

# Fingerprinting Codes and the Price of Approximate Differential Privacy

June 1, 2014

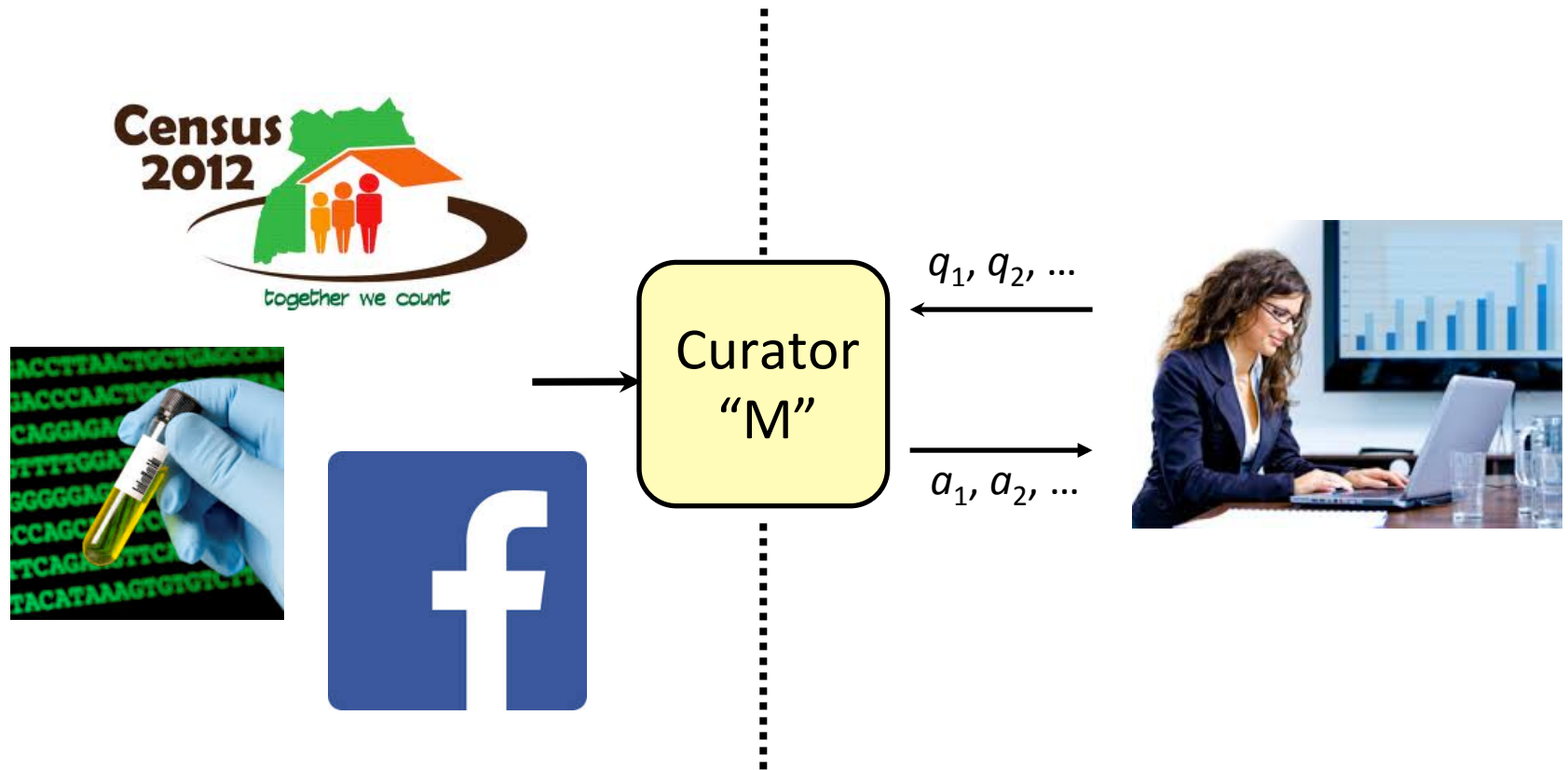
*Mark Bun*

Jonathan Ullman

Salil Vadhan

Harvard University

# Privacy-Preserving Data Analysis



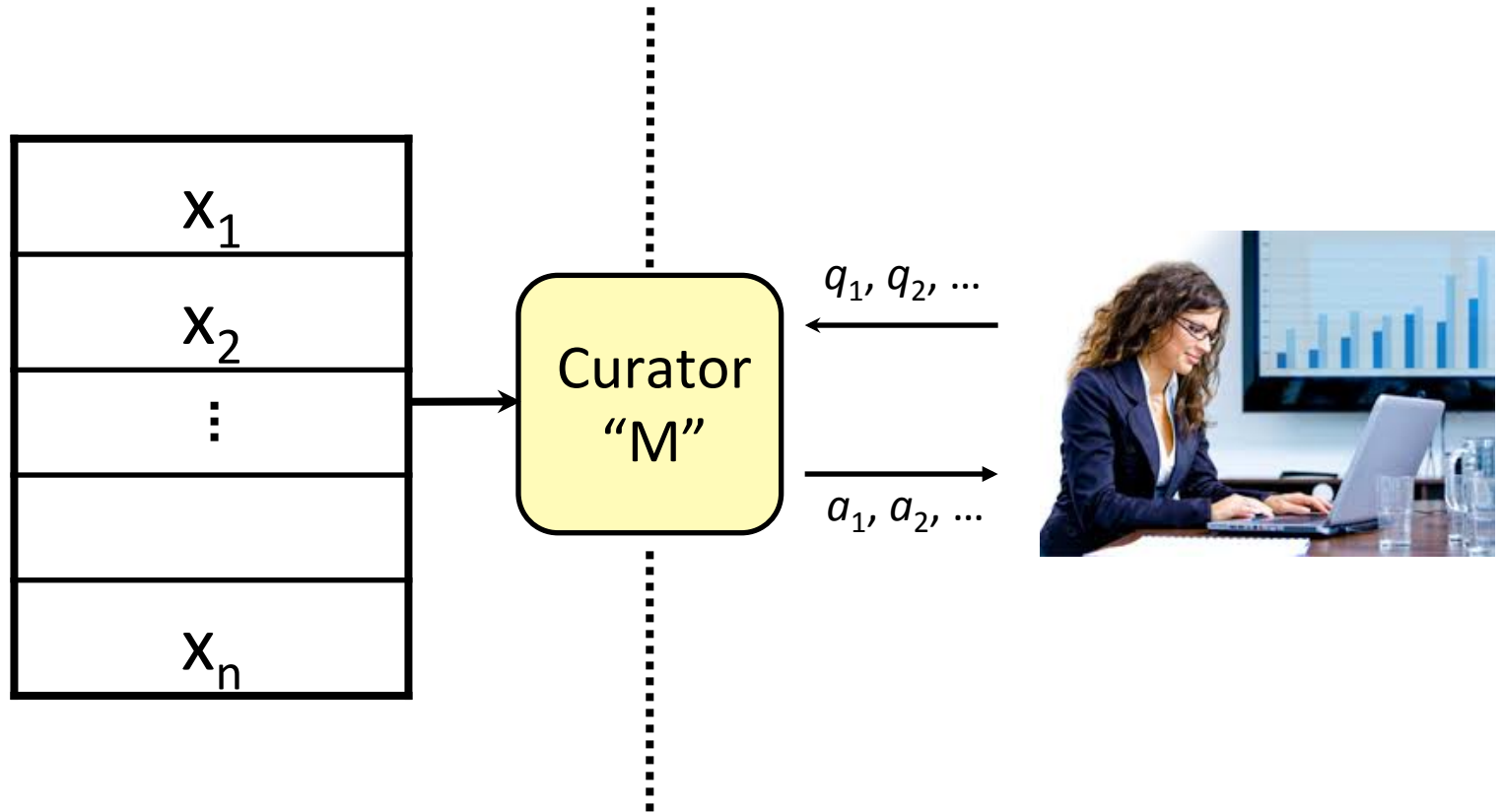
Want curators that are:

◆ Private

◆ Accurate

◆ Efficient

# Privacy-Preserving Data Analysis



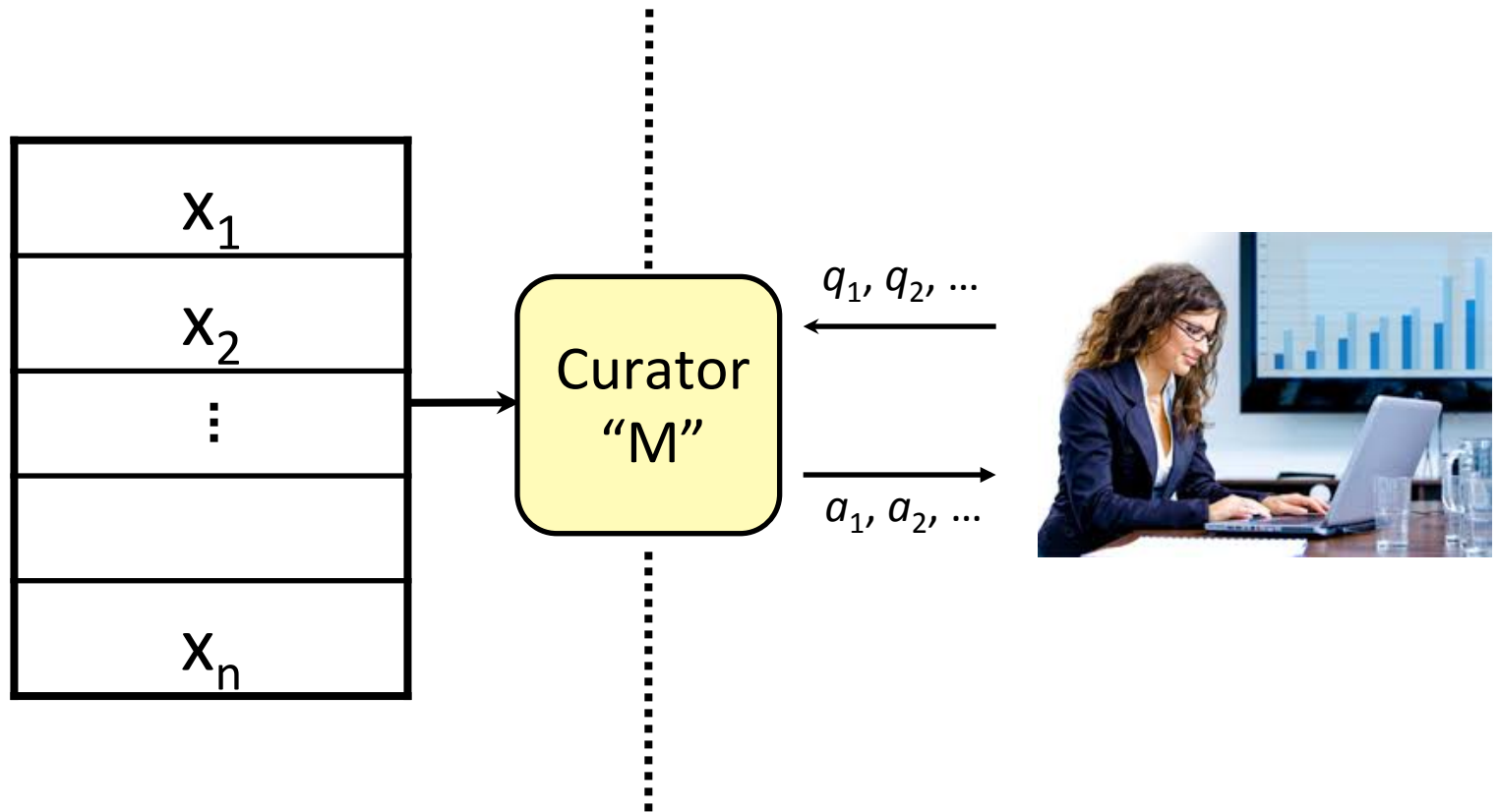
Want curators that are:

◆ Private

◆ Accurate

◆ Efficient

# Privacy-Preserving Data Analysis



Want curators that are:

- ◆ Differentially Private
- ◆ Statistically Accurate
- ◆ Sample Efficient

# What This Talk is About

- **Sample complexity** for approx. differential privacy
- **MAIN RESULT**: For **high-dimensional data**, **Privacy + Accuracy** requires more samples than **Accuracy** alone

e.g.  $d$  attribute means

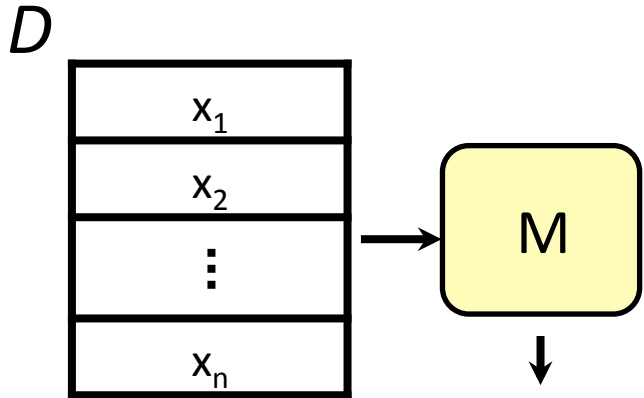
Accuracy:  $\Theta(\log d)$

Privacy + Accuracy:  $\tilde{\Theta}(d^{1/2})$

- New techniques for privacy lower bounds

# Differential Privacy

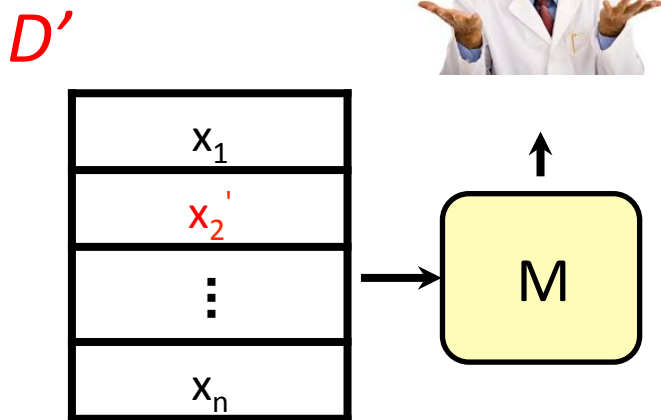
[DN03+Dwork, DN04, BDMN05, DMNS06, DKMMN06]



$D$  and  $D'$  are **neighbors** if they differ on one row

small const., e.g.  $\epsilon = 0.1$

“cryptographically small”  
need  $\delta \ll 1/n$ , often  $\delta = \text{negl}(n)$



$M$  is  **$(\epsilon, \delta)$ -differentially private** if for all neighbors  $D, D'$  and  $T \subseteq \text{Range}(M)$ :

$$\Pr[M(D') \in T] \leq (1+\epsilon)\Pr[M(D) \in T] + \delta$$

◆ Privacy

◆ Accuracy

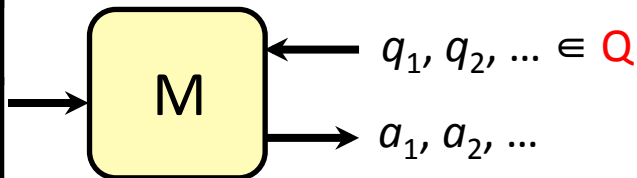
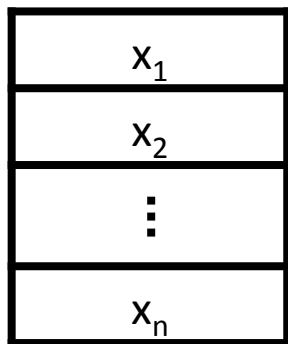
◆ Sample Complexity

# Counting Queries

“What fraction of the rows of  $D$  satisfy some property  $q$ ?”

E.g. **attribute** means  
 $q = \text{Skywalker?}$   
 $q(D) = 3/4$

	DarkSide?	Twin?	Skywalker?	< 3ft?
	0	0	0	1
	0	1	1	0
	0	1	1	0
	1	0	1	0



M is  **$\alpha$ -accurate** for  **$Q$**  if  
 $|a_i - q_i(D)| < \alpha$  for every  $i$

◆ Privacy

◆ Accuracy

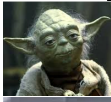
◆ Sample Complexity

# (Privately) Answering Attribute Means

[DN03, DN04, BDMN05, DMNS06]

$d$  binary attributes

$n$  rows

	DarkSide?	Twin?	Skywalker?	< 3ft?
	0	0	0	1
	0	1	1	0
	0	1	1	0
	1	0	1	0

$$\frac{3}{4} + \text{Noise}(O(1/n))$$

( $\alpha$ -accuracy requires  $n \geq 1/\alpha$ )

◆ Privacy

◆ Accuracy

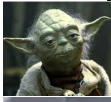



◆ Sample Complexity



# (Privately) Answering Attribute Means

[DN03, DN04, BDMN05, DMNS06]

$d$  binary attributes

		DarkSide?	Twin?	Skywalker?	< 3ft?
$n$ rows		0	0	0	1
		0	1	1	0
		0	1	1	0
		1	0	1	0
		1/4 +	1/2 +	3/4 +	1/4 +
		Noise( $O(d^{1/2}/n)$ )	Noise( $O(d^{1/2}/n)$ )	Noise( $O(d^{1/2}/n)$ )	Noise( $O(d^{1/2}/n)$ )

( $\alpha$ -accuracy requires  $n \geq d^{1/2}/\alpha$ )

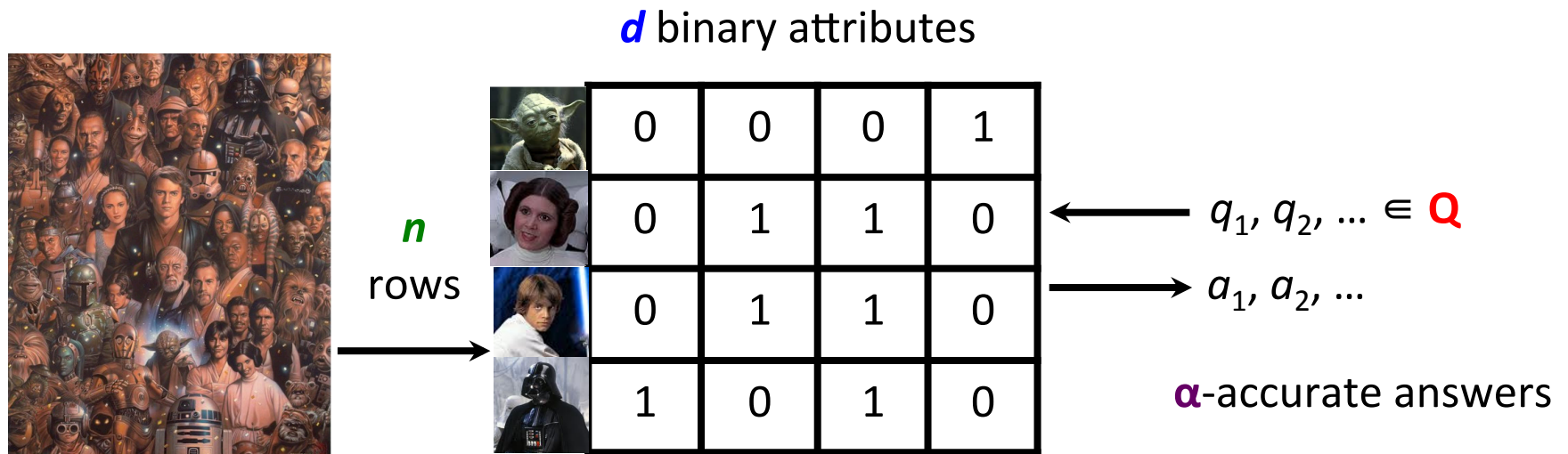
◆ Privacy

◆ Accuracy

◆ Sample Complexity

# Sample Complexity

How big does  $n$  have to be to guarantee statistical accuracy *on the population*?



◆ Privacy

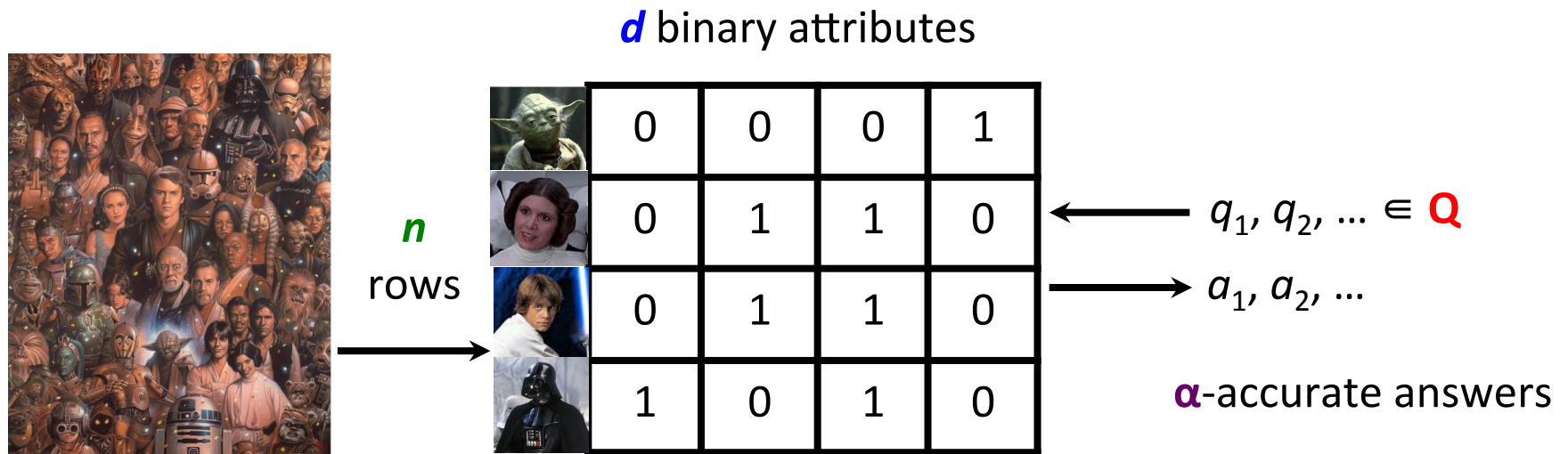
◆ Accuracy

◆ Sample Complexity

# Sample Complexity

Answer:  $n = \Theta(\log |\mathbf{Q}| / \alpha^2)$  [Vap98]

e.g.  $\Theta(\log d)$  for attribute means with  $\alpha = 0.05$



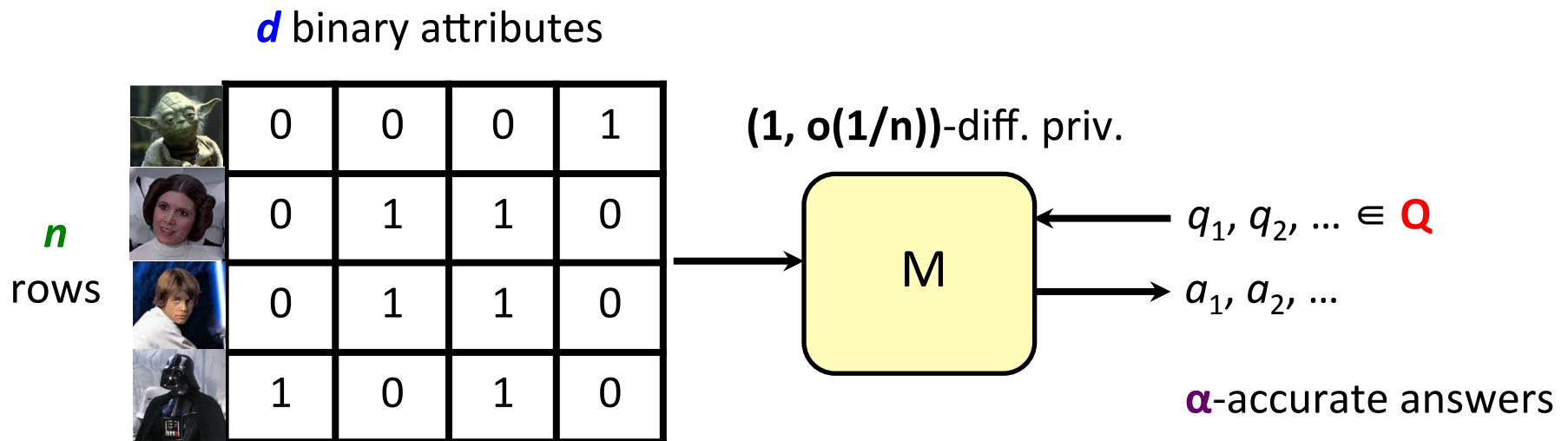
◆ Privacy

◆ Accuracy

◆ Sample Complexity

# Sample Complexity for Diff. Privacy

How big does  $n$  have to be to guarantee accuracy *and* privacy?



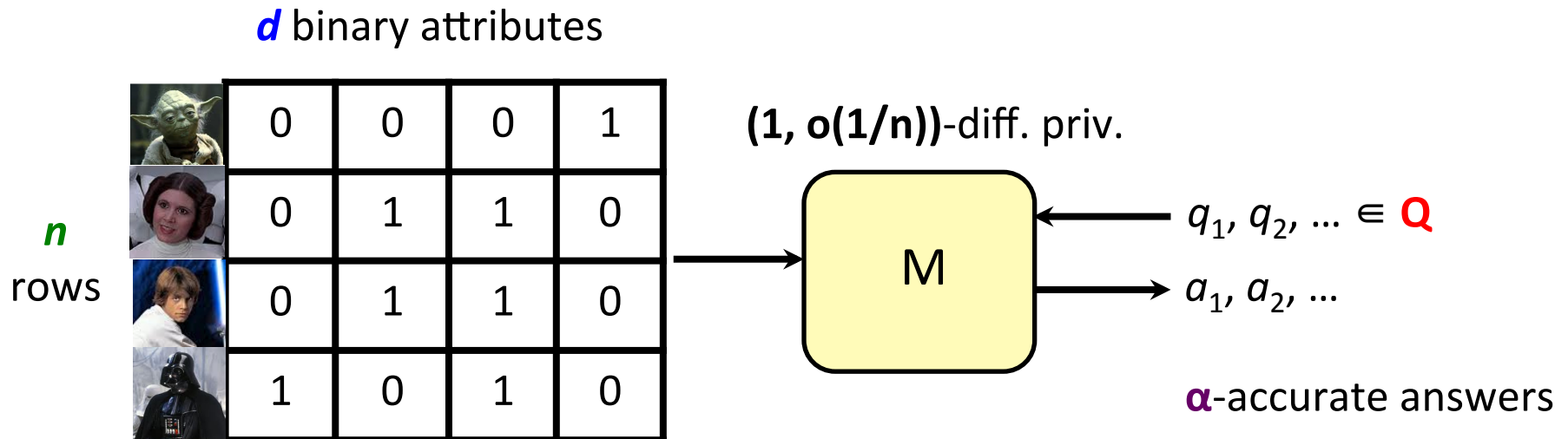
◆ Privacy

◆ Accuracy

◆ Sample Complexity

# Sample Complexity for Diff. Privacy

Question: Is there an additional **price of diff. privacy** over statistical accuracy alone?



◆ Privacy

◆ Accuracy

◆ Sample Complexity

# Sample Complexity for Diff. Privacy

No privacy

$\mathcal{Q}$  = attribute means  
 $\alpha = 0.05$

$\mathcal{Q}, \alpha$  arbitrary

$n = \Theta(\log d)$ [Vap98]	$n = \Theta(\log  \mathcal{Q}  / \alpha^2)$ [Vap98]
---------------------------------	--

$(0.1, o(1/n))$ -  
 diff. privacy

Upper bound:

$\tilde{O}(d^{1/2})$ [...DMNS06]	$\forall \mathcal{Q}: \tilde{O}(\log  \mathcal{Q}  \cdot d^{1/2} / \alpha^2)$ [HR10]
-------------------------------------	---

Lower bound:

$\tilde{\Omega}(\log d)$ [DN03, Rot10]	$\exists \mathcal{Q}: \max \tilde{\Omega}(\log  \mathcal{Q}  / \alpha), \tilde{\Omega}(1 / \alpha^2)$ [DN03]
---	---

**OUR WORK:**

$\tilde{\Omega}(d^{1/2})$	$\exists \mathcal{Q}: \tilde{\Omega}(\log  \mathcal{Q}  \cdot d^{1/2} / \alpha^2)$
---------------------------	--

◆ Privacy

◆ Accuracy

◆ Sample Complexity

# Beyond Reconstruction Attacks

- Tight lower bounds known for  $(\epsilon, 0)$ -diff. privacy [HT10, Har11], but break even for  $\delta = \text{negl}(n)$  [De11, BNS13]
- Prior lower bounds for  $(\epsilon, \delta)$ -diff. privacy gave **reconstruction attacks** [DN03, Rot10], which hold even for  $\delta = \text{constant}$
- **This work:** Fingerprinting codes enable optimal lower bounds for  $(\epsilon, \delta=o(1/n))$ -diff. privacy (followed by [DTTZ14, BST14])

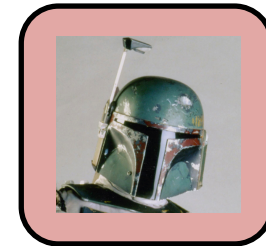
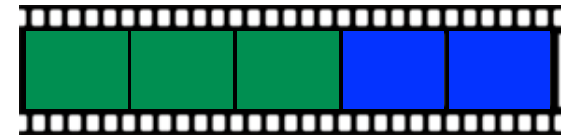
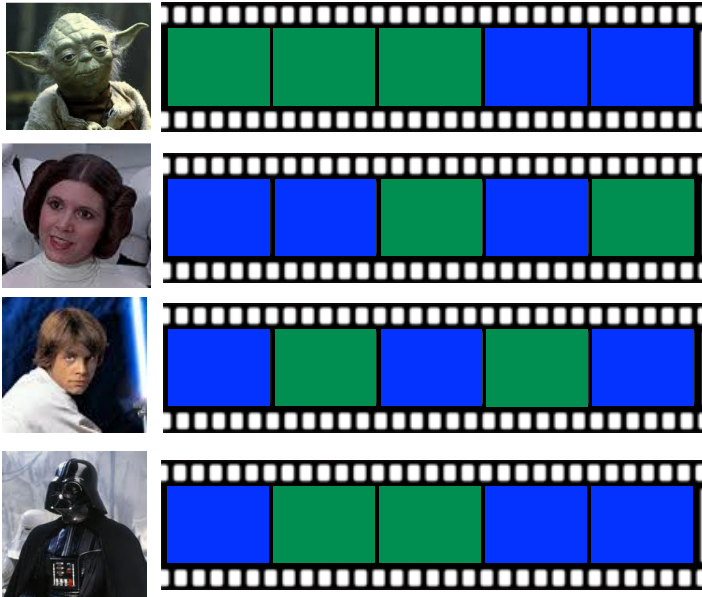
# New Techniques

- Fingerprinting codes  $\rightarrow$  diff. privacy lower bounds
  - $\tilde{\Omega}(d^{1/2})$  for attribute means ( $\alpha$  const.)
- Composition of sample complexity lower bounds
  - $\tilde{\Omega}(kd^{1/2})$  for  $k$ -way conjunctions ( $\alpha$  const.)
  - $\tilde{\Omega}(\log |Q| \cdot d^{1/2} / \alpha^2)$  for arbitrary queries

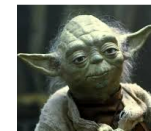


# Fingerprinting Codes [BS95]

I want to distribute my new movie



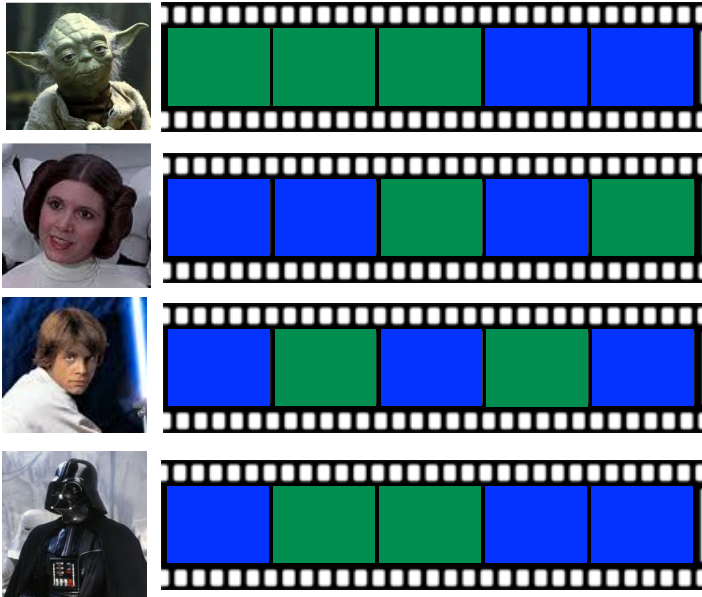
Trace Algorithm



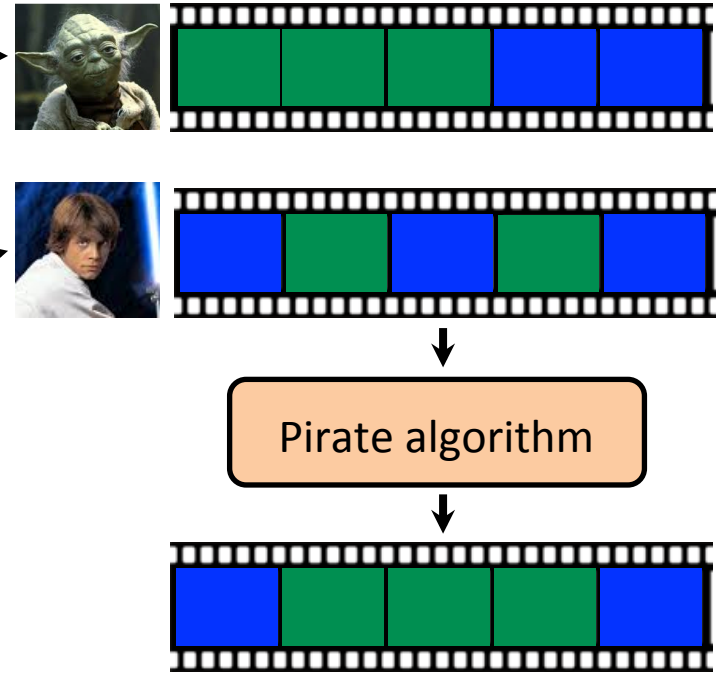
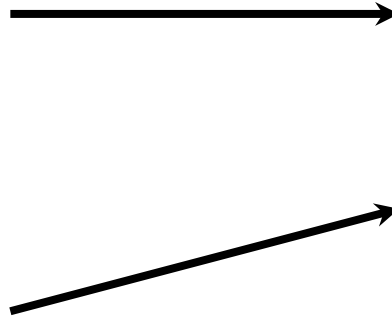
...but the galaxy is full of pirates!

# Fingerprinting Codes [BS95]

I want to distribute my new movie

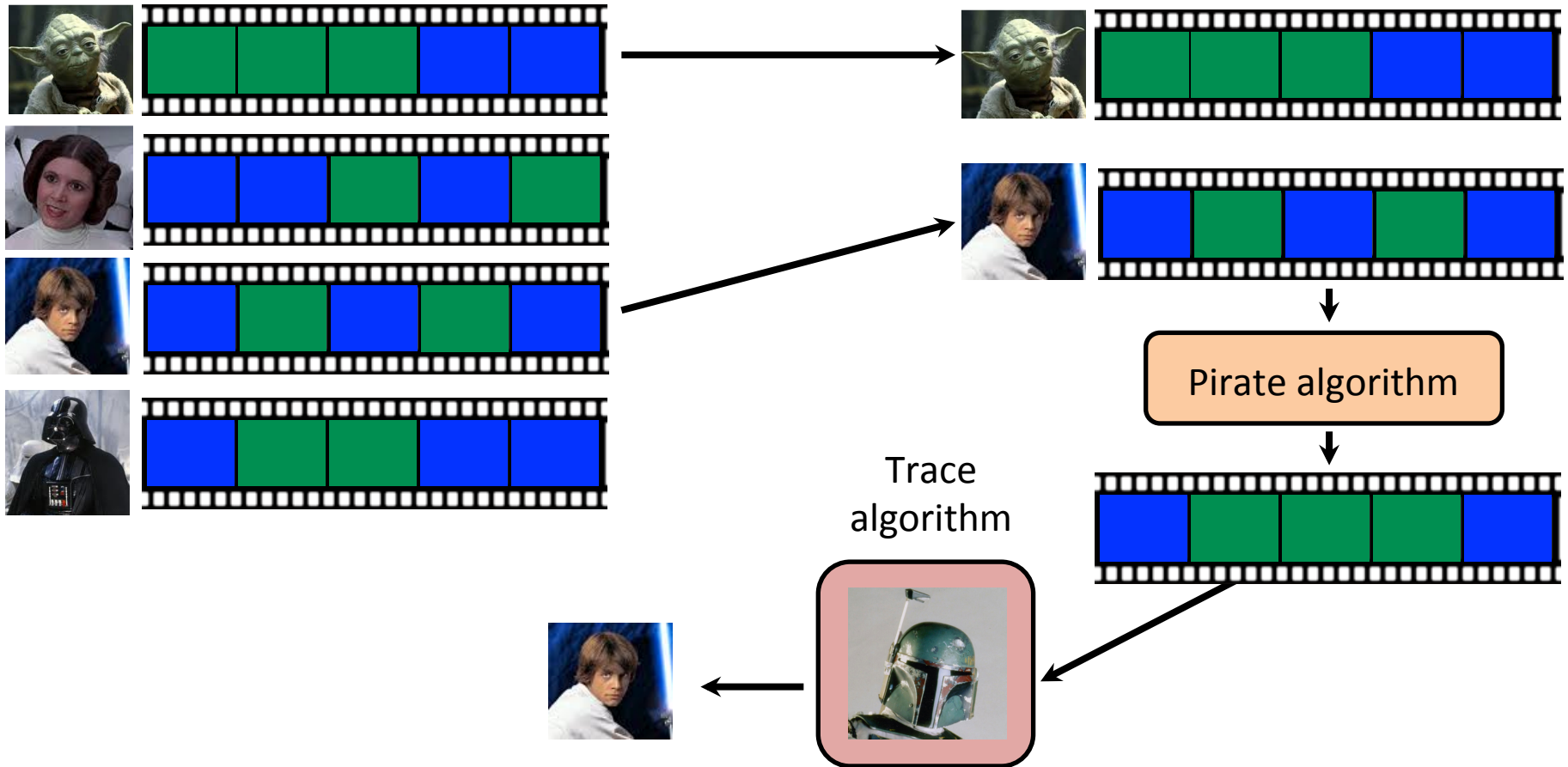


...but the galaxy is full of pirates!



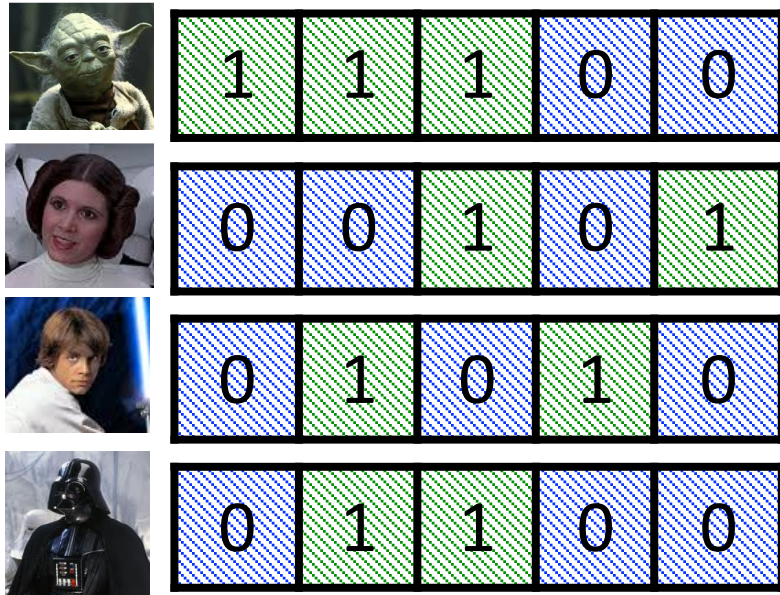
Who collude against me!

# Fingerprinting Codes [BS95]

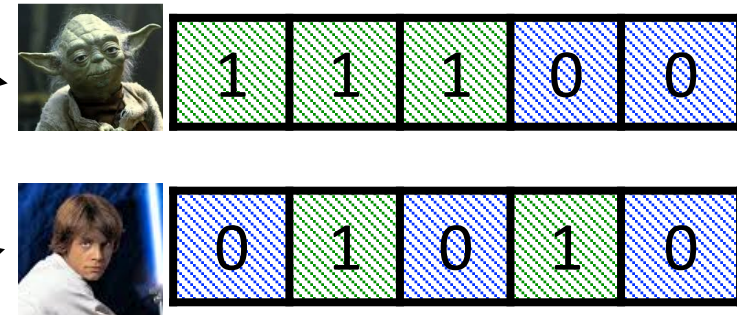


# Fingerprinting Codes [BS95]

Gen( $1^n$ ) outputs  $C \in (\{0,1\}^d)^n$



Pirate coalition  $S \subseteq [n]$

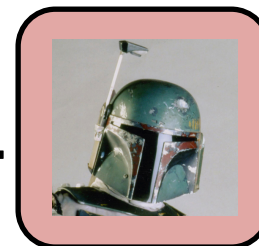


Pirate algorithm



Feasible pirate codeword  $w$

Trace algorithm



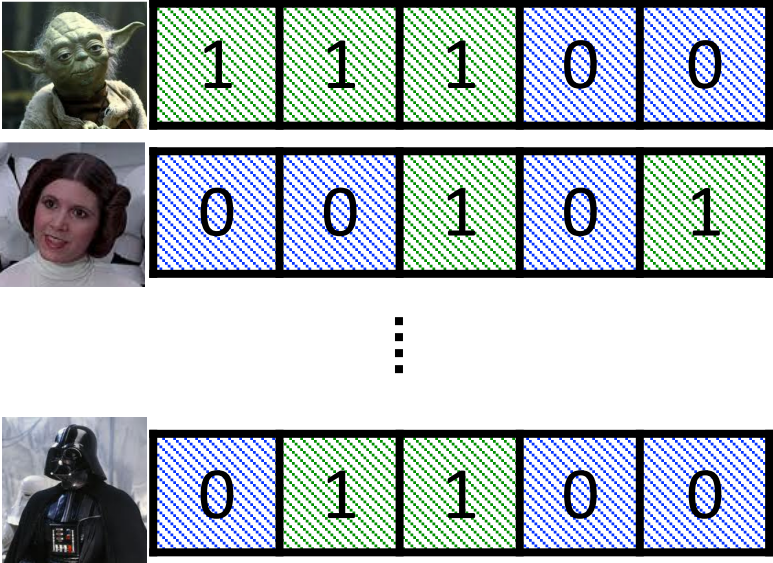
For all coalitions  $S$  and all pirate alg. for producing  $w$ ,

$$\Pr[\text{Trace}(w, C) \in S] \approx 1$$

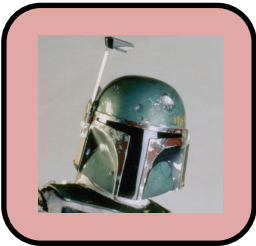


# FP Codes vs. Diff. Privacy

Coalition of  $n$  pirates

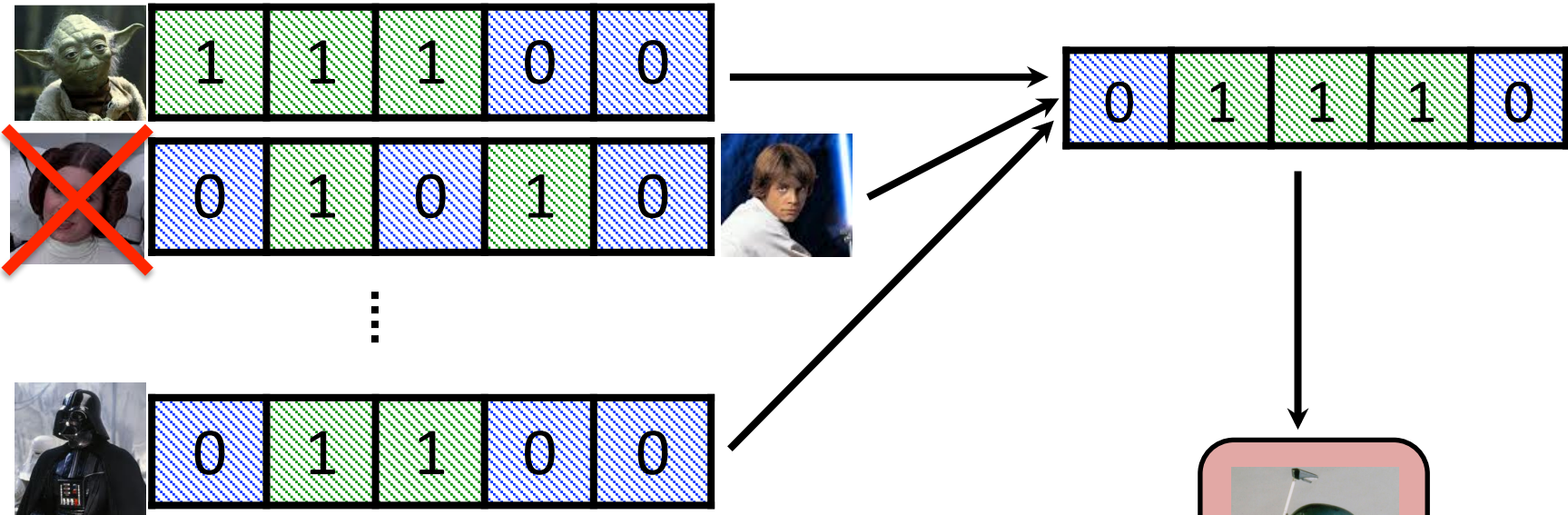


$$\Pr[\text{Trace}(w, C) = \text{Leia}] \geq 1/n$$



# FP Codes vs. Diff. Privacy

Coalition of  $n$  pirates



$$\Pr[\text{Trace}(w, C) = \text{[Image of Leia's face]}] \ll 1/n$$

# FP Codes vs. Diff. Privacy

Trace behaves very differently depending on whether  is in the coalition

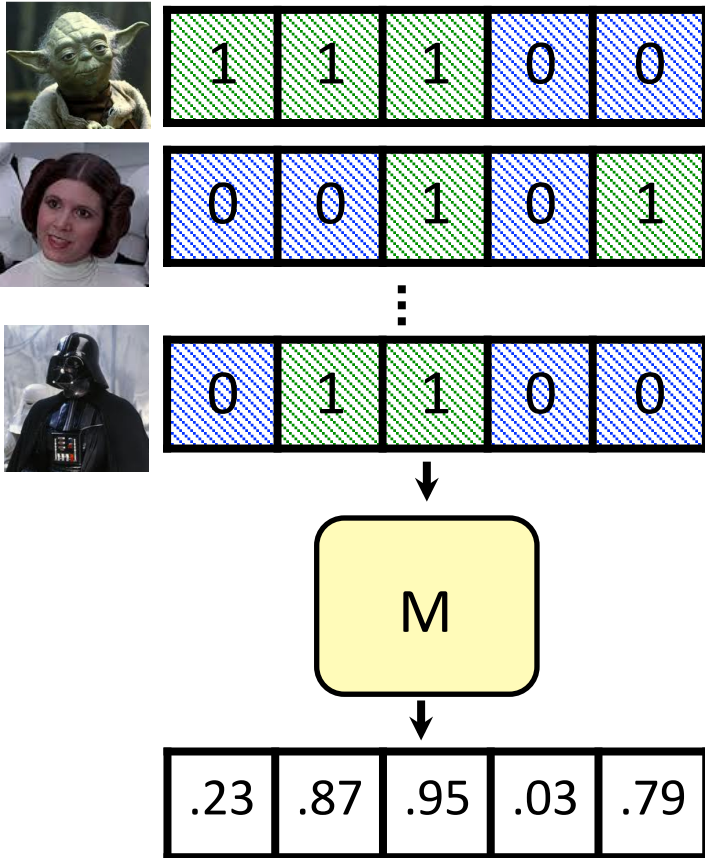


Fingerprinting codes are the “opposite” of differential privacy!

(Parallels computational lower bounds via **traitor-tracing schemes** [DNRRV09, U13])

# Lower Bound for Attribute Means

Database of  $n$  users



Suppose (for contradiction) we have

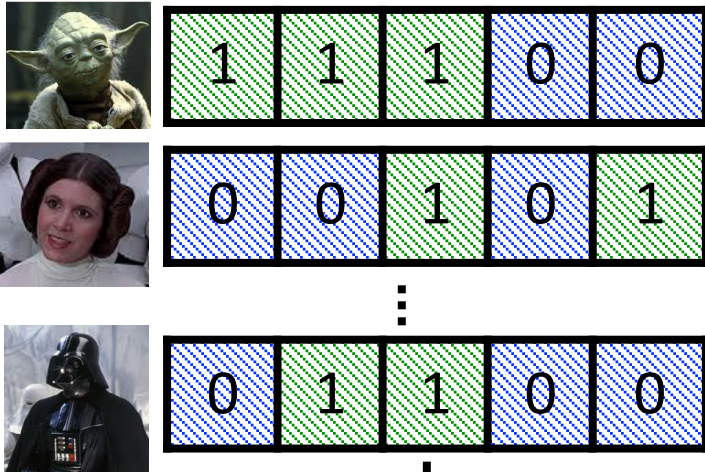
- A FP code of length  $d$  for  $(n+1)$  users
- A diff. private  $M$  that is accurate for attribute means on  $(\{0,1\}^d)^n$

Reduction: Use  $M$  to break security of the FP code



# Lower Bound for Attribute Means

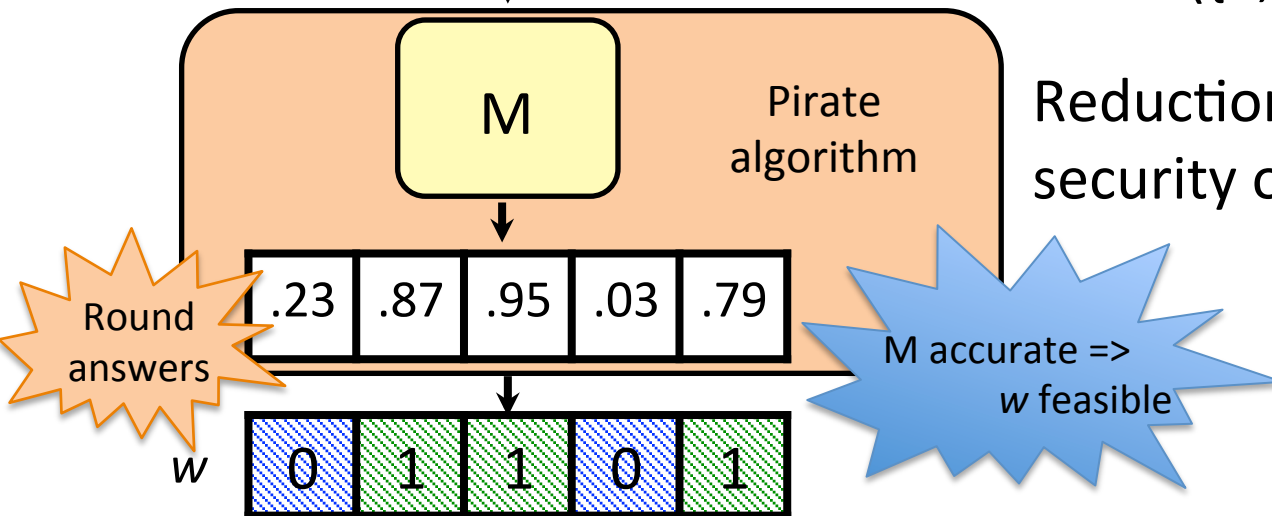
Database of  $n$  users = Coalition of  $n$  pirates



Suppose (for contradiction) we have

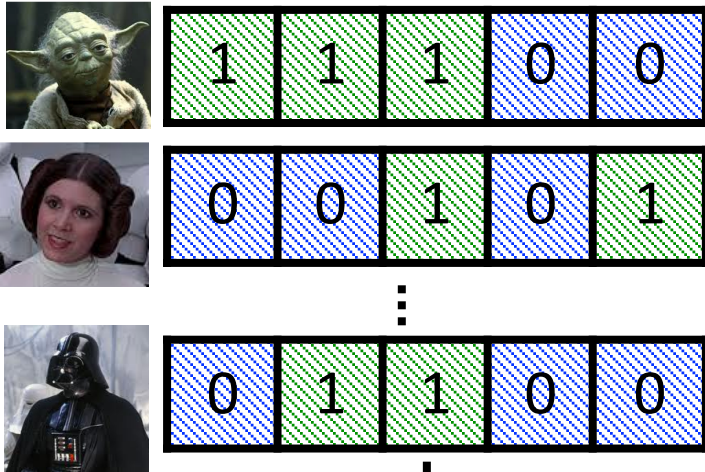
- A FP code of length  $d$  for  $(n+1)$  users
- A diff. private  $M$  that is accurate for attribute means on  $(\{0,1\}^d)^n$

Reduction: Use  $M$  to break security of the FP code

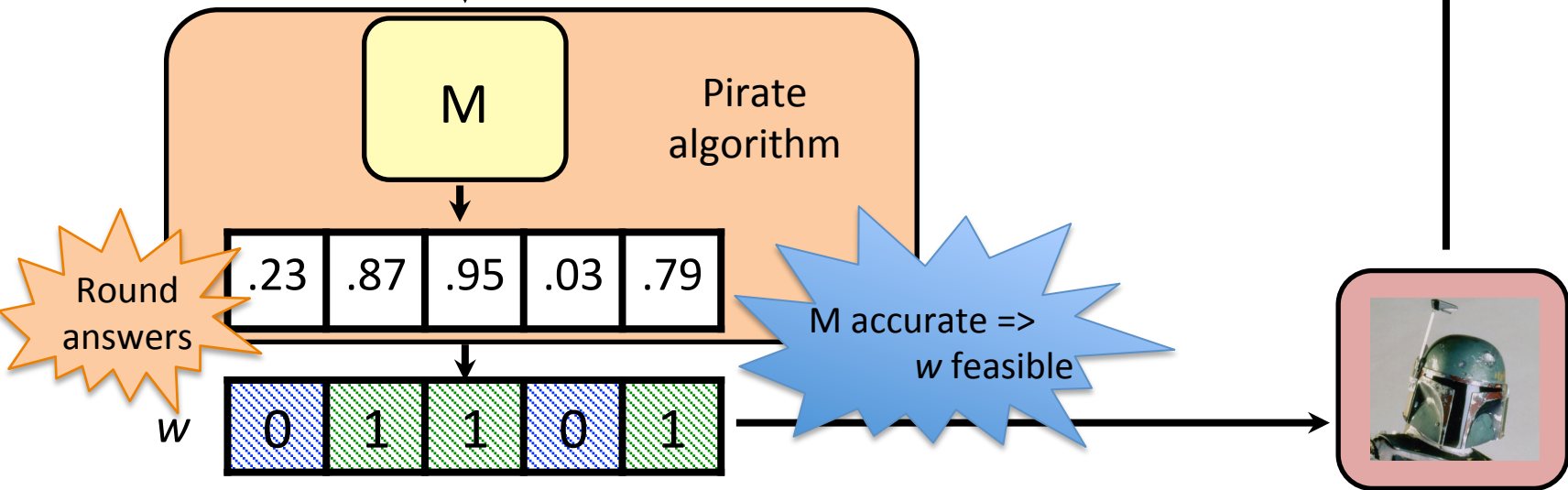


# Lower Bound for Attribute Means

Database of  $n$  users = Coalition of  $n$  pirates

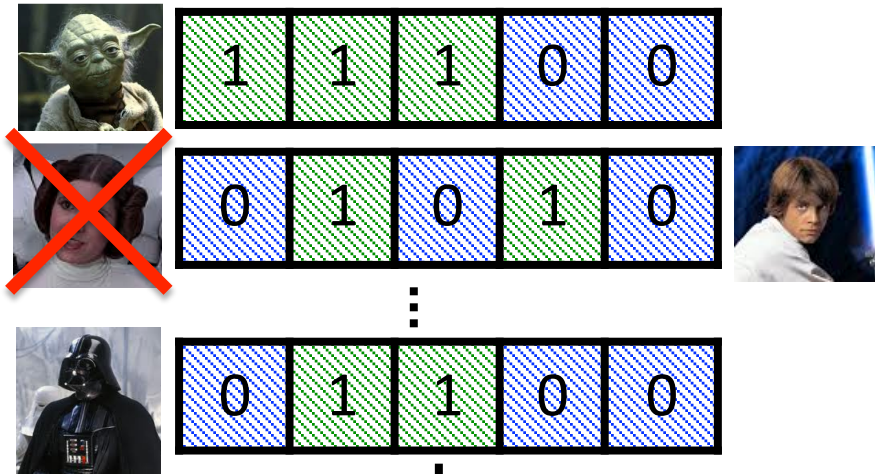


$$\Pr[\text{Trace}(w) = \text{Princess Leia}] \geq 1/n$$



# Lower Bound for Attribute Means

Database of  $n$  users = Coalition of  $n$  pirates

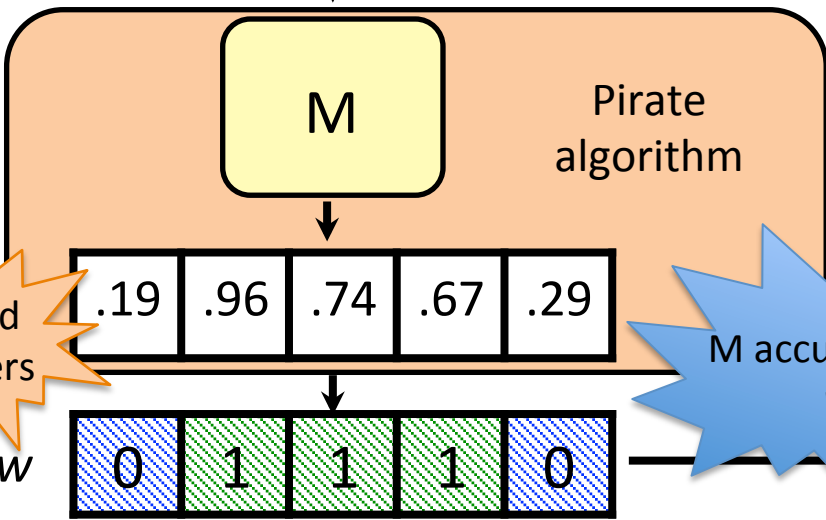


Contradicts security of FP code!

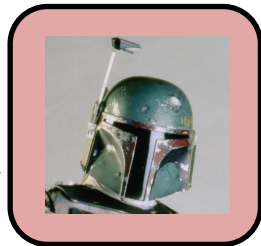
$$\Pr[\text{Trace}(w) = \text{Ava Kaur}] \geq \frac{(1/n) - \delta}{1 + \epsilon}$$

M private => Trace fails

$$\geq \frac{1}{3n}$$



M accurate => w feasible



Round answers  
w

# Lower Bound for Attribute Means

- $\exists$  FP code for  $n$  users with length  $d$   
 $\Rightarrow d$  attribute means require  $n$  samples
- [Tar03]  $\exists$  FP code for  $\tilde{\Omega}(d^{1/2})$  users of length  $d$   
 $\therefore$  attribute means require  $n \geq \tilde{\Omega}(d^{1/2})$

# Sample Complexity for Diff. Privacy

No privacy

$\mathcal{Q}$  = attribute means  
 $\alpha = 0.05$

$\mathcal{Q}, \alpha$  arbitrary

$n = \Theta(\log d)$ [Vap98]	$n = \Theta(\log  \mathcal{Q}  / \alpha^2)$ [Vap98]
---------------------------------	--

$(1, o(1/n))$ -  
 diff. privacy

Upper bound:

$\tilde{O}(d^{1/2})$ [...DMNS06]	$\forall \mathcal{Q}: \tilde{O}(\log  \mathcal{Q}  \cdot d^{1/2} / \alpha^2)$ [HR10]
-------------------------------------	---

Lower bound:

$\tilde{\Omega}(\log d)$ [DN03, Rot10]	$\exists \mathcal{Q}: \max \tilde{\Omega}(\log  \mathcal{Q}  / \alpha), \tilde{\Omega}(1 / \alpha^2)$ [DN03]
---	---

**OUR WORK:**

$\tilde{\Omega}(d^{1/2})$	$\exists \mathcal{Q}: \tilde{\Omega}(\log  \mathcal{Q}  \cdot d^{1/2} / \alpha^2)$
---------------------------	--

◆ Privacy

◆ Accuracy

◆ Sample Complexity

# Conclusions

- Fingerprinting codes yield privacy violations beyond reconstruction attacks
- Price of  $(\epsilon, \delta)$ -diff. privacy for **high-dimensional data**
- Open questions:
  - Sample complexity of **computationally efficient** algorithms for  $k$ -way conjunctions?  
[e.g. BCD+07, GHRU11, UV11, TUV12, DNT13, CTUW14]
  - Combinatorial characterization of sample complexity?  
[e.g. HT10, Har11, NTZ13, BNS13]

Thank you!