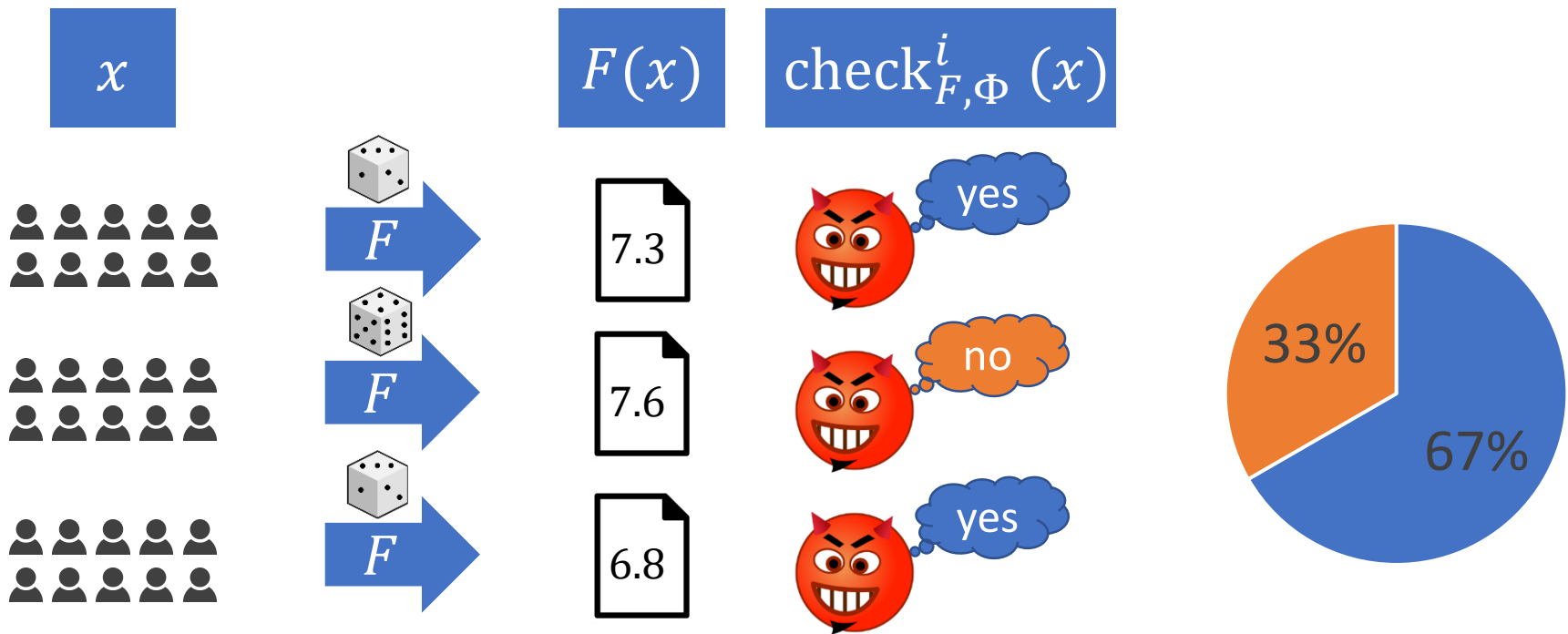
$$\varepsilon(x, x', \Phi) := \log \frac{\Pr[F(x) \in \Phi]}{\Pr[F(x') \in \Phi]}$$

# Estimating $\varepsilon$

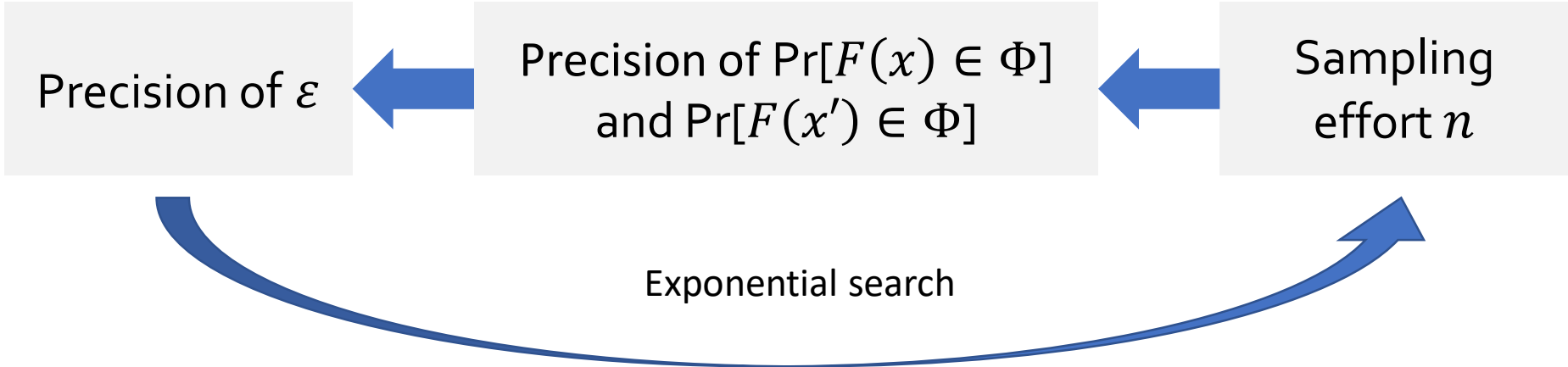
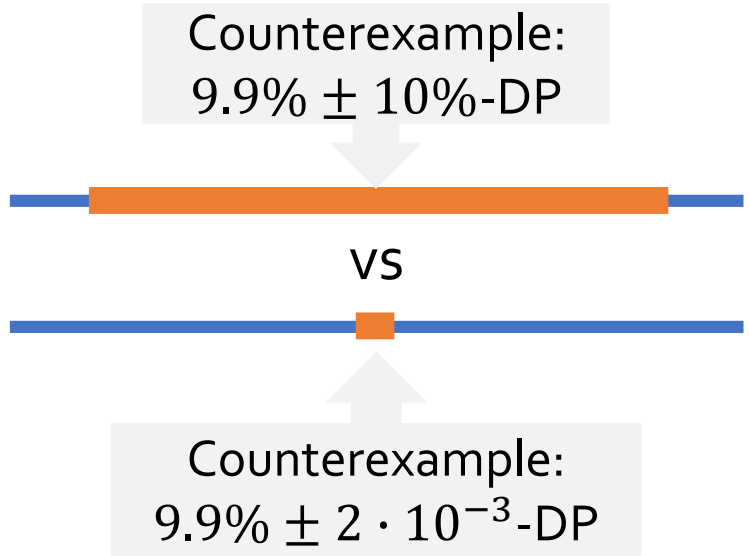
Estimate

$$\varepsilon(x, x', \Phi) := \log \frac{\Pr[F(x) \in \Phi]}{\Pr[F(x') \in \Phi]}$$

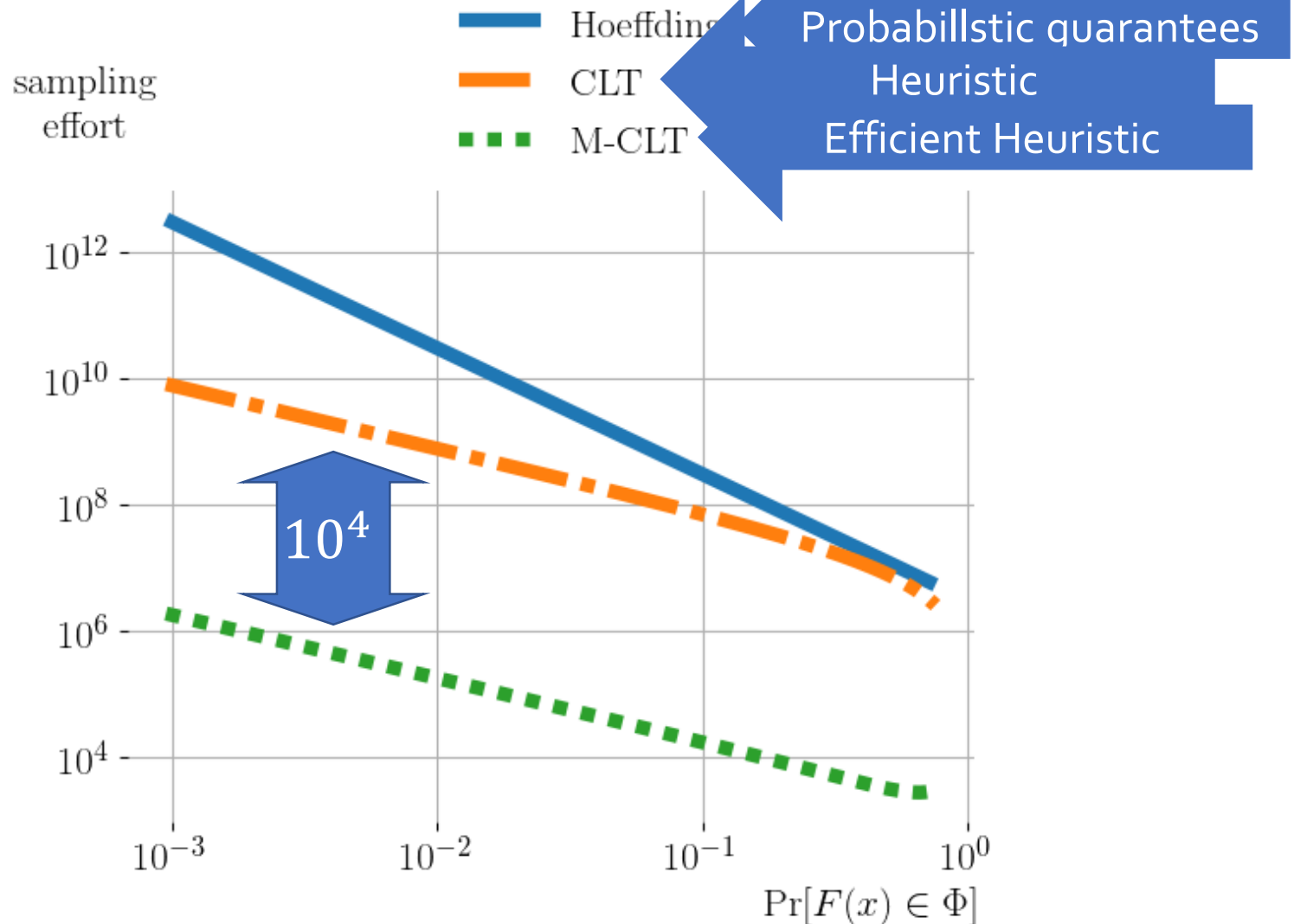
$$\widehat{\Pr}[F(x) \in \Phi] = \frac{1}{n} \sum_{i=1}^n \text{check}_{F, \Phi}^i(x)$$



# How precise is our estimate?

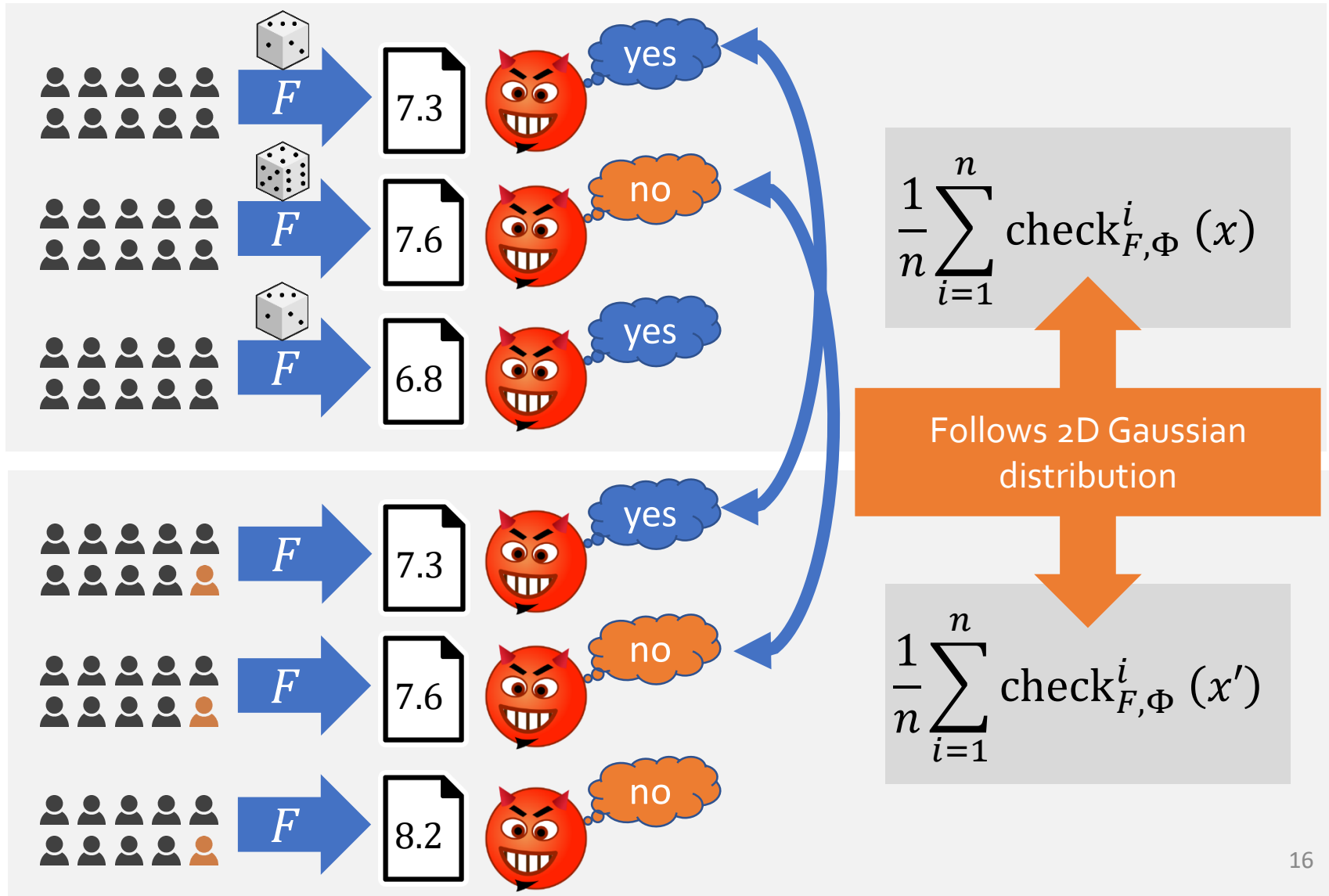


# Estimating precisely is expensive



Estimating  $\varepsilon$  up to an error of  $2 \cdot 10^{-3}$  with confidence of 90%

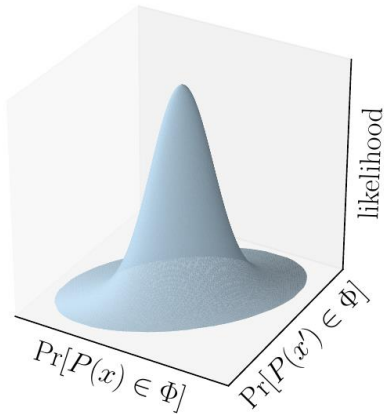
# Applying the M-CLT (Correlation)



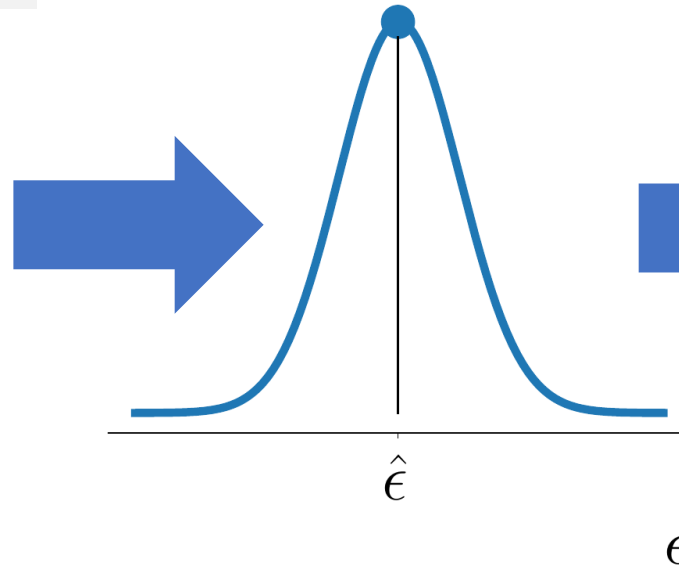


# Obtaining a Confidence Interval for $\varepsilon$

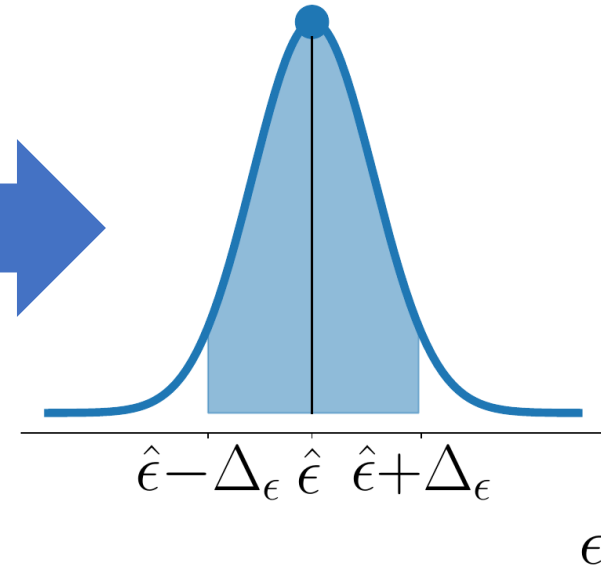
Joint likelihood of  
 $\left( \Pr[F(x) \in \Phi] \right)$   
 $\left( \Pr[F(x') \in \Phi] \right)$



Likelihood of  
 $\varepsilon(x, x', \Phi)$



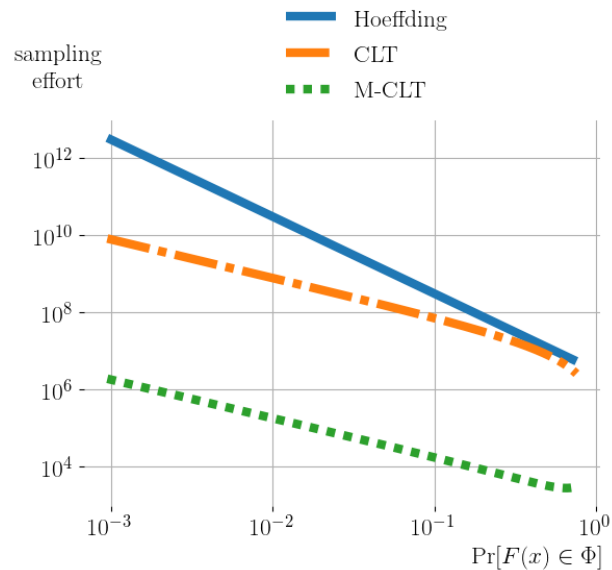
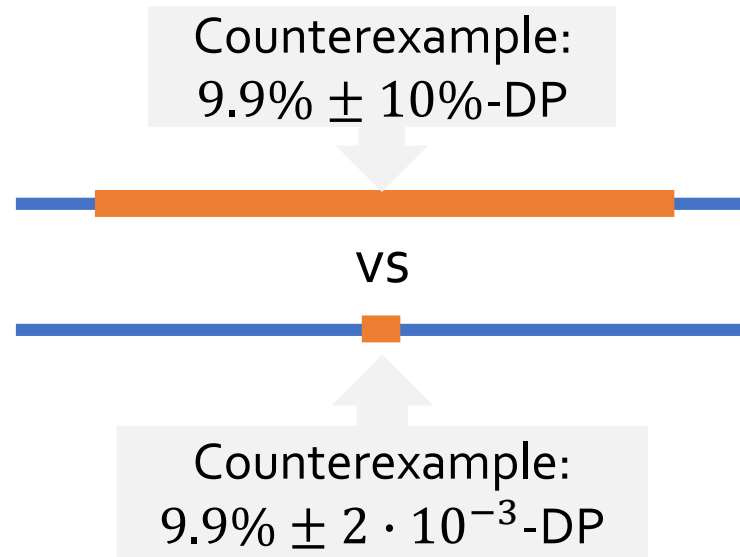
Confidence Interval  
for  $\varepsilon(x, x', \Phi)$



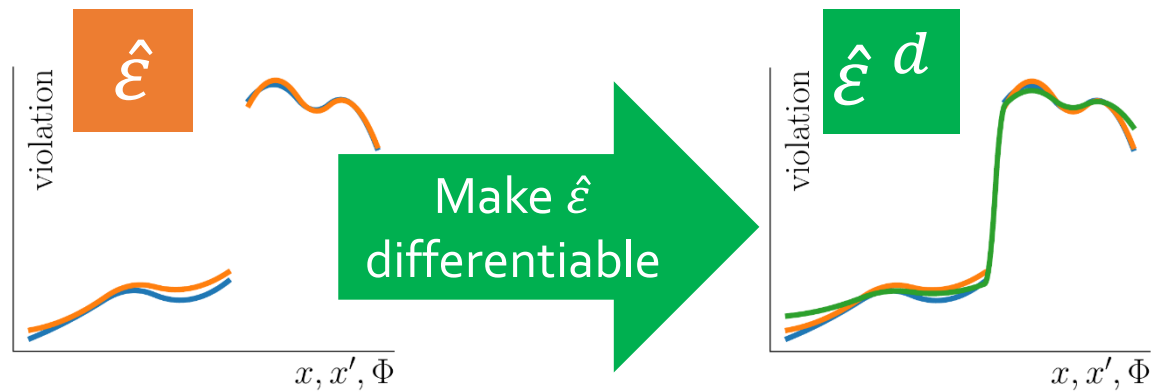
Distribution of  $\frac{\text{Gauss}}{\text{Gauss}}$  (correlated):

D. V. Hinkley. 1969. On the Ratio of Two Correlated Normal Random Variables. *Biometrika* 56, 3 (1969), 635–639. <http://www.jstor.org/stable/2334671>

# How precise is our estimate?



# Step 2: Finding Counterexamples



# How can we optimize our estimate?

maximize

$$\hat{\epsilon}(x, x', \Phi) = \log \frac{\frac{1}{n} \sum_{i=1}^n \text{check}_{F, \Phi}^i}{\frac{1}{n} \sum_{i=1}^n \text{check}_{F, \Phi}^i(x')}$$

Not differentiable

## Goals

- Make differentiable
- Preserve semantics

$$\neg B \sim 1 - B$$

$$B_1 \wedge B_2 \sim B_1 \cdot B_2$$

$$\text{if } (B) : \{x = E_1\} \text{ else } : \{x = E_2\} \sim x = B \cdot E_1 + (1 - B) \cdot E_2$$

# How can we optimize our estimate?

maximize

$$\hat{\epsilon}(x, x', \Phi) = \log \frac{\frac{1}{n} \sum_{i=1}^n \text{check}_{F, \Phi}^i}{\frac{1}{n} \sum_{i=1}^n \text{check}_{F, \Phi}^i(x')}$$

Not differentiable

- Maximize using SLSQP (supports hard constraints for neighborhood)
- Random starting point (+ restart)
- What about division by zero?
- What about very small denominators?

# Main differences to Ding et al.

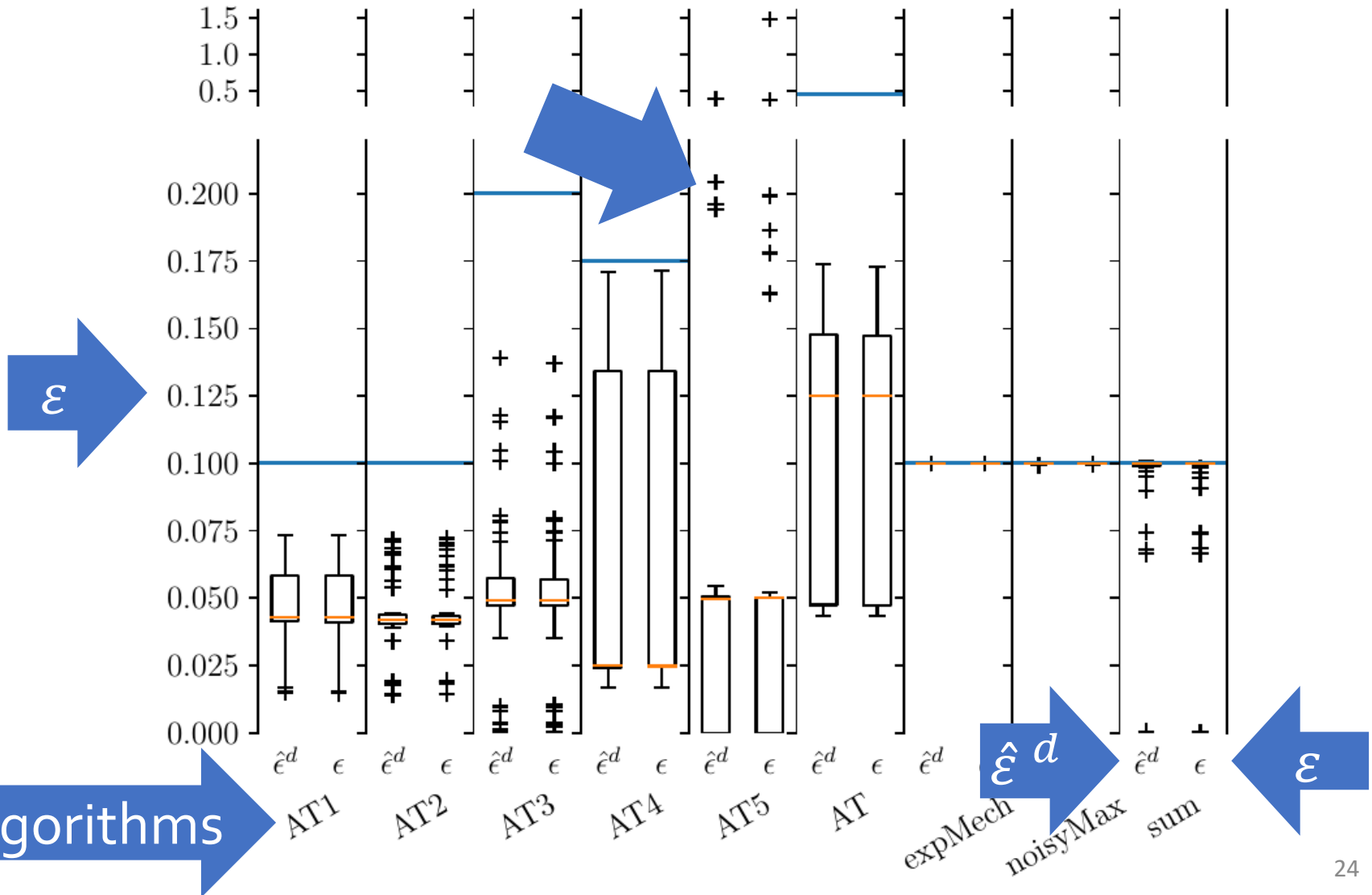
Dimension	Ding et al.	This work
Problem statement	$\varepsilon(x, x', \Phi) > \varepsilon_0?$	Maximize $\varepsilon(x, x', \Phi)$
Approach	Statistical tests	Estimate + confidence interval
Search	By patterns	Gradient descent (incremental)

# Evaluation

- How **precise** is the differentiable estimate?
- How **efficient** is DP-Finder in finding violations compared to random search?

Exact solver (PSI)  
for ground truth

# Precision of Differentiable Estimate

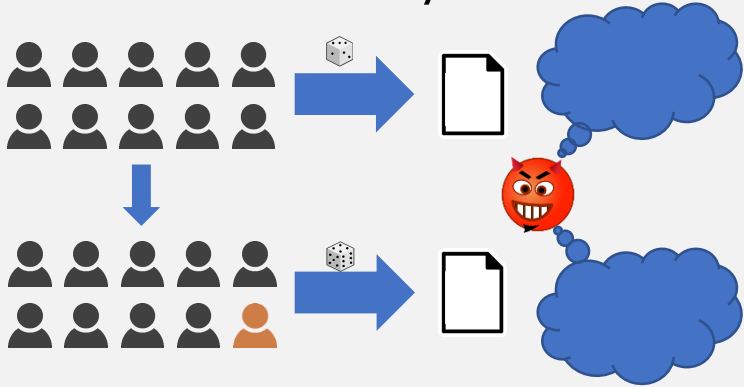






# Conclusion

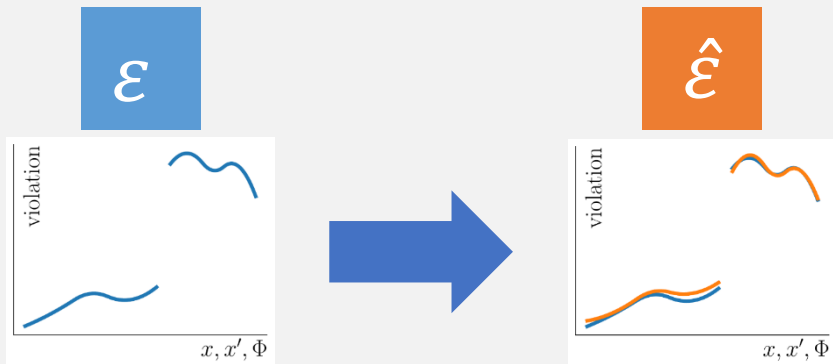
## Differential Privacy



## $\epsilon$ -DP Counterexamples



## Estimate $\epsilon$



## Finding Counterexamples

