# Sample Complexity of Differential Privacy

## Mark Bun*

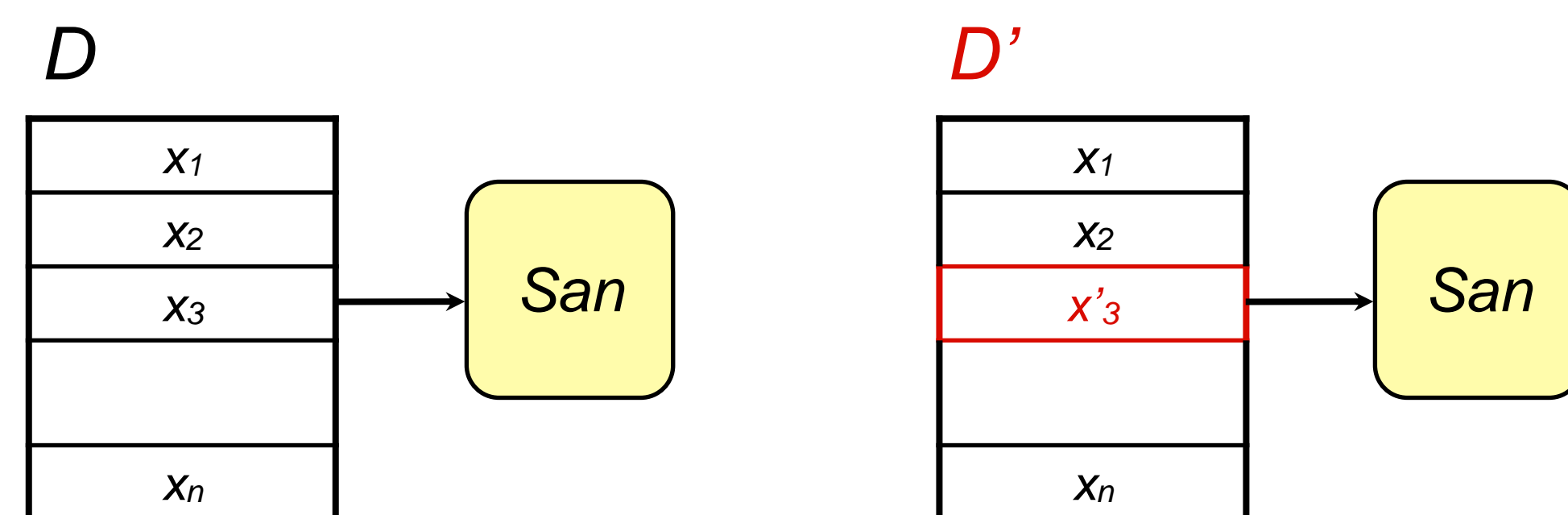### Harvard University (2nd year Ph.D., supported by an NDSEG fellowship)

---

## MAIN QUESTION

How many **data samples** do we need to achieve both **differential privacy** and **statistical accuracy**?

i.e. How big a study do we need to conduct to answer our questions and preserve privacy?

## DIFFERENTIAL PRIVACY



$D$ and $D'$ are neighbors if they differ only on one user's data

An algorithm *San* is (ε,δ)-differentially private if for all neighbors $D$, $D'$ and every $S \subseteq$ Range(*San*),

$$\Pr[San(D) \in S] \leq e^{\varepsilon} \Pr[San(D') \in S] + \delta$$

Think of $\varepsilon = \Theta(1)$ and $\delta = o(1/n)$

## ACCURACY FOR COUNTING QUERIES

Counting queries: What fraction of rows in a database satisfy property $q$?

e.g. $q(x) = $ LikesBread AND LikesToast

| LikesBread? | LikesButter? | LikesToast? | LikesJam? | |
|---|---|---|---|---|
| 0 | 0 | 1 | 1 | $q(x_1)=0$ |
| 1 | 1 | 1 | 1 | $q(x_2)=1$ |
| 1 | 0 | 1 | 0 | $q(x_3)=1$ |

$d$ (=4) attributes per record

$q(D)=2/3$

Answers $a_q$ are α-accurate if $|a_q - q(D)| < \alpha$ for every $q \in Q$
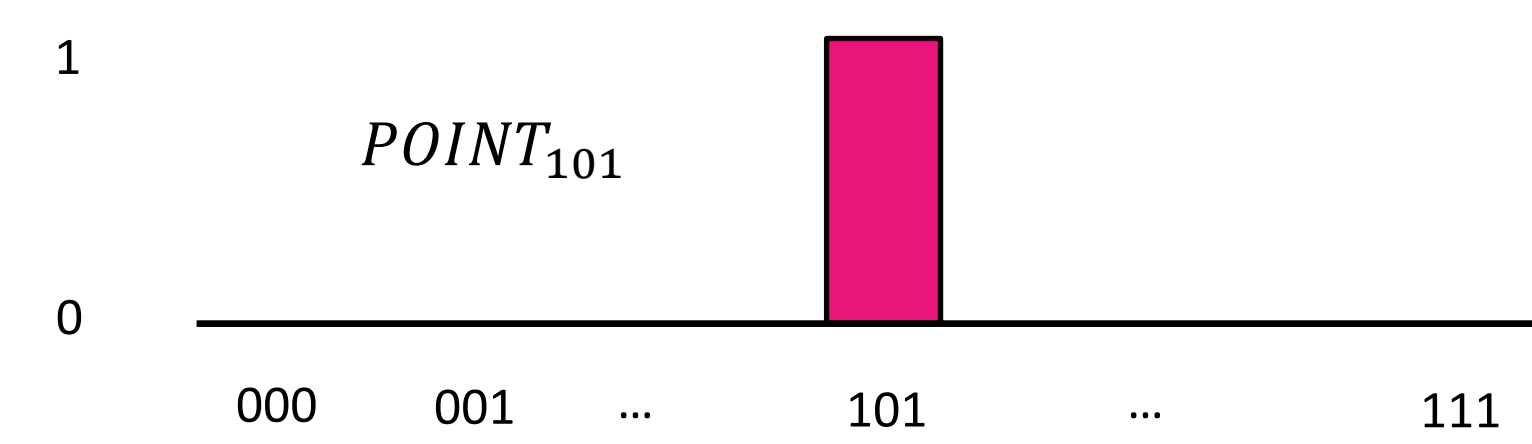
## SAMPLE COMPLEXITY UPPER BOUNDS

For general queries $Q$,

$$O(\sqrt{d} \log |Q| / \alpha^2)$$

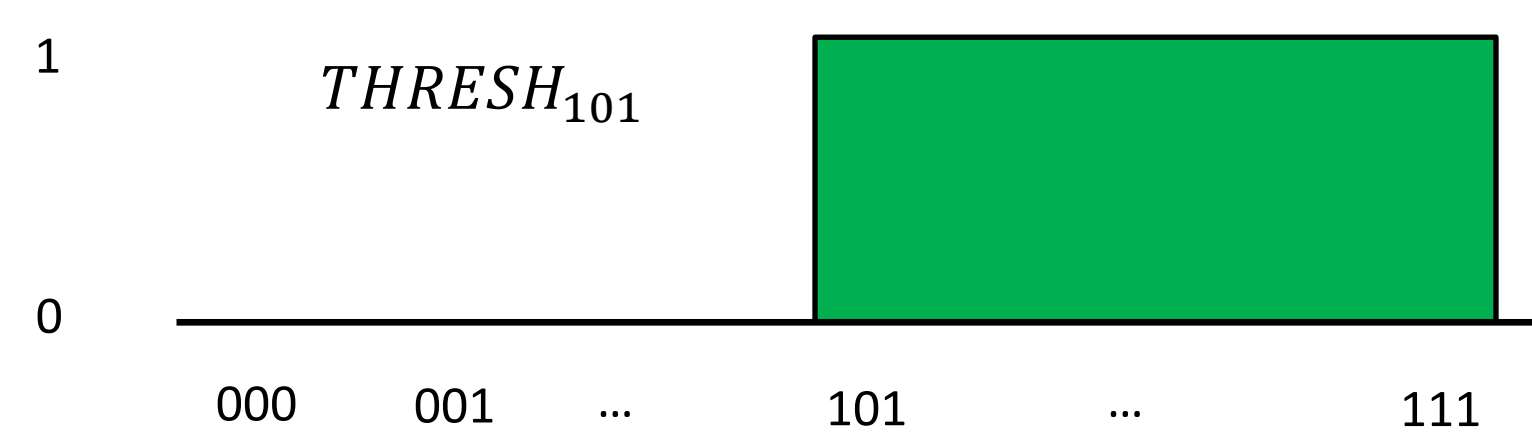samples suffice [HR10], using the analysis in [GRU12]

But for certain $Q$, the sample complexity can be much lower:

**Point queries:** $POINT_y(x) = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{otherwise} \end{cases}$



$\log |Q| = d$, but just $O(1/\alpha)$ samples suffice

**Threshold queries:** $THRESH_y(x) = \begin{cases} 1 & \text{if } x \geq y \\ 0 & \text{otherwise} \end{cases}$



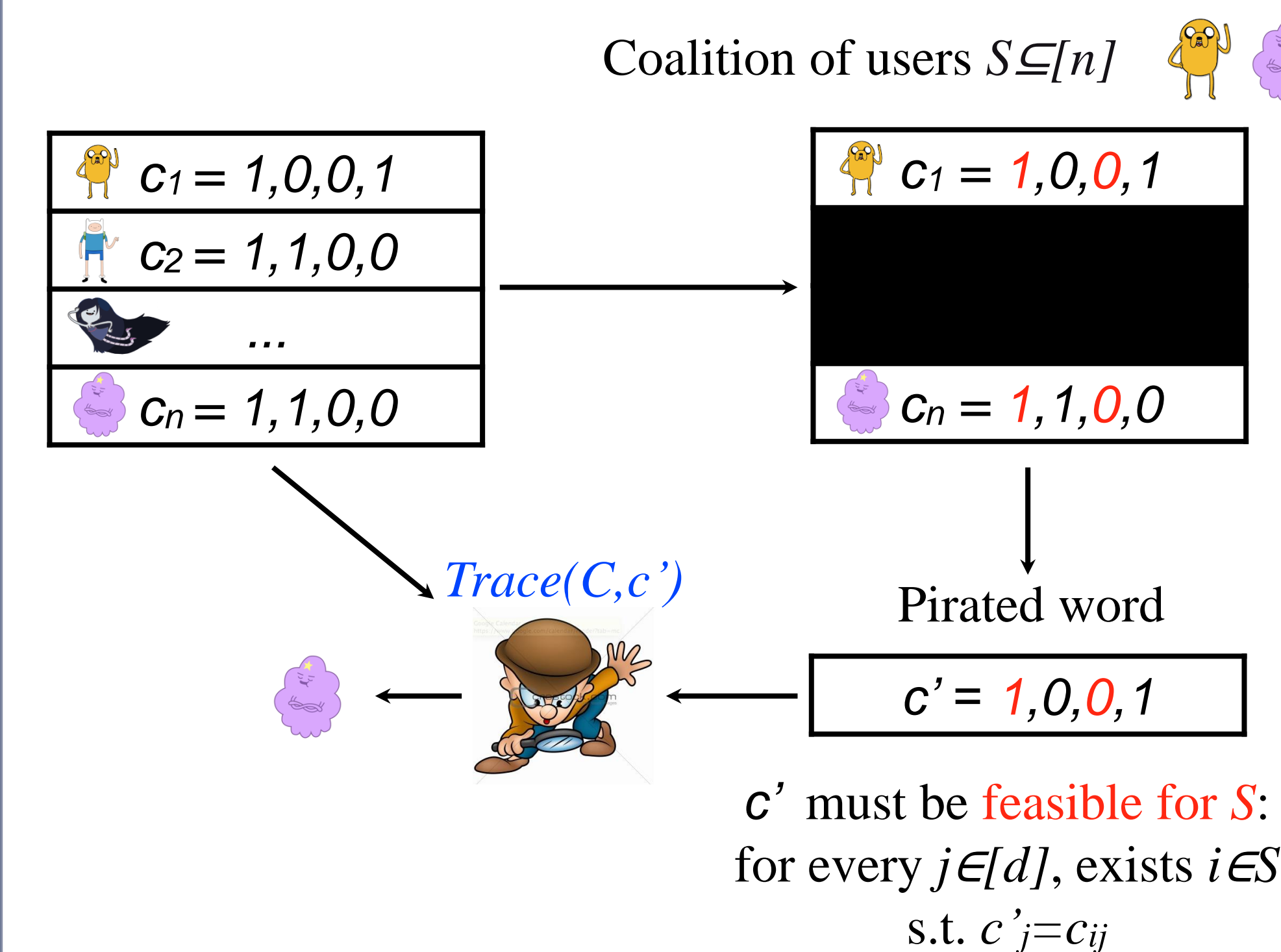Again, $\log |Q| = d$, but $<< d /\alpha^{2.5}$ samples suffice.[BNS13]

- Extend to upper bounds on the sample complexity of differentially private *PAC learning*.

- Sample complexity is much smaller than what is needed for pure (i.e. $\delta = 0$) privacy.

- Relevant quantity seems to be the VC-Dimension of $Q$

## SAMPLE COMPLEXITY LOWER BOUNDS
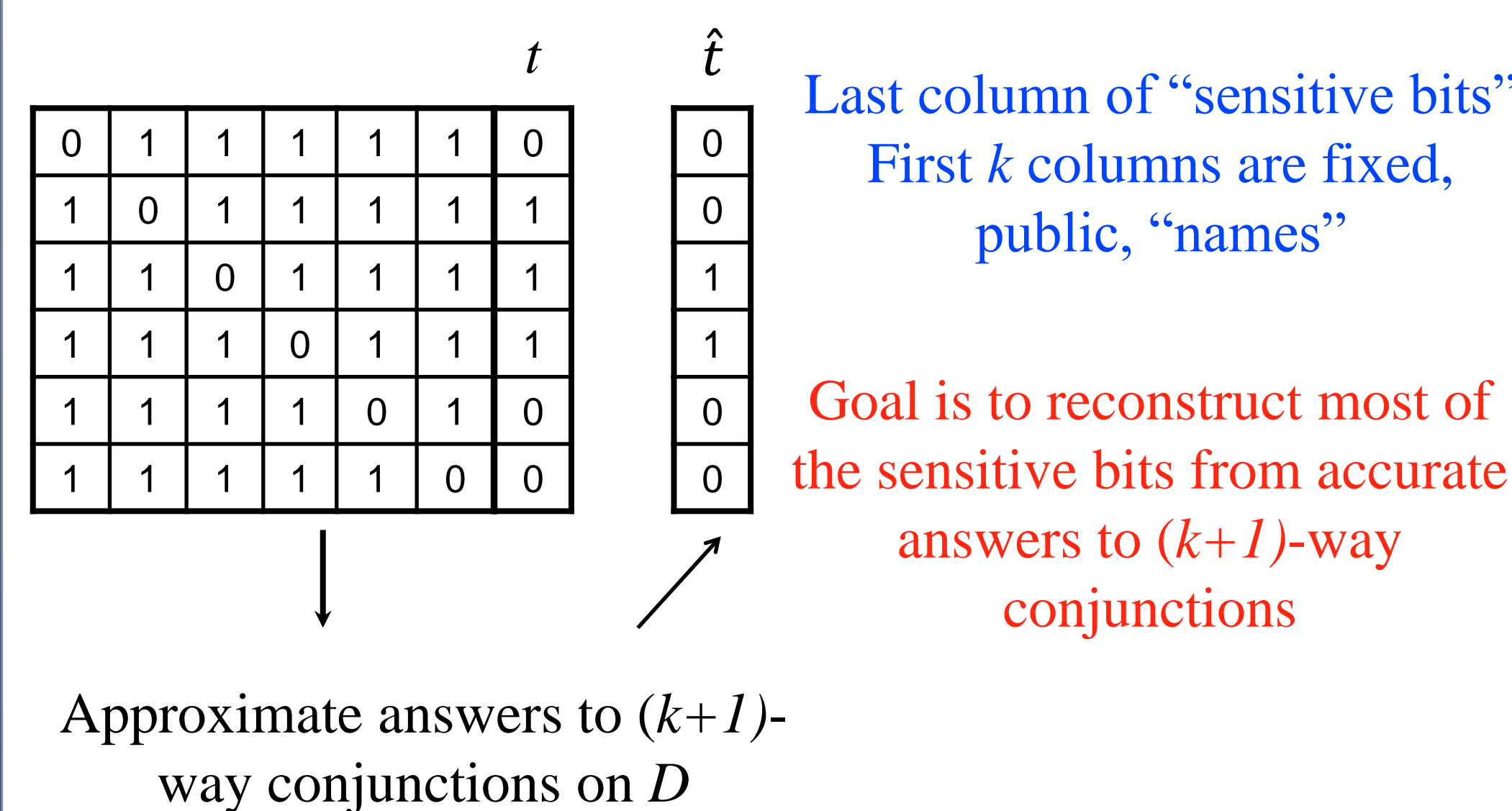
Our contributions [BUV13]

- To answer *arbitrary* queries, $\Omega(\sqrt{d} \log |Q| / \alpha^2)$ samples are necessary (nearly tight)

- If $\alpha$ is a constant, this lower bound still holds for *conjunction* queries

### Tool 1: Fingerprinting Codes

Coalition of users $S \subseteq [n]$



$c'$ must be feasible for $S$: for every $j \in [d]$, exists $i \in S$ s.t. $c'_j = c_{ij}$

- Sensitive database = traceable codebook
- Traceability is the "opposite" of privacy
- Yields a lower bound of $\Omega(\sqrt{d})$ for estimating the mean of each column

### Tool 2: Reconstruction Attacks [DN03]



Last column of "sensitive bits" First $k$ columns are fixed, public, "names"

Goal is to reconstruct most of the sensitive bits from accurate answers to $(k+1)$-way conjunctions

Approximate answers to $(k+1)$-way conjunctions on $D$

## COMPOSITION OF LOWER BOUNDS



$D_1 \in (\{0,1\}^d)^m$

$D_2 \in (\{0,1\}^d)^m$

Random stack of "sensitive databases"
First $k$ columns are public, fixed "names" for each $D_i$

$D_k \in (\{0,1\}^d)^m$

Goal is to answer (most) *1*-way conj's on at least one $D_i$ ⇒ privacy breach

$(k+1)$-way conj's compute "subset sums of *1*-way conj's"

### REFERENCES

[BUV13] Mark Bun, Jonathan Ullman, and Salil Vadhan. Fingerprinting codes and the price of approximate differential privacy. *Manuscript*, 2013.

[BNS13] Amos Beimel, Kobbi Nissim, and Uri Stemmer. Private learning and sanitization: pure vs. approximate differential privacy. In *RANDOM*, 2013.

[DN03] Irit Dinur and Kobbi Nissim. Revealing information while preserving privacy. In PODS, 2003.

[GRU10] Anupam Gupta, Aaron Roth, and Jonathan Ullman. Iterative constructions and private data release. In TCC, 2012.

[HR10] Moritz Hardt and Guy Rothblum. A multiplicative weights mechanism for privacy-preserving data analysis. In FOCS, 2010.

### CONTACT

mbun@seas.harvard.edu

Harvard School of Engineering and Applied Sciences
Maxwell Dworkin 138
33 Oxford St.
Cambridge, MA 02138