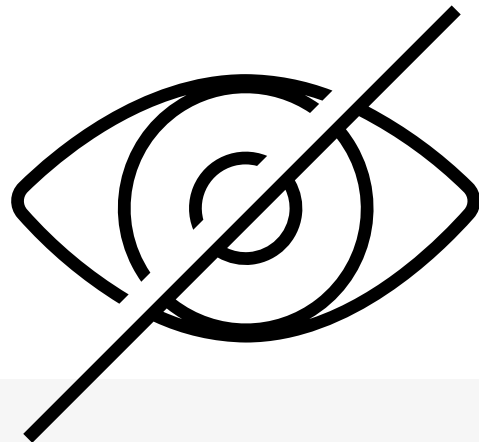
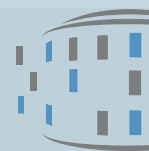


# Locally Private Histograms in All Privacy Regimes

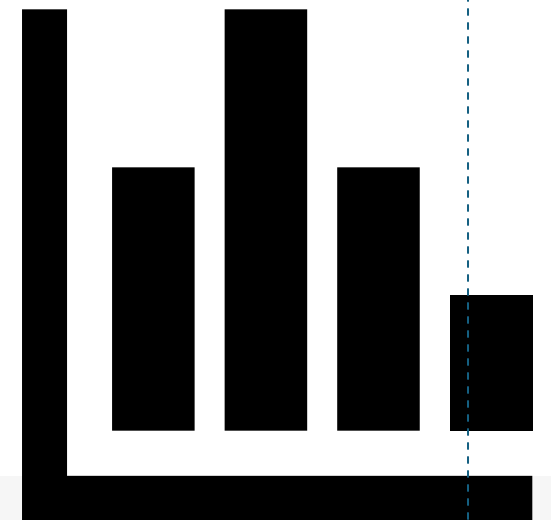


Clément Canonne  
(University of Sydney)

Joint work with Abigail Gentle

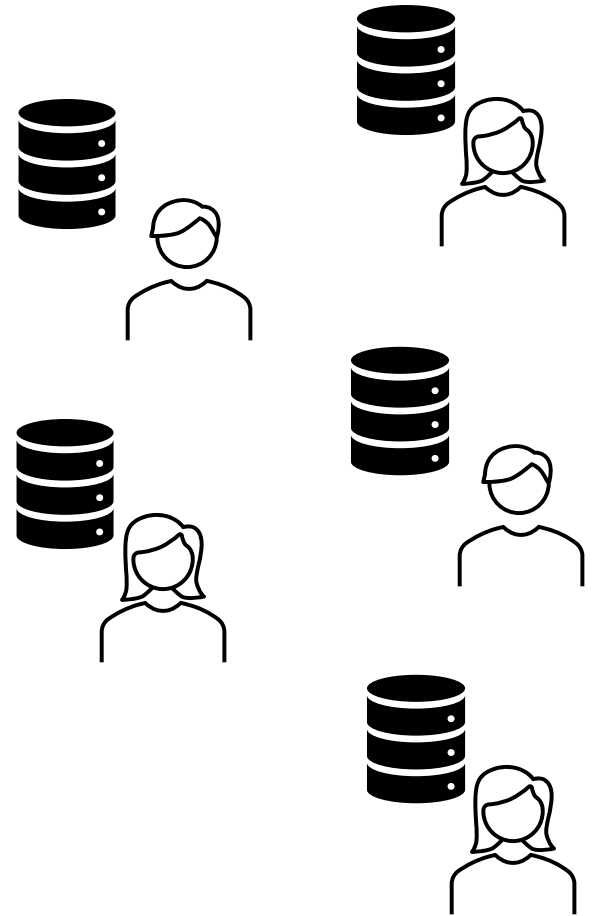
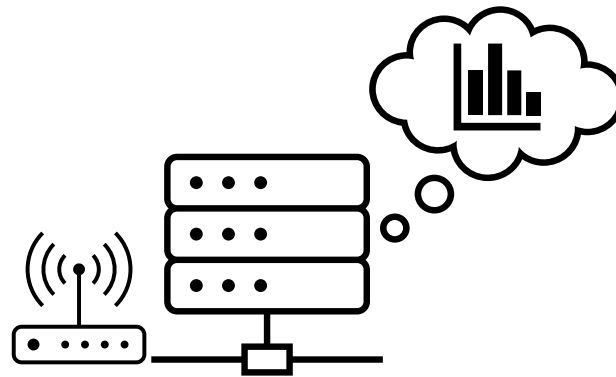


**SIMONS**  
INSTITUTE  
for the Theory of Computing



# The (distributed) setting

- Users have **data** (observations)
- Center wants to **learn from this data**



# The (distributed) setting: example

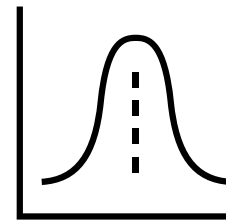
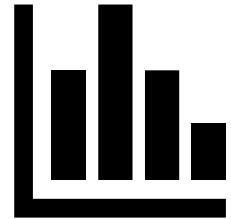
- Users have **data**  $X_1, \dots, X_n$



- Personal data
- Observations i.i.d. from some distribution  $p$

- Can we learn the counts/frequencies/histogram?

- To learn about the users' preferences
- To learn about the underlying data distribution



# The (distributed) setting: example

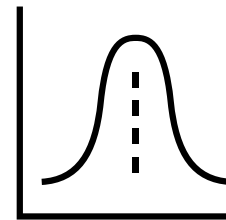
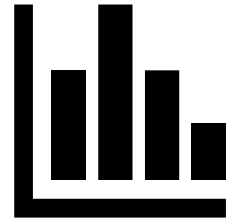
- Users have **data**  $X_1, \dots, X_n$



- Personal data
- Observations i.i.d. from some distribution  $p$

- Can we learn the counts/frequencies/histogram?

- To learn about the users' preferences
- To learn about the underlying data distribution





# The (distributed) setting

- Users have **data** (observations)  $X_1, \dots, X_n$  (fixed or i.i.d.)
- Users have **constraints**: privacy, bandwidth, ...
- Center has **goals**: e.g., maximize utility given **n** users (minimize number)

# The (distributed) setting

- Users have **data** (observations)  $X_1, \dots, X_n$  (fixed or i.i.d.)
- Users have **constraints**: **privacy**, bandwidth, ...
- Center has **goals**: e.g., maximize utility given **n** users (minimize number)

# What is **private**?

- (Central) Privacy: Trust the **Center**
- Local Privacy: Trust **Nobody**
- Shuffle Privacy: Trust The **Middle Box**

Three variants of  
**Differential Privacy**



# What is **private**?



- (Central) Privacy: Trust the Center
- Local Privacy: Trust **Nobody**
- Shuffle Privacy: Trust The Middle Box

Three variants of  
**Differential Privacy**



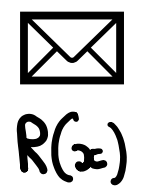
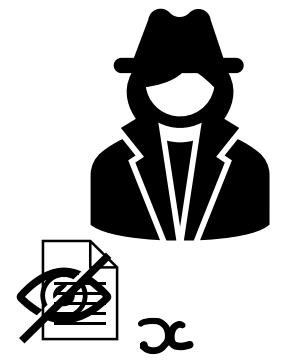
Local

# Differential privacy [DMNS06, KLNRS11]

$\forall x, x', \forall S$


$$\mathbb{P}\{R(x) \in S\} \leq e^\epsilon \mathbb{P}\{R(x') \in S\}$$

$\approx 1 + \epsilon$  (?)



# Histograms

$x_1, \dots, x_n \in [k]$

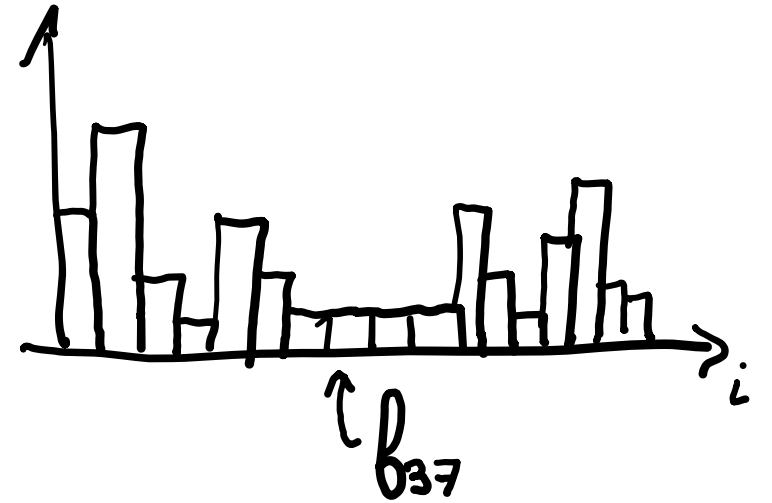
private: user  $j$  holds  $x_j$  

$$f_i = \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{x_j=i}, \quad i \in [k]$$

Goal: center outputs  $\hat{f}_1, \dots, \hat{f}_k$  minimizing

$$\mathbb{E}[\|\hat{f} - f\|_\infty]$$

over protocol's randomness



# Histograms

$$\text{Best } \mathbb{E}[\|\hat{\beta} - \beta\|_{\infty}] = \Theta\left(\sqrt{\frac{\log k}{n \epsilon^2}}\right)$$

# Histograms

$$\text{Best } \mathbb{E}[\|\hat{f} - f\|_\infty] = \Theta\left(\sqrt{\frac{\log k}{n \epsilon^2}}\right)$$

(and this is achieved by several protocols, including some with  $O(\log k)$  bits from each user)

\* e.g., [Acharya-Sun '19]

Done!



# Histograms



$$\text{Best } \mathbb{E}[\|\hat{\beta} - \beta\|_{\infty}] = \Theta\left(\sqrt{\frac{\log k}{n \epsilon^2}}\right)$$

for  $\epsilon \ll 1$

("high privacy")

# Histograms



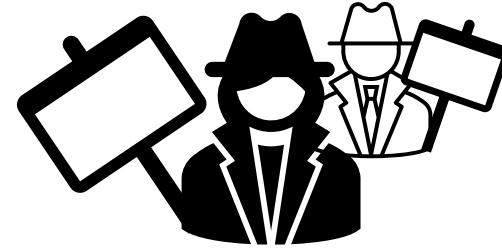
Best  $E[\|\hat{\beta} - \beta\|_{\infty}] = ?$

for  $\epsilon \gg 1$

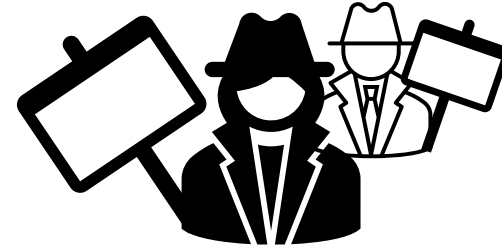
(and why should we care?)

# Histograms: what do we need?

- Simple, efficient protocols
- Low communication
- Non-interactive
- No public randomness
- Good error...
- ... including in the **low-privacy regime**



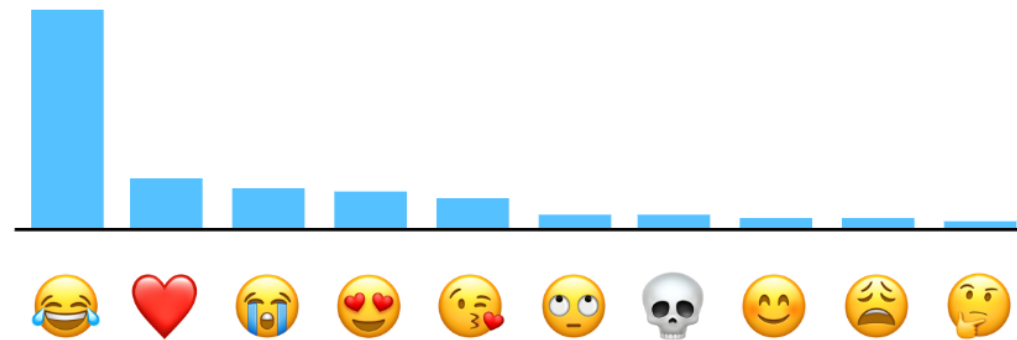
# Histograms: what do we need?



For Lookup Hints, Apple uses a privacy budget with epsilon of 4, and limits user contributions to two donations per day. For emoji, Apple uses a privacy budget with epsilon of 4, and submits one donation per day. For QuickType, Apple uses a privacy budget with epsilon of 8, and submits two donations per day.

For Health types, Apple uses a privacy budget with epsilon of 2 and limits user contributions to one donation per day. The donations do not include health information itself, but rather which health data types are being edited by users.

For Safari, Apple limits user contributions to 2 donations per day. For Safari domains identified as causing high energy use or crashes, Apple uses a single privacy budget with epsilon of 4. For Safari Auto-play intent detection, Apple uses a privacy budget with epsilon of 8.



A detour (?)

Let  $X_1, \dots, X_k$  be  $\sigma^2$ -subgaussian r.v.s.

Then

$$E\left[\max_{1 \leq i \leq k} |X_i|\right] \leq \sqrt{2\sigma^2 \log k}$$

## A detour (?)

Let  $X_1, \dots, X_k$  be  $\sigma^2$ -subgaussian r.v.s.

Then

$$\mathbb{E} \left[ \max_{1 \leq i \leq k} X_i \right] \leq \sqrt{2 \sigma^2 \log k}$$

Proof.  $\forall t > 0, \mathbb{E} \max_i X_i \leq \frac{1}{t} \mathbb{E} \log e^{t \max_i X_i} \leq \frac{1}{t} \log \mathbb{E} \max_i e^{t X_i}$

$$\leq \frac{1}{t} \log \sum_i \mathbb{E} e^{t X_i} \leq \frac{1}{t} \log (n e^{\frac{t^2 \sigma^2}{2}})$$
$$= \frac{\log n}{t} + \frac{\sigma^2}{2} t.$$

Choose  $t = \sqrt{\frac{2 \log n}{\sigma^2}}$ .  $\square$

A detour (?)

Hoeffding's Lemma.

If  $x \in [0, 1]$ , then  $\mathbb{E} e^{t(x - \mathbb{E}x)} \leq e^{\frac{t^2}{8}}$ .

"Bernoulli's are  $\frac{1}{4}$ -subgaussian."

A detour (?)

Hoeffding's Lemma.

If  $x \in [0, 1]$ , then  $\mathbb{E} e^{t(x - \mathbb{E}x)} \leq e^{\frac{t^2}{8}}$ .

"Bernoulli's are  $\frac{1}{4}$ -subgaussian."

Corollary. If  $X = \sum_{i=1}^n X_i$ ,  $X_i$  is  $\frac{n}{4}$ -subgaussian.  
↑  
indep<sup>t</sup> Bernoulli



# Back in track

- RAPPOR:
- user  $j$  encodes  $x_j \in [k]$  as  $e_{x_j} \in \{0, 1\}^k$
  - flips each bit independently w.p.  $\frac{1}{e^{\epsilon/2} + 1}$
  - sends  $Y_j \in \{0, 1\}^k$  to the center

Center: Set  $\hat{b}_i = \left( \frac{1}{n} \sum_{j=1}^n Y_j - \frac{1}{e^{\epsilon/2} + 1} \right) \frac{e^{\epsilon/2} + 1}{e^{\epsilon/2} - 1}, i \in [k]$

# Back in track

- RAPPOR:
- user  $j$  encodes  $x_j \in [k]$  as  $e_{x_j} \in \{0,1\}^k$
  - flips each bit independently w.p.  $\frac{1}{e^{\epsilon/2} + 1}$
  - sends  $Y_j \in \{0,1\}^k$  to the center

Center:

$$\text{Set } \hat{b}_i = \left( \underbrace{\frac{1}{n} \sum_{j=1}^n Y_j}_{\frac{n}{4} \text{ - subgaussian!}} - \frac{1}{e^{\epsilon/2} + 1} \right) \frac{e^{\epsilon/2} + 1}{e^{\epsilon/2} - 1}, i \in [k]$$

# Back in track

- RAPPOR:
- user  $j$  encodes  $x_j \in [k]$  as  $e_{x_j} \in \{0,1\}^k$
  - flips each bit independently w.p.  $\frac{1}{e^{\epsilon/2} + 1}$
  - sends  $y_j \in \{0,1\}^k$  to the center

Center:

$$\text{Set } \hat{\beta}_i = \left( \frac{1}{n} \sum_{j=1}^n y_j - \frac{1}{e^{\epsilon/2} + 1} \right) \frac{e^{\epsilon/2} + 1}{e^{\epsilon/2} - 1}, i \in [k]$$

$$\mathbb{E} \|\hat{\beta} - \beta\|_{\infty} \leq \sqrt{\frac{\log k}{n}} \cdot \frac{e^{\epsilon/2} + 1}{e^{\epsilon/2} - 1}$$

# Back in track

- RAPPOR:
- user  $j$  encodes  $x_j \in [k]$  as  $e_{x_j} \in \{0,1\}^k$
  - flips each bit independently w.p.  $\frac{1}{e^{\epsilon/2} + 1}$
  - sends  $y_j \in \{0,1\}^k$  to the center

Center:

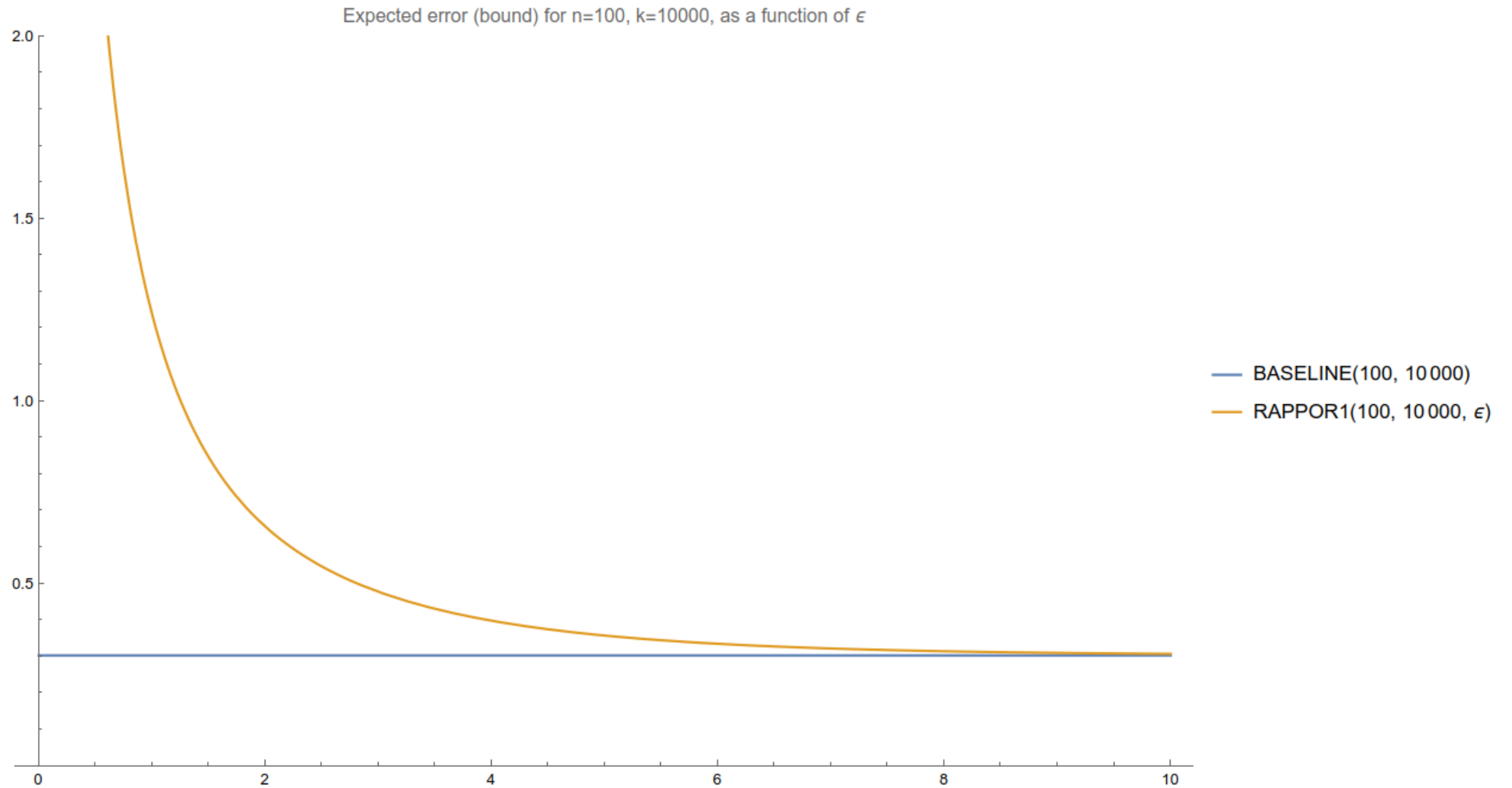
$$\text{Set } \hat{b}_i = \left( \frac{1}{n} \sum_{j=1}^n y_j - \frac{1}{e^{\epsilon/2} + 1} \right) \frac{e^{\epsilon/2} + 1}{e^{\epsilon/2} - 1}, \quad i \in [k]$$

$$\mathbb{E} \|\hat{b} - b\|_\infty \leq \sqrt{\frac{\log k}{n}}$$

$$\frac{e^{\epsilon/2} + 1}{e^{\epsilon/2} - 1}$$

$\frac{1}{\epsilon}$  when  $\epsilon \ll 1$   
1 when  $\epsilon \gg 1$

# It doesn't go to 0...



What can we do?

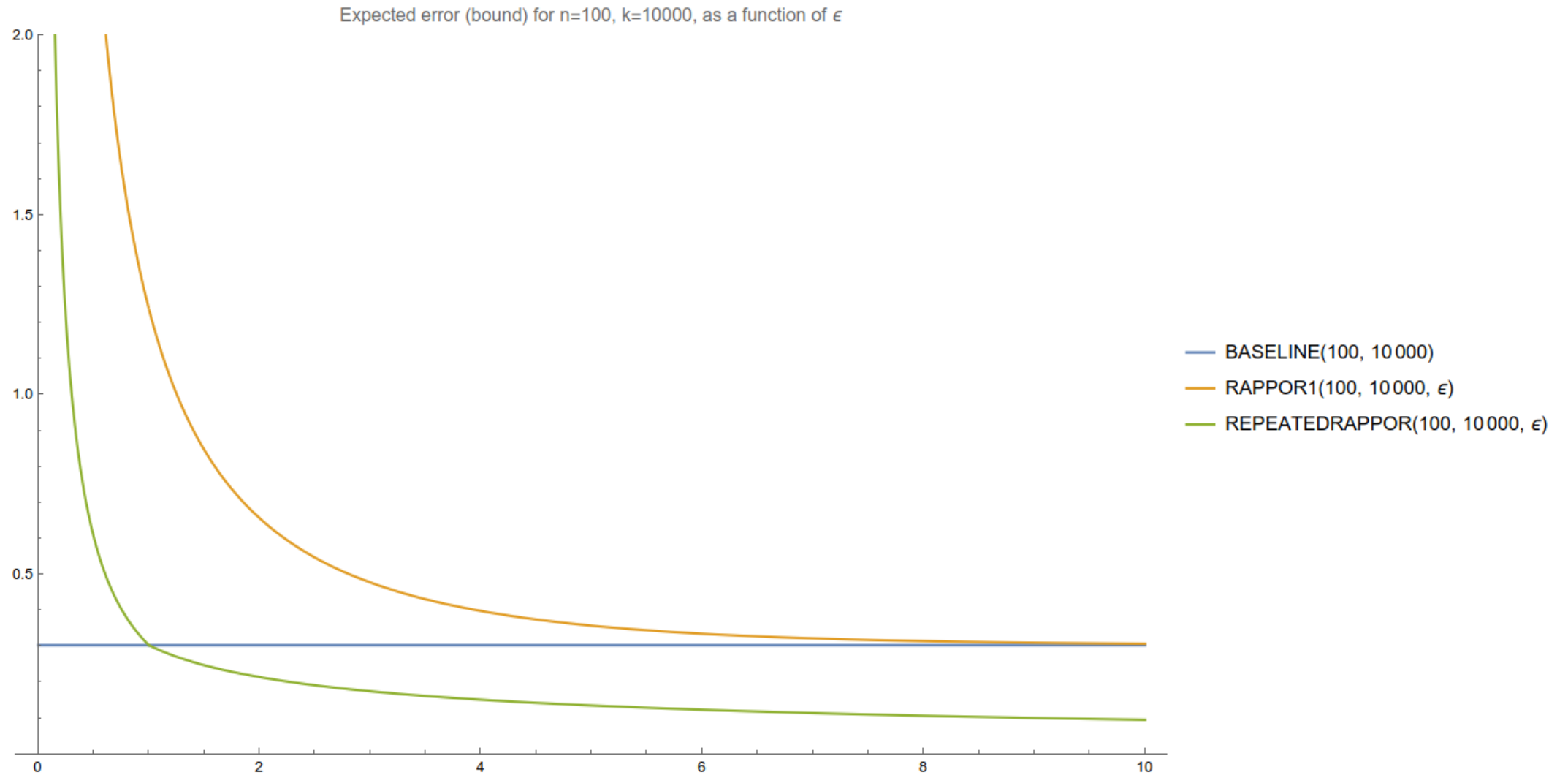
What can we do?

Modify the algorithm: can get

$$O\left(\sqrt{\frac{\log k}{n \min(\epsilon, \epsilon^2)}}\right)$$

via a general "repetition scheme". (Blows up communication by  $\frac{1}{\epsilon}$ )

# It goes to 0!





What else can we do?

Don't Modify the algorithm: can get

$$O\left(\sqrt{\frac{\log k}{n \min(\epsilon, \epsilon^2)}}\right)$$

just by analyzing RAPFOR better!

# Remember Hoeffding?

Kearns - Saul inequality.

If  $X \sim \text{Bern}(p)$ , then

$$\mathbb{E} e^{t(X - \mathbb{E}X)} \leq e^{\frac{1-2p}{4 \log \frac{1-p}{p}} t^2}$$

“Bernoullis are  $\frac{1-2p}{2 \log \frac{1-p}{p}}$  - subgaussian.”

# Remember Hoeffding?

Kearns - Saul inequality.

If  $X \sim \text{Bern}(p)$ , then

$$\mathbb{E} e^{t(X - \mathbb{E}X)} \leq e^{\frac{1-2p}{4 \log \frac{1-p}{p}} t^2}$$

“Bernoullis are  $\frac{1-2}{2 \log \frac{1-p}{p}}$  - subgaussian.”  
 $\underbrace{\hspace{10em}}_{\sigma^2(p)}$

# Remember RAPPOR?

In RAPPOR, each user flips each of its  $k$  bits with proba.  $\frac{1}{e^{\epsilon/2} + 1}$

Each bit is either  $\text{Bern}(p)$  or  $\text{Bern}(1-p)$

Each bit is  $\sigma^2$ -sub-gaussian for

$$\sigma^2 = \sigma^2(p) = \sigma^2(1-p)$$

$$\underbrace{\frac{1}{e^{\epsilon/2} + 1}}_{=p}$$

# Remember RAPPOR?

In RAPPOR, each user flips each of its  $k$  bits with proba.  $\frac{1}{e^{\epsilon/2} + 1}$

Each bit is either  $\text{Bern}(p)$  or  $\text{Bern}(1-p)$

Each bit is  $\sigma^2$ -sub-gaussian for

$$\sigma^2 = \sigma^2(p) = \sigma^2(1-p) = \frac{e^{\epsilon/2} - 1}{e^{\epsilon/2} + 1} \cdot \frac{1}{\epsilon}$$

## Remember RAPPOR?

In RAPPOR, each user flips each of its  $k$  bits with proba.  $\frac{1}{e^{\epsilon/2} + 1}$

This gives error bound

$$\mathbb{E} \|\hat{\beta} - \beta\|_{\infty} \leq \sqrt{\frac{e^{\epsilon/2} + 1}{(e^{\epsilon/2} - 1)\epsilon} \cdot \frac{\log k}{n}}$$

# Remember RAPPOR?

In RAPPOR, each user flips each of its  $k$  bits with proba.  $\frac{1}{e^{\epsilon/2} + 1}$

This gives error bound

$$\mathbb{E} \|\hat{\beta} - \beta\|_{\infty} \leq \sqrt{\frac{e^{\epsilon/2} + 1}{(e^{\epsilon/2} - 1)\epsilon} \cdot \frac{\log k}{n}}$$

$$\sqrt{\frac{\log k}{n\epsilon^2}} \quad \text{when } \epsilon \ll 1$$

$$\sqrt{\frac{\log k}{n\epsilon}} \quad \text{when } \epsilon \gg 1$$

So...



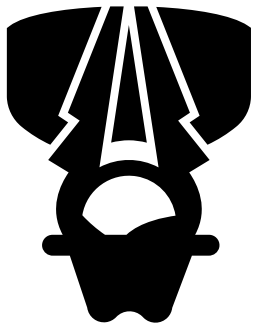


We couldn't prove a matching lower bound...



We couldn't prove a matching lower bound...

So we tried to improve the upper bound.



# Local Glivenko-Cantelli [CK23]

Let  $\mu$  be a product distribution over  $\{0,1\}^d$  with mean vector  $p \in [0,1]^d$ . Define  $\tilde{p}_i = \min(p_i, 1-p_i)$ , and assume  $\log \tilde{p}_i \geq -\tilde{p}_i$ .

Given  $X_1, \dots, X_n \sim \mu$ , setting

$$\hat{p} = \frac{1}{n} \sum_{j=1}^n X_j$$

we have

$$\mathbb{E}[\|\hat{p} - p\|_\infty] \lesssim \max_{1 \leq i \leq d} \sqrt{\frac{\tilde{p}_i \log(i+1)}{n}} + \frac{\log n}{n} \max_{1 \leq i \leq d} \frac{\log(i+1)}{\log \frac{1}{\tilde{p}_i}}$$

# Local Glivenko-Cantelli [CK23] + RAPPOR

For  $\omega$ ,  $p_i \in \{p, 1-p\} \forall i$ , so  $\tilde{p}_i = p = \frac{1}{e^{\epsilon/2} + 1}$ . We immediately\* get

$$\mathbb{E} \|\hat{f} - f\|_{\infty} \leq \frac{e^{\epsilon/2} + 1}{e^{\epsilon/2} - 1} \left( \sqrt{\frac{\log k}{n(e^{\epsilon/2} + 1)}} + \frac{\log k}{n \max(\epsilon, 1)} \right)$$

$$\stackrel{\epsilon \gg 1}{\sim} \sqrt{\frac{\log k}{n e^{\epsilon/2}}} + \frac{\log k}{n \epsilon}$$

# Local Glivenko-Cantelli [CK23] + RAPPOR

For  $\omega$ ,  $p_i \in \{p, 1-p\} \forall i$ , so  $\tilde{p}_i = p = \frac{1}{e^{\epsilon/2} + 1}$ . We immediately\* get

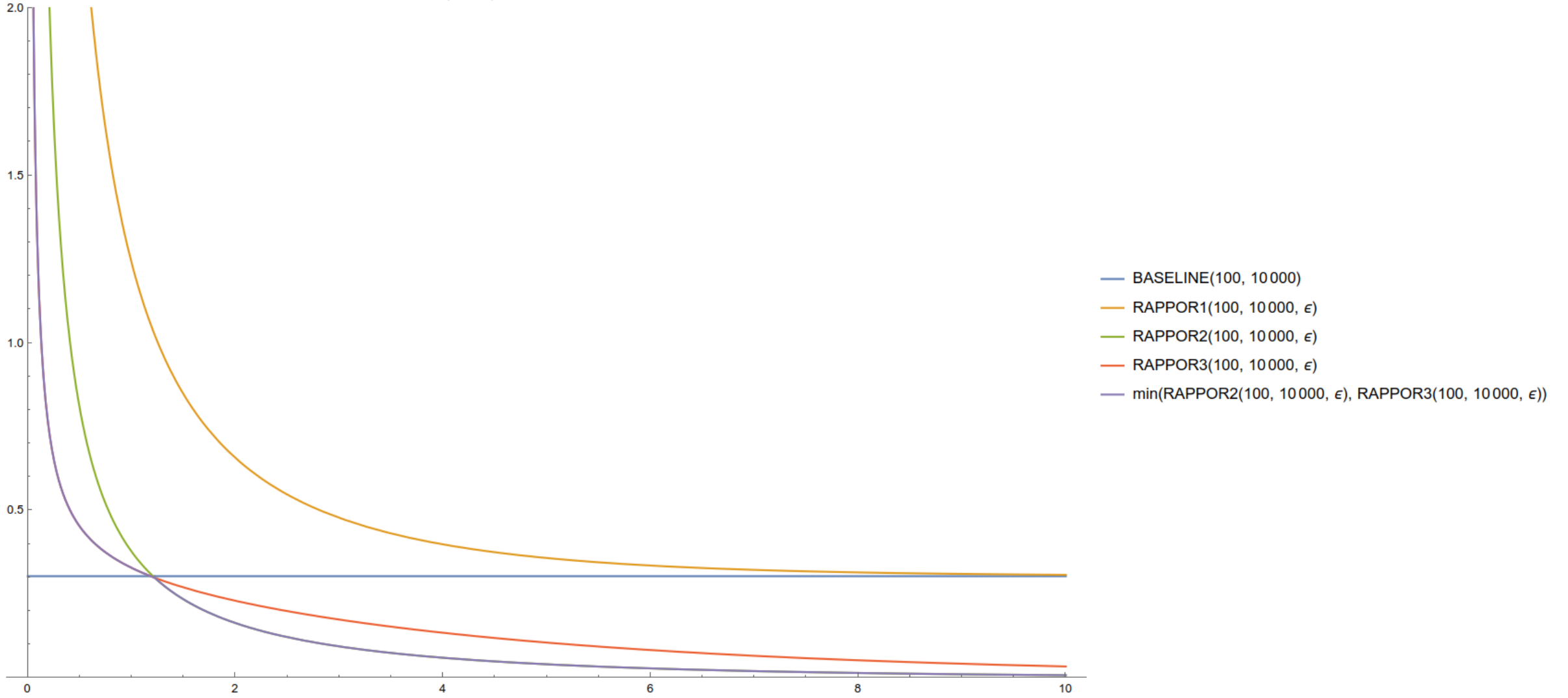
$$\mathbb{E} \|\hat{f} - f\|_{\infty} \leq \frac{e^{\epsilon/2} + 1}{e^{\epsilon/2} - 1} \left( \sqrt{\frac{\log k}{n(e^{\epsilon/2} + 1)}} + \frac{\log k}{n \max(\epsilon, 1)} \right)$$

$$\stackrel{\epsilon \gg 1}{\sim} \sqrt{\frac{\log k}{n e^{\epsilon/2}}} + \frac{\log k}{n \epsilon}$$

\*This is a lie.

# We did nothing, and it got better!

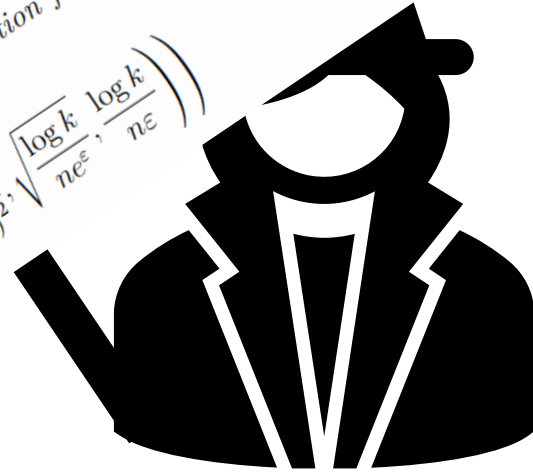
Expected error (bound) for  $n=100$ ,  $k=10000$ , as a function of  $\epsilon$



Also, this is optimal\*

**Theorem 7.** Fix any  $\epsilon > 0$ . Any non-interactive (public- or private-coin) protocol  $\Pi$  for distribution estimation from  $n$  users must have  $\min_{\max}$  expected  $l_\infty$  error

$$\Omega\left(\max\left(\sqrt{\frac{\log k}{n(e^\epsilon - 1)^2}}, \sqrt{\frac{\log k \log k}{n e^\epsilon}}, \frac{\log k}{n \epsilon}\right)\right)$$



\*This is a lie.

# Summary

- Algorithms can "improve" without any modification
- Can we go beyond RAPPOR? (**yes**: work in progress)
- Even for "solved" problems, there is **so much** we don't know!
- Trying to prove lower bounds is useful
- Not mentioned here: beyond worst-case, and other protocols

Thank you.







# The plot thickens

RAPPOR  $l(\infty)$  by  $\epsilon$

