



# MATHEMATICAL FRONTIERS

*The National  
Academies of* | SCIENCES  
ENGINEERING  
MEDICINE

[nas.edu/MathFrontiers](http://nas.edu/MathFrontiers)



**Board on  
Mathematical Sciences & Analytics**

# MATHEMATICAL FRONTIERS

## 2019 Monthly Webinar Series, 2-3pm ET

**February 12:** *Machine Learning for Materials Science*

**March 12:** *Mathematics of Privacy*

**April 9:** *Mathematics of Gravitational Waves*

**May 14:** *Algebraic Geometry*

**June 11:** *Mathematics of Transportation*

**July 9:** *Cryptography & Cybersecurity*

**August 13:** *Machine Learning in Medicine*

**September 10:** *Logic and Foundations*

**October 8:** *Mathematics of Quantum Physics*

**November 12:** *Quantum Encryption*

**December 10:** *Machine Learning for Text*

*Made possible by support for BMSA from the  
National Science Foundation  
Division of Mathematical Sciences  
and the  
Department of Energy  
Advanced Scientific Computing Research*

# MATHEMATICAL FRONTIERS

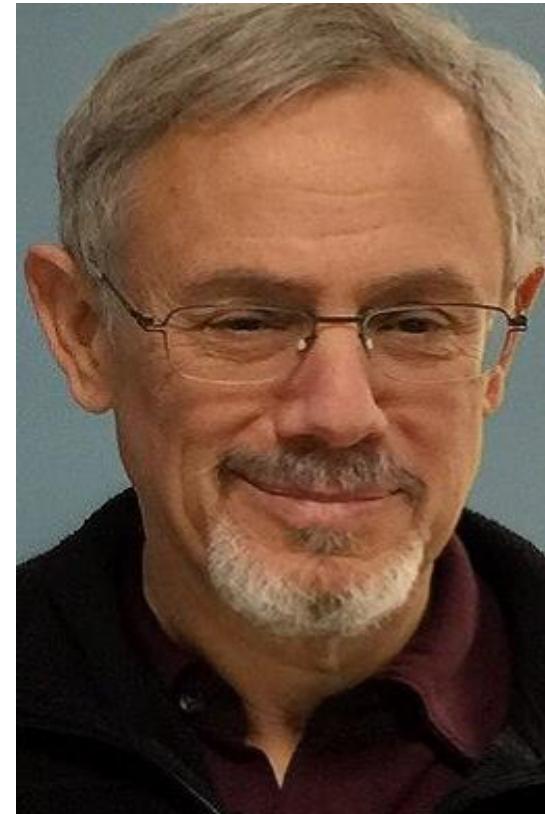
## Mathematics of Privacy



**Kamalika Chaudhuri,  
UC San Diego**



**Katrina Ligett,  
Hebrew University**



**Mark Green,  
UCLA (moderator)**

# MATHEMATICAL FRONTIERS

## Mathematics of Privacy



*Associate Professor of Computer Science  
and Engineering*

### The Mathematics of Differential Privacy

**Kamalika Chaudhuri,  
UC San Diego**

*View webinar videos and learn more about BMSA at [www.nas.edu/MathFrontiers](http://www.nas.edu/MathFrontiers)*

# The Mathematics of Differential Privacy

Kamalika Chaudhuri  
(UCSD)

UC San Diego

# Sensitive Data

Medical Records



Genetic Data



Search Logs



# **AOL Violates Privacy**

# AOL Violates Privacy

## A Face Is Exposed for AOL Searcher No. 4417749

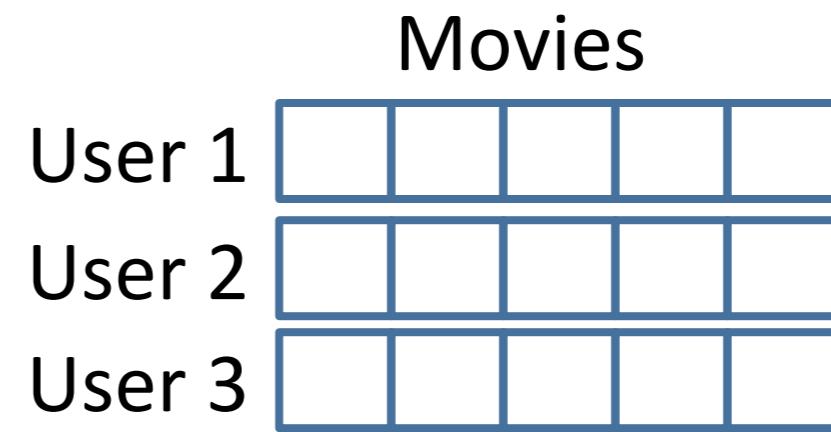
By MICHAEL BARBARO and TOM ZELLER Jr.  
Published: August 9, 2006

Buried in a list of 20 million Web search queries collected by AOL and recently released on the Internet is user No. 4417749. The number was assigned by the company to protect the searcher's anonymity, but it was not much of a shield.



No. 4417749 conducted hundreds of searches over a three-month period on topics ranging from "numb fingers" to "60 single men" to "dog that urinates on

# Netflix Violates Privacy [NS08]



2-8 movie-ratings and dates for Alice reveals:  
Whether Alice is in the dataset or not  
Alice's other movie ratings

# High-dimensional Data is Unique

## Example: UCSD Employee Salary Table

Position	Gender	Department	Ethnicity	Salary
Faculty	Female	CSE	SE Asian	-

One employee (Kamalika) fits description!

**Simply anonymizing data is unsafe!**

# Disease Association Studies [WLWTZ09]



**Cancer**

1.00
.190 1.00
.216 .251 1.00
.186 .117 .047 1.00
.154 .011 .170 .083 1.00
.190 .140 .102 .095 .139 1.00
.270 .215 .294 .248 .140 .141 1.00
.101 .085 .170 .056 .234 .099 .175 1.00
.239 .071 .163 .111 .161 .093 .199 .157 1.00
.471 .117 .243 .094 .144 .123 .283 .216 .274 1.00
.179 .202 .132 .094 .087 .159 .207 .108 .092 .294 1.00

**Healthy**

1.00
.141 1.00
.099 .175 1.00
.093 .199 .157 1.00
.123 .283 .216 .274 1.00
.159 .207 .108 .092 .294 1.00
.088 .152 .075 .163 .156 .220 1.00
.046 .161 .092 .072 .157 .143 .147 1.00
.078 .392 .122 .229 .160 .172 .145 .177 1.00
.045 .155 .135 .139 .110 .048 .126 .104 .169 1.00
.178 .135 .102 .258 .314 .165 .147 .158 .131 .074 1.00

Correlations

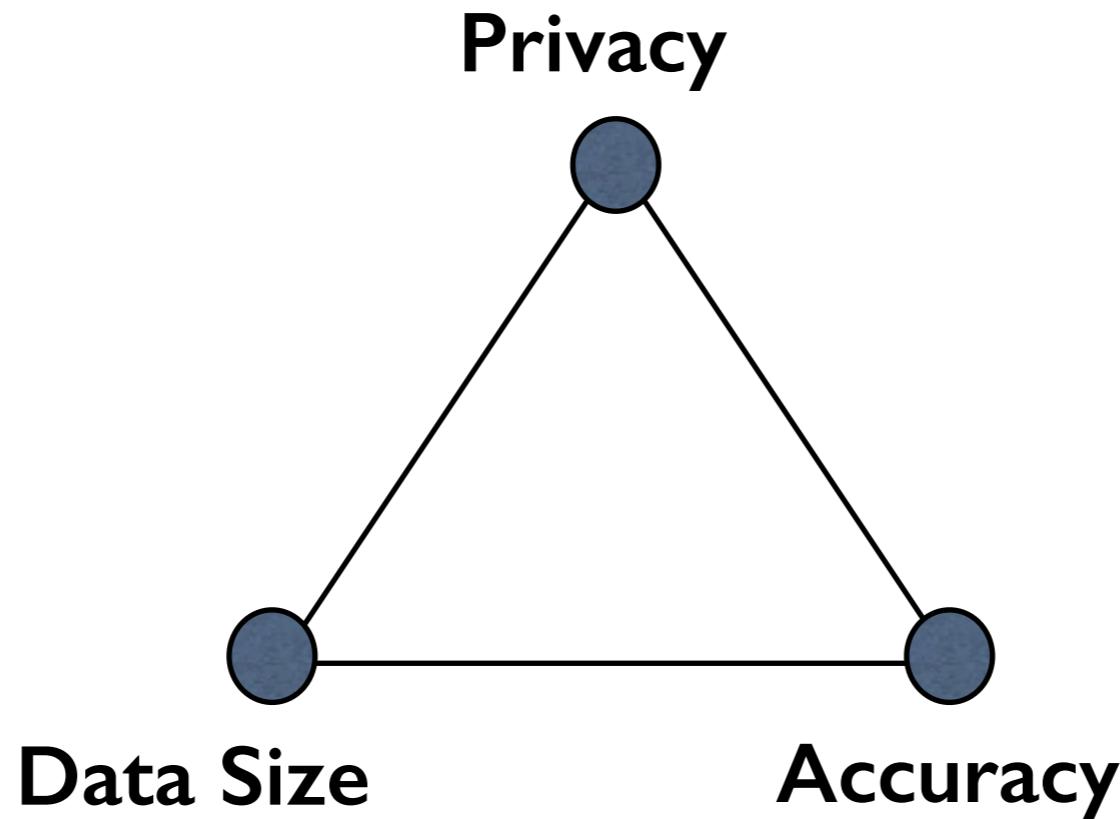
Correlations

Correlation ( $R^2$  values), Alice's DNA reveals:

If Alice is in the **Cancer** set or **Healthy** set

**Simply anonymizing data is unsafe!**

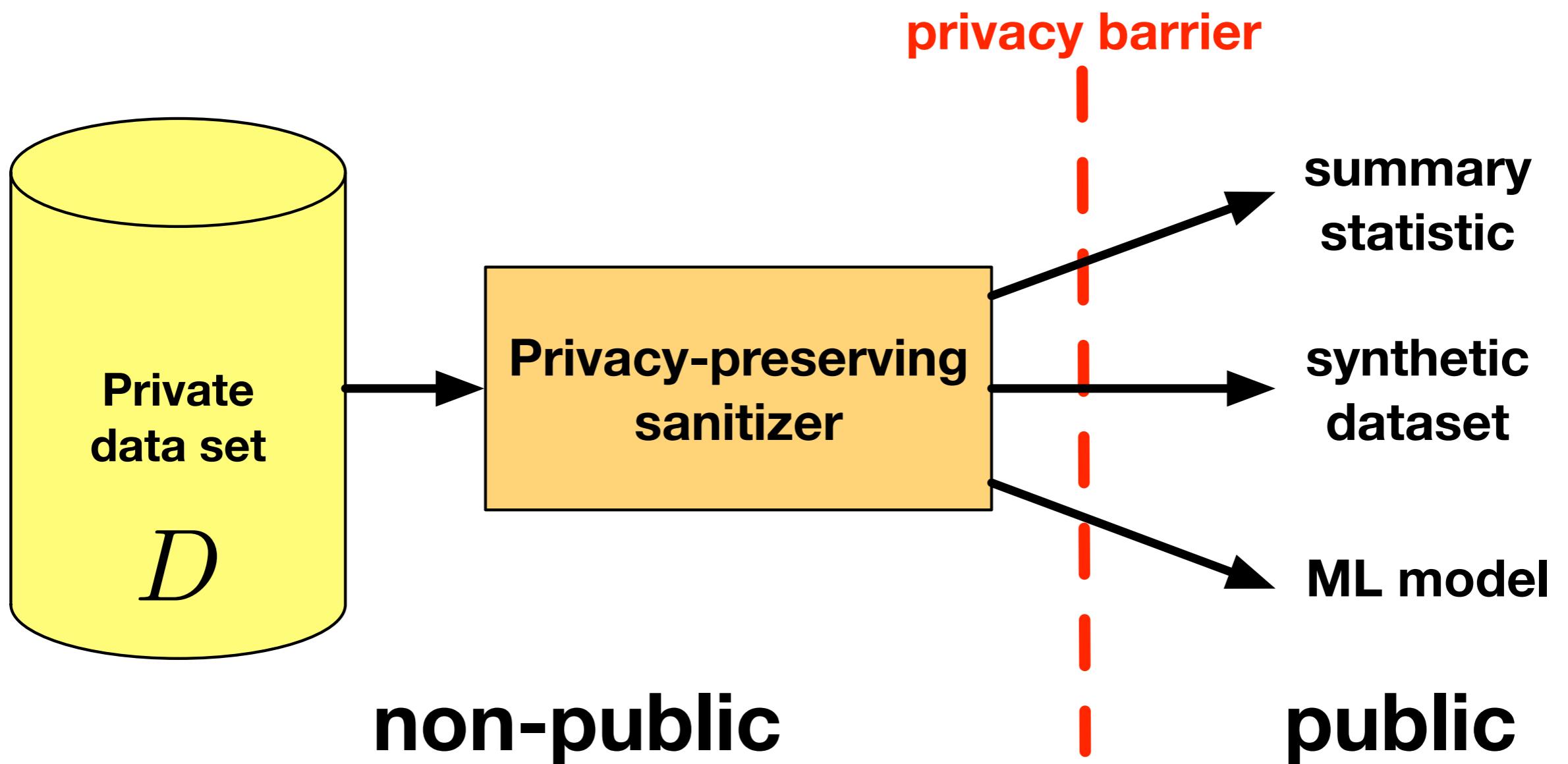
**Statistics on small data sets is unsafe!**



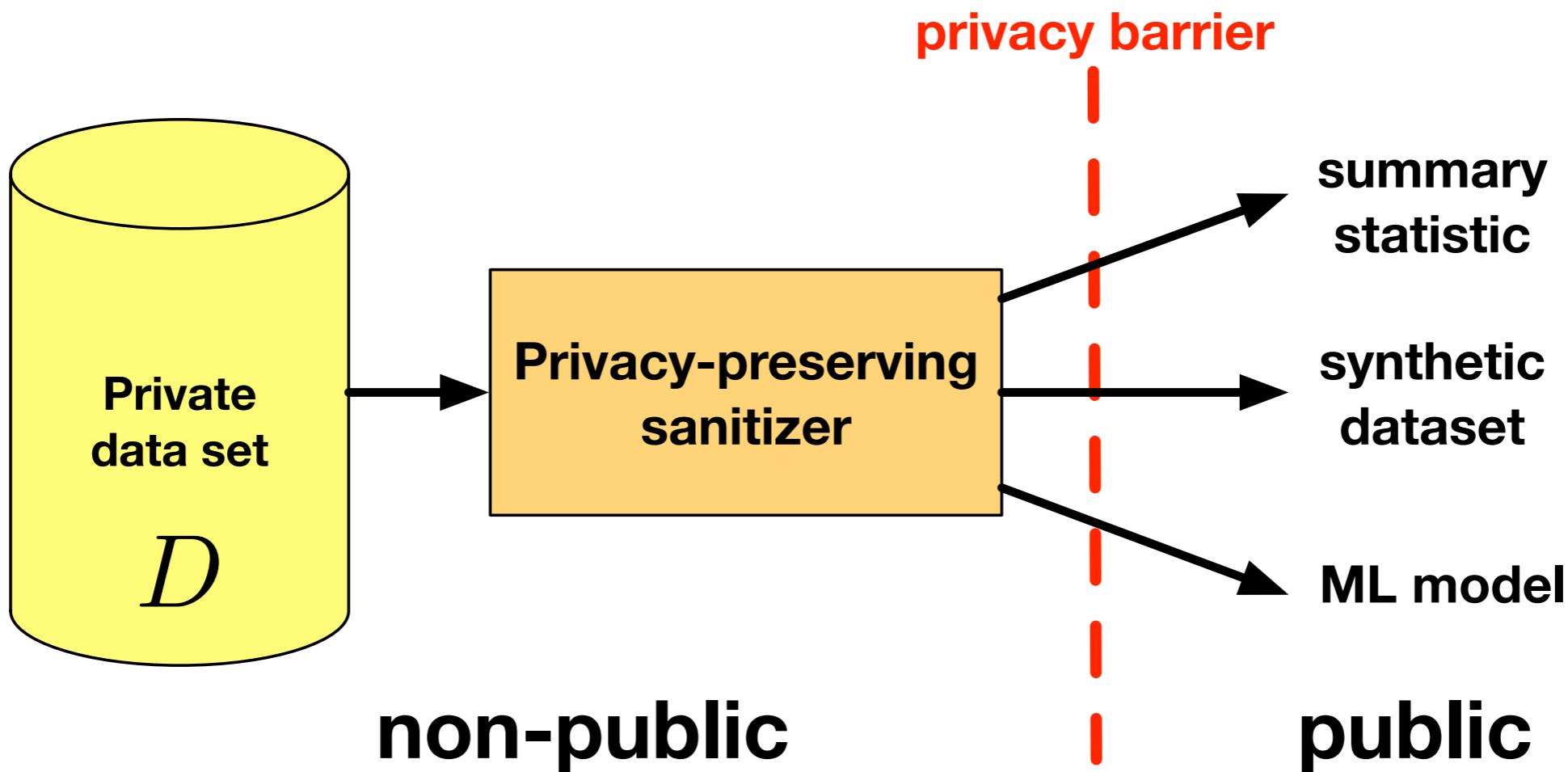
**Need: rigorous definition of privacy**

# Privacy Definition

# The Setting



# Property of Sanitizer



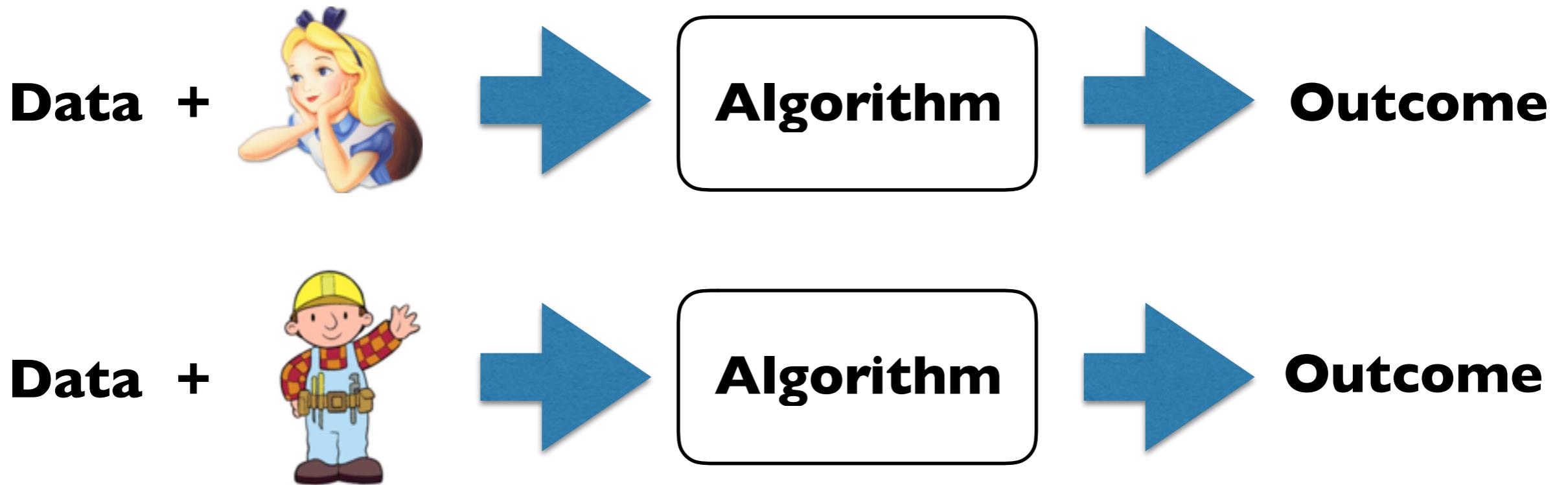
Aggregate information computable

Individual information protected  
(robust to side-information)

# Differential Privacy

[Dwork-McSherry-Nissim-Smith 2006]

# Differential Privacy [DMNS06]



Participation of a person does not change outcome

Since a person has agency, they can decide to participate in a dataset or not

## Adversary



**Prior Knowledge:**

A's Genetic profile

A smokes

# Adversary



Prior Knowledge:  
A's Genetic profile  
A smokes

## Case I: Study

1.00										
.190	1.00									
.216	.251	1.00								
.186	.117	.047	1.00							
.154	.011	.170	.083	1.00						
.190	.140	.102	.095	.139	1.00					
.270	.215	.294	.248	.140	.141	1.00				
.101	.085	.170	.056	.234	.099	.175	1.00			
.239	.071	.163	.111	.161	.093	.199	.157	1.00		
.471	.117	.243	.094	.144	.123	.283	.216	.274	1.00	
.179	.202	.132	.094	.087	.159	.207	.108	.092	.294	1.00

# Cancer

[ Study violates A's privacy ]

A has  
cancer

# Adversary



Prior Knowledge:  
A's Genetic profile  
A smokes

## Case I: Study

1.00										
.190	1.00									
.216	.251	1.00								
.186	.117	.047	1.00							
.154	.011	.170	.083	1.00						
.190	.140	.102	.095	.139	1.00					
.270	.215	.294	.248	.140	.141	1.00				
.101	.085	.170	.056	.234	.099	.175	1.00			
.239	.071	.163	.111	.161	.093	.199	.157	1.00		
.471	.117	.243	.094	.144	.123	.283	.216	.274	1.00	
.179	.202	.132	.094	.087	.159	.207	.108	.092	.294	1.00

# Cancer

## [ Study violates A's privacy ]

A has  
cancer

## Case 2: Study

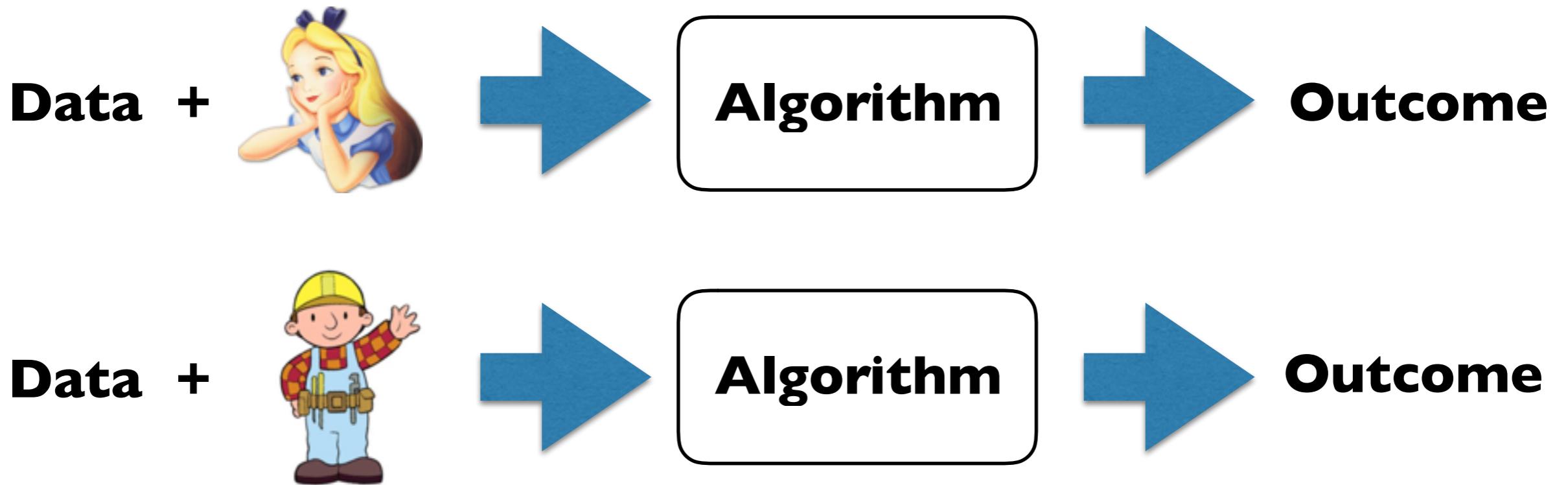


# Smoking causes cancer

A probably  
has cancer

[ Study does not violate privacy]

# Differential Privacy [DMNS06]



Participation of a person does not change outcome

Since a person has agency, they can decide to participate in a dataset or not

# How to ensure this?

...through randomness

$A(\mathbf{Data} + \mathbf{ }$



Random  
variables

have close  
distributions

$A(\mathbf{Data} + \mathbf{ }$



# How to ensure this?

Random variables

$A(\mathbf{Data} + )$



have close distributions

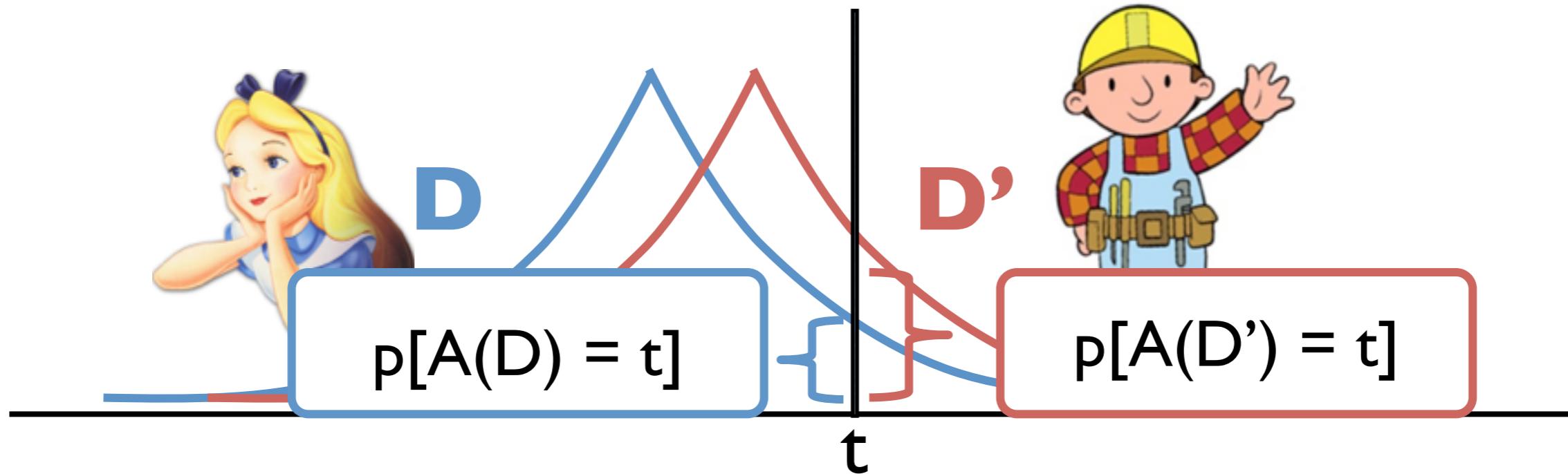
$A(\mathbf{Data} + )$



**Randomness:** Added by randomized algorithm A

**Closeness:** Probability of every event is close

# Differential Privacy [DMNS06]



For all  $D, D'$  that differ in one person's value,

If  $A = \epsilon$ -differentially private randomized algorithm, then for all  $t$ ,

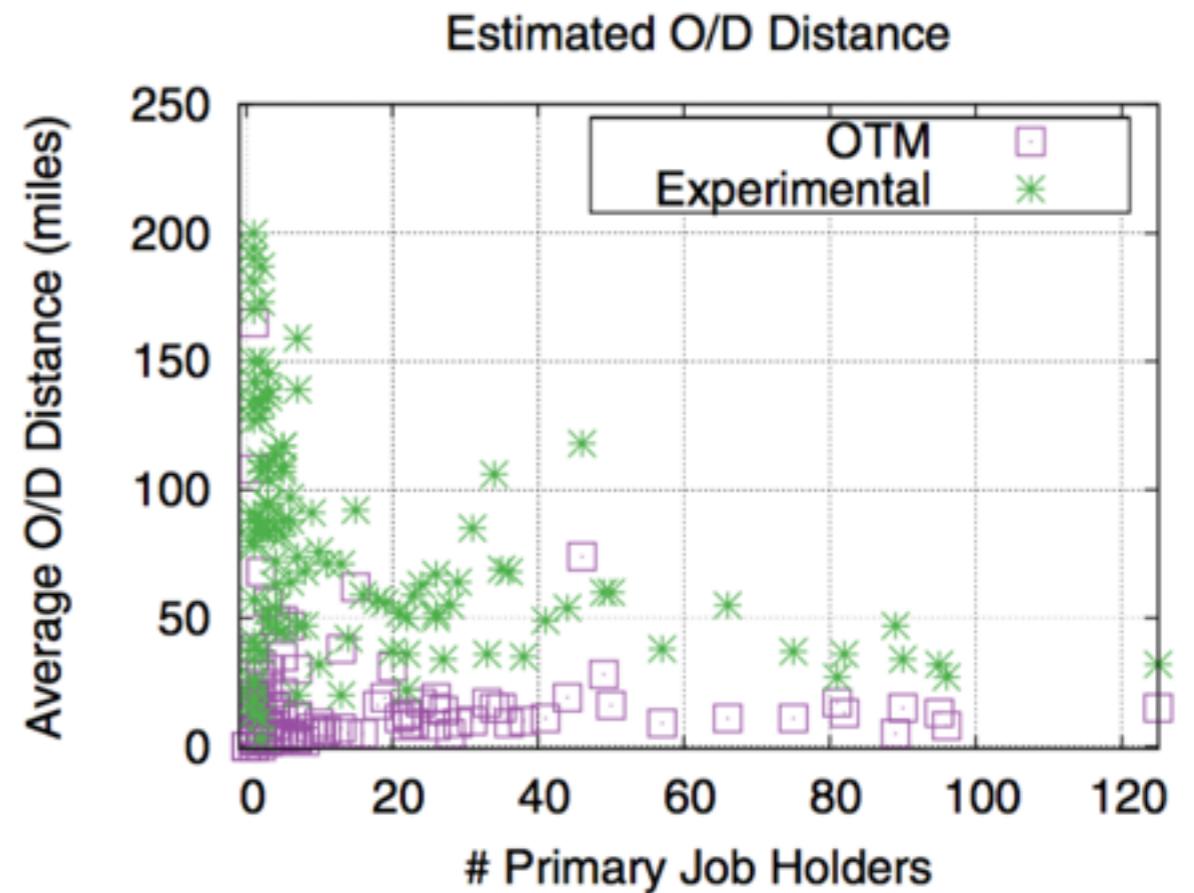
$$p(A(D) = t) \leq e^\epsilon p(A(D') = t)$$

$\epsilon$  = privacy parameter

What can we do with  
Differential Privacy?

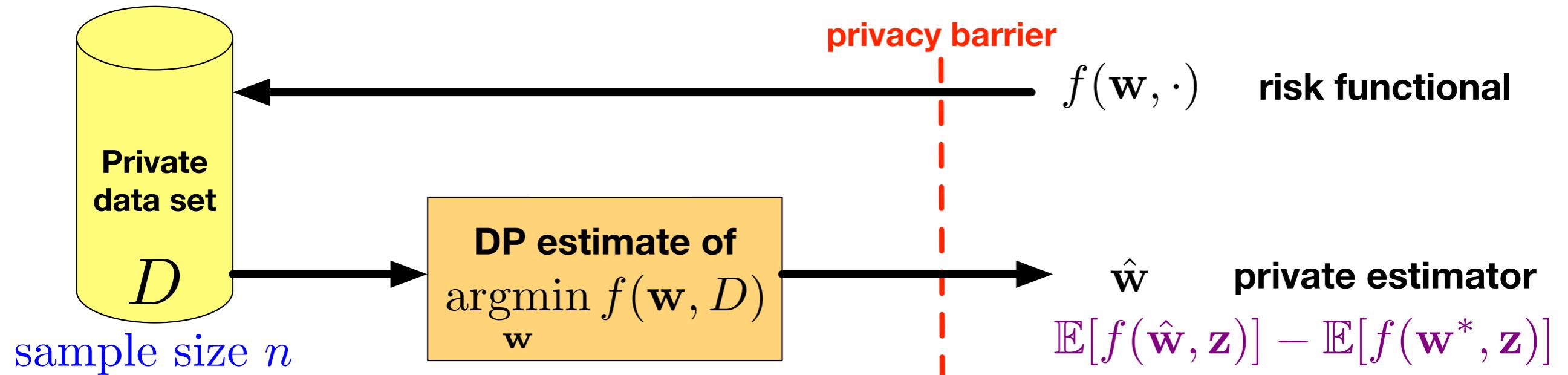
# Statistics

- Counts, means, variances
- Contingency tables, histograms
- 2020 Census will use differential privacy



[MKAGV'08]

# Estimation and prediction problems



**Statistical estimation:** estimate a parameter or predictor using private data that has good expected performance on future data.

# Examples: Estimation and Prediction



# Language Models

[MRTZ'18]

# HIV Epidemiology

# [SMVLC'18]

# References

- [DMNS06] Calibrating Noise to Sensitivity in Private Data Analysis, C. Dwork, F. McSherry, K. Nissim, A. Smith, TCC 2006
- [NS08] Robust Deanonymization of Large Sparse Datasets, A. Narayanan, V. Shmatikov, IEEE-S&P, 2008
- [WLWTZ09] Learning your identity and disease from research papers: Information leaks in GWAS, Wang et al, CCS 2009
- [MKAGV08] Privacy: Theory meets practice on the map, Machanavajjhala et al, ICDE 2008
- [SMVLC18] Differentially Private Continual Release of Graph Statistics, Song et al, Arxiv 2018
- [MRTZ18] Learning Differentially Private Recurrent Language Models, McMahon et al, Arxiv 2018

# MATHEMATICAL FRONTIERS

## Mathematics of Privacy



**Katrina Ligett,**  
**Hebrew University**

*Associate Professor of Computer Science*

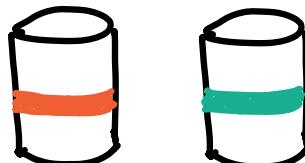
### What Does Differential Privacy Offer?

What does  
differential privacy  
offer?

Katrina Ligett  
Hebrew University of Jerusalem

## $\epsilon$ -Differential Privacy [DMNS 06]

For all



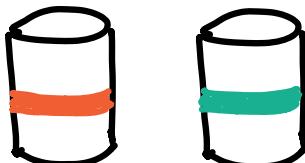
,  $t$  :

$$D \sim D'$$

$$\sup_t \left| \log \frac{P(A(D) = t)}{P(A(D') = t)} \right| \leq \epsilon$$

## $\epsilon$ -Differential Privacy [DMNS 06]

For all



$$D \sim D'$$

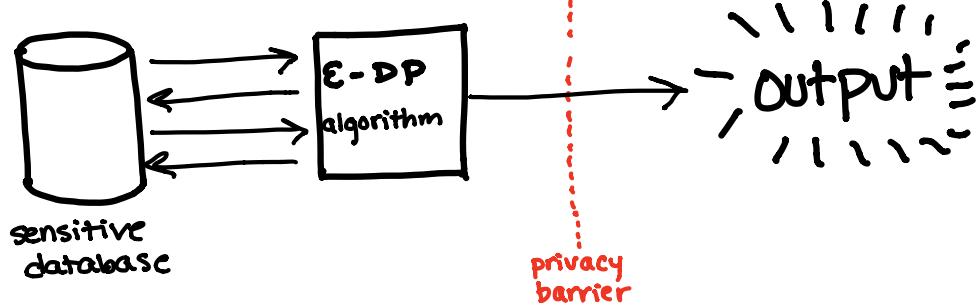
,  $t$  :

$$\sup_t \left| \log \frac{P(A(D) = t)}{P(A(D') = t)} \right| \leq \epsilon$$

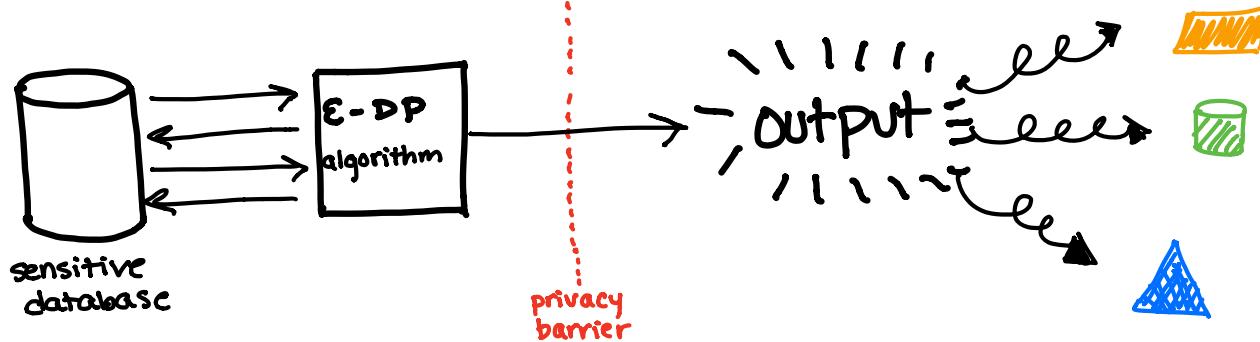
property of mechanism



# Properties of $\epsilon$ -Differential Privacy

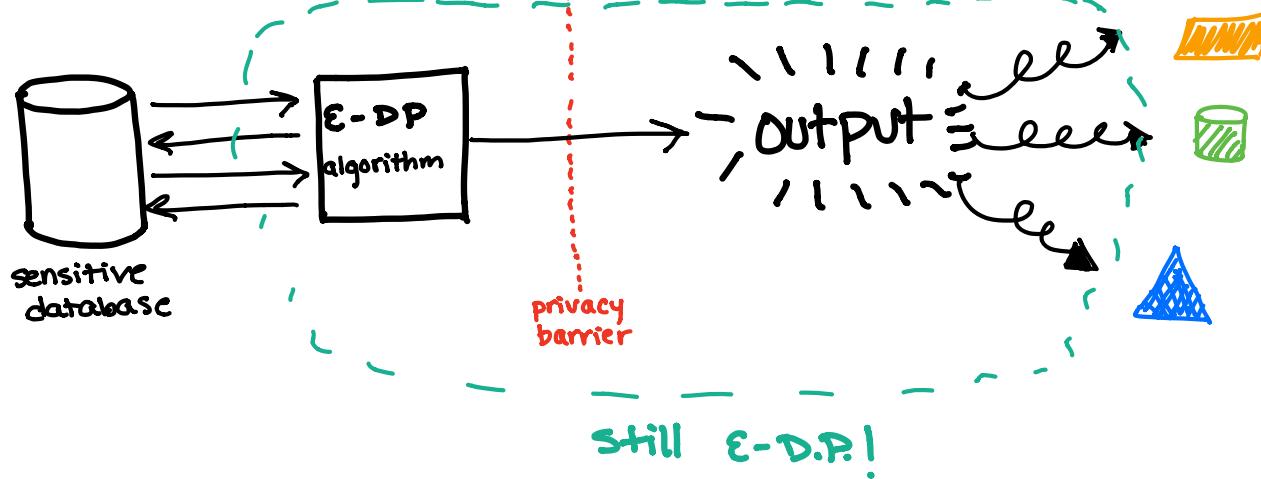


# Properties of $\epsilon$ -Differential Privacy



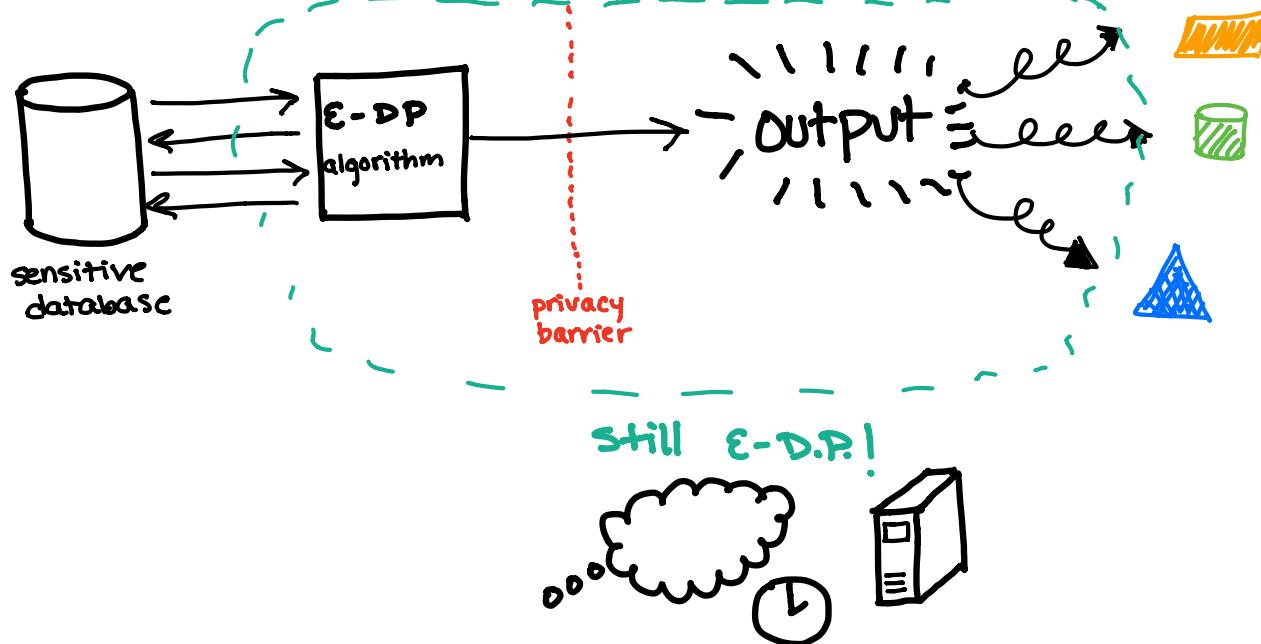
1. post-processing

# Properties of $\epsilon$ -Differential Privacy



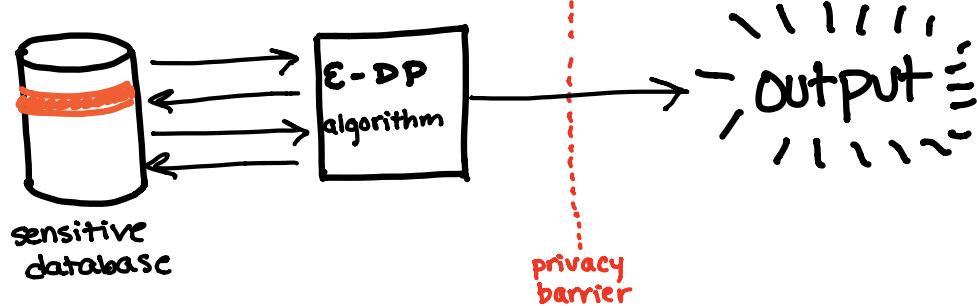
1. post-processing

# Properties of $\epsilon$ -Differential Privacy



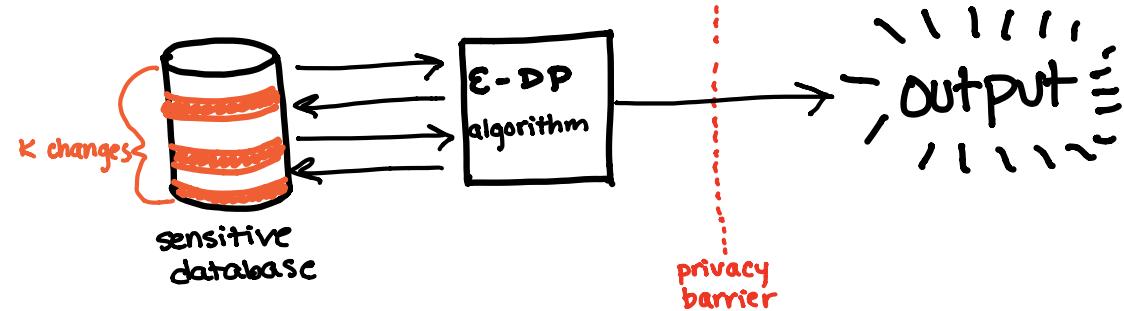
1. post-processing

# Properties of $\epsilon$ -Differential Privacy



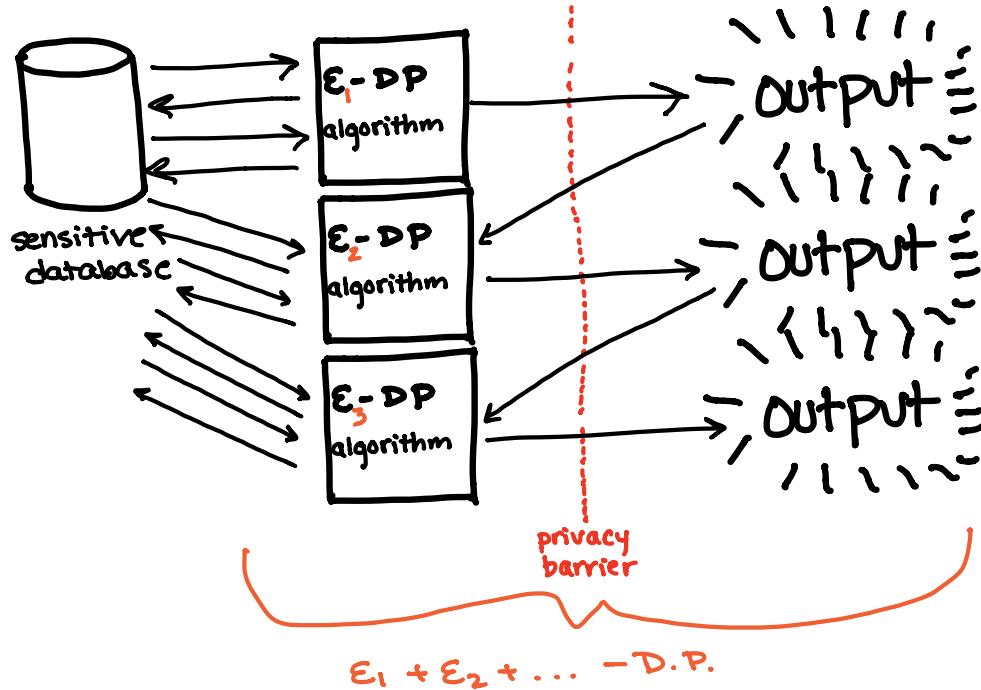
2. group privacy

# Properties of $\epsilon$ -Differential Privacy



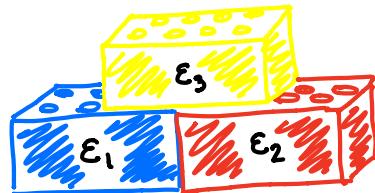
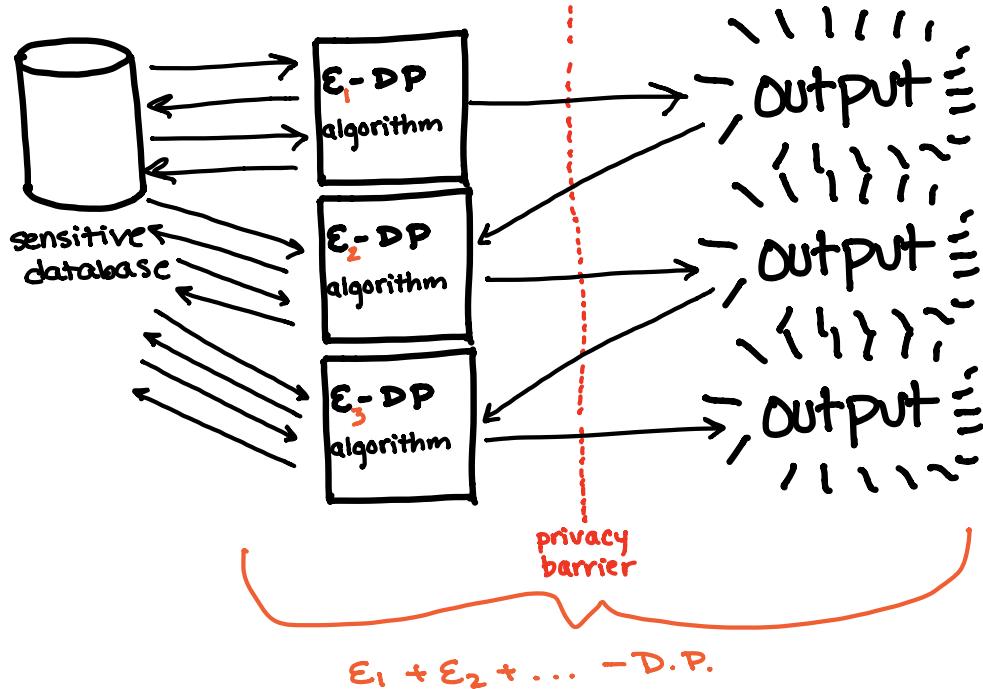
2.  $K\epsilon$  - group privacy

# Properties of $\epsilon$ -Differential Privacy



3. composition

# Properties of $\epsilon$ -Differential Privacy



3. composition

# What does D.P. protect against?

- attacker knows all of database except you

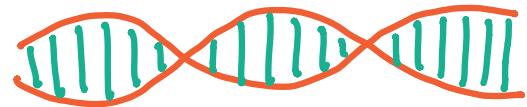


# What does D.P. protect against?

- attacker knows all of database except you



- your data is correlated with a few other people



# What does D.P. protect against?

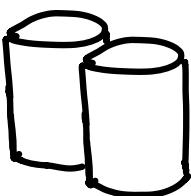
- attacker knows all of database except you



- your data is correlated with a few other people



- attacker gains complementary database



Ori Heffetz and Katrina Ligett  
"Privacy and Data-Based Research"  
Journal of Economic Perspectives  
2014

# MATHEMATICAL FRONTIERS

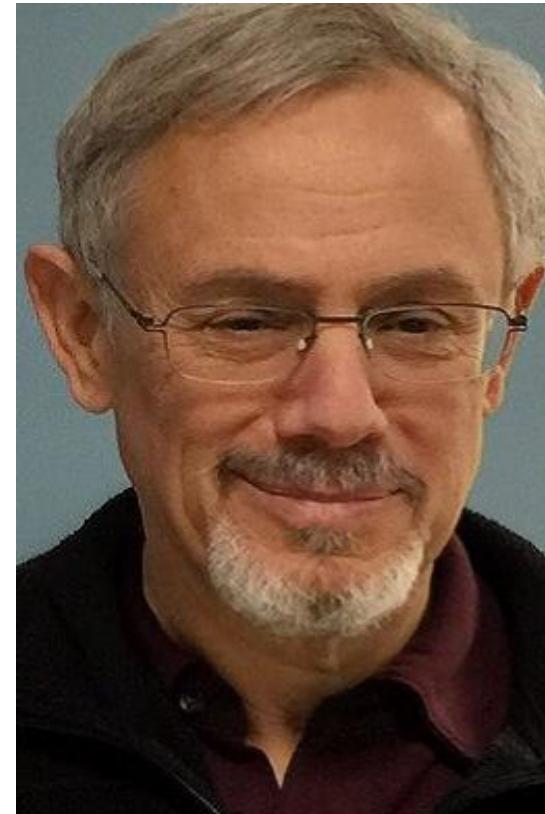
## Mathematics of Privacy



**Kamalika Chaudhuri,  
UC San Diego**



**Katrina Ligett,  
Hebrew University**



**Mark Green,  
UCLA (moderator)**

# MATHEMATICAL FRONTIERS

## 2019 Monthly Webinar Series, 2-3pm ET

**February 12:** *Machine Learning for Materials Science*

**March 12:** *Mathematics of Privacy*

**April 9:** *Mathematics of Gravitational Waves*

**May 14:** *Algebraic Geometry*

**June 11:** *Mathematics of Transportation*

**July 9:** *Cryptography & Cybersecurity*

**August 13:** *Machine Learning in Medicine*

**September 10:** *Logic and Foundations*

**October 8:** *Mathematics of Quantum Physics*

**November 12:** *Quantum Encryption*

**December 10:** *Machine Learning for Text*

*Made possible by support for BMSA from the  
National Science Foundation  
Division of Mathematical Sciences  
and the  
Department of Energy  
Advanced Scientific Computing Research*

