# 2550 Intro to cybersecurity

L15: Data Privacy
(Anonymous Data Isn't!)

abhi shelat/Ran Cohen

# The era of big data

# Predict our preferences

# Social networks

# Medical & Genomic data

# Contact tracing

# Statistical data

# Big Data is Invaluable

Schizophrenia Genome-Wide Association Studies

3,500 cases
⇔ 0 loci

10,000 cases
⇔ 5 loci

35,000 cases
⇔ 62 loci!

Data courtesy of Manolis Kellis

Increasing sample sizes for schizophrenia association studies has led to increases in the number of risk genes discovered

new biological insights

# Outline

- Popular ideas that do not work
  + privacy horror stories

- An approach that works

# Popular idea #1

Remove Personally Identifiable Information (PII)

we do not collect any personal information

🔍 All     🖾 News     ▶ Videos     🖼 Images     ⚲ Maps     ⋮ More

About 2,060,000,000 results (0.61 seconds)

# Anonymizing data

**NIST**
National Institute of
Standards and Technology
U.S. Department of Commerce

Special Publication 800-122

## Guide to Protecting the Confidentiality of Personally Identifiable Information (PII)

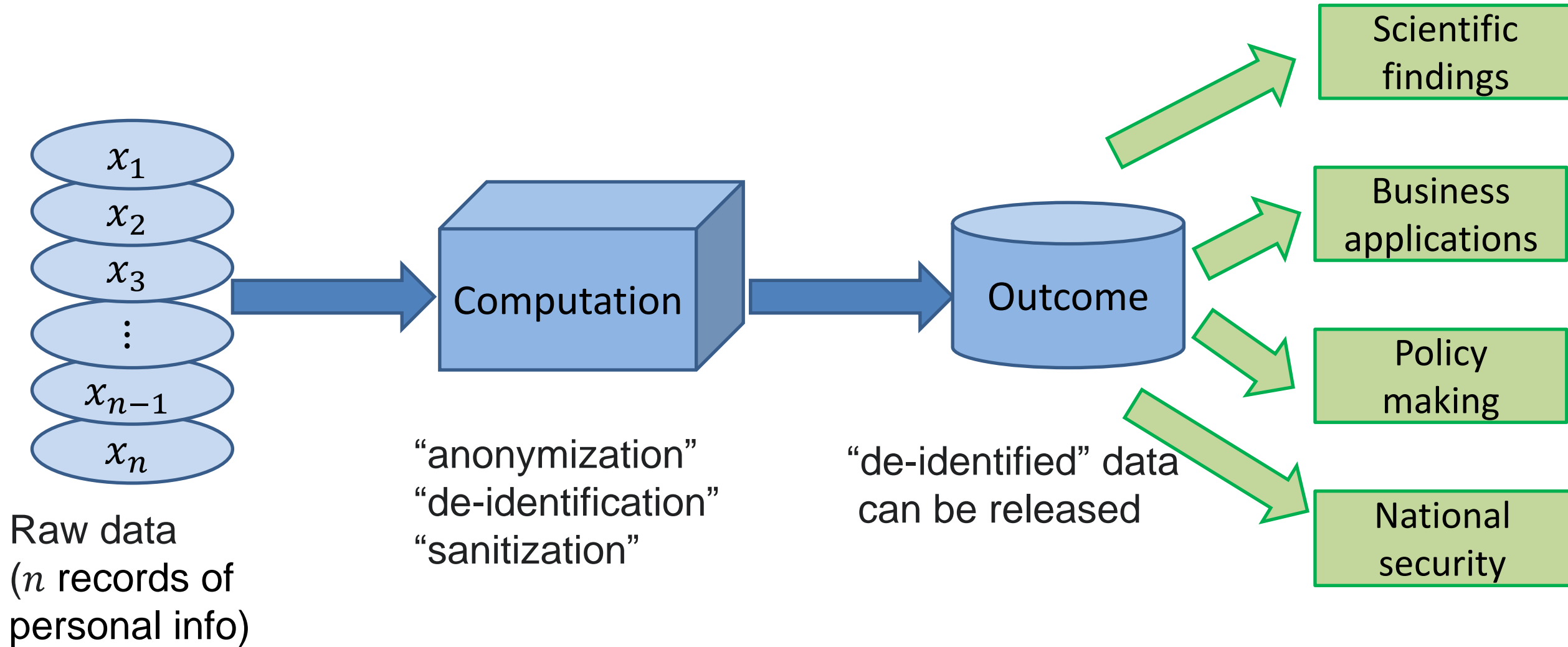Recommendations of the National Institute of Standards and Technology

Erika McCallister
Tim Grance
Karen Scarfone

# "Privacy-preserving" data release



Raw data
($n$ records of
personal info)

$x_1$
$x_2$
$x_3$
⋮
$x_{n-1}$
$x_n$

Computation

"anonymization"
"de-identification"
"sanitization"

Outcome

"de-identified" data
can be released

Scientific findings

Business applications

Policy making

National security

# Linkage attack (Sweeney '97)

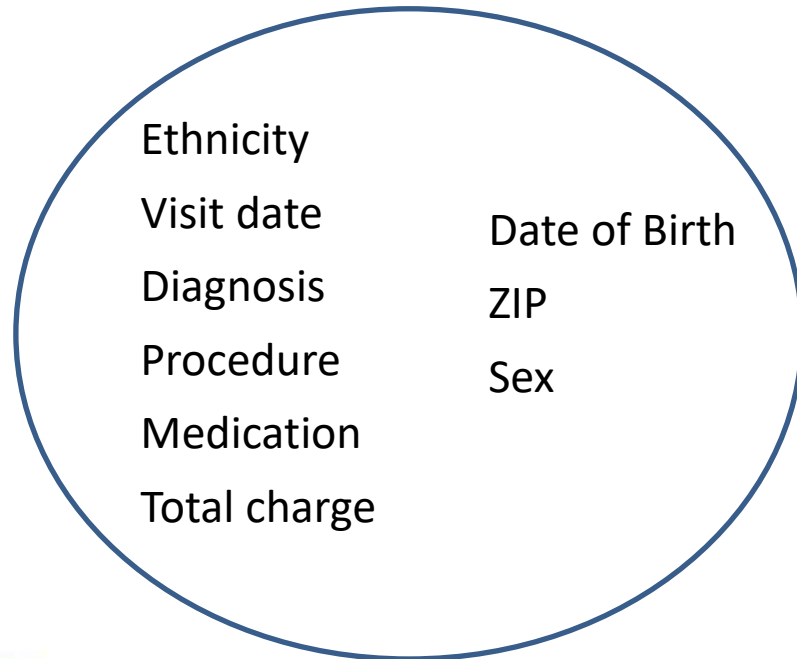Massachusetts Group Insurance Commission (GIC)

- In mid-1990s GIC released "anonymized" data of state employees that showed every single hospital visit

- Goal: provide real data for researchers

- Privacy?
  Removed personally identifiable information (PII): Name, SSN, Address

- William Weld, then Governor of Massachusetts, assured the public that GIC had protected patient privacy by deleting identifiers

# Linkage attack (Sweeney '97)

## MA Group Insurance Commission

- Contained ~135,000 patients
- Anonymized: Name, SSN removed

Ethnicity

Visit date

Diagnosis          Date of Birth

Procedure          ZIP

Medication         Sex

Total charge

## Voters registration of Cambridge MA

- Public information

**Auxiliary information**

Date of Birth     Name

ZIP               Address

Sex               Date registered

                  Party affiliation

                  Date last voted

Commonwealth of Massachusetts
Group Insurance Commission

Register to VOTE

14

# Linkage attack (Sweeney '97)

- A unique record fully de-anonymize the record
- (DoB, ZIP, Sex) uniquely identifies 87% of US population

Ethnicity

Visit date

Diagnosis

Procedure

Medication

Total charge

Date of Birth

ZIP

Sex

Name

Address

Date registered

Party affiliation

Date last voted

Quasi-identifiers

Commonwealth of Massachusetts
Group Insurance Commission

Register to VOTE
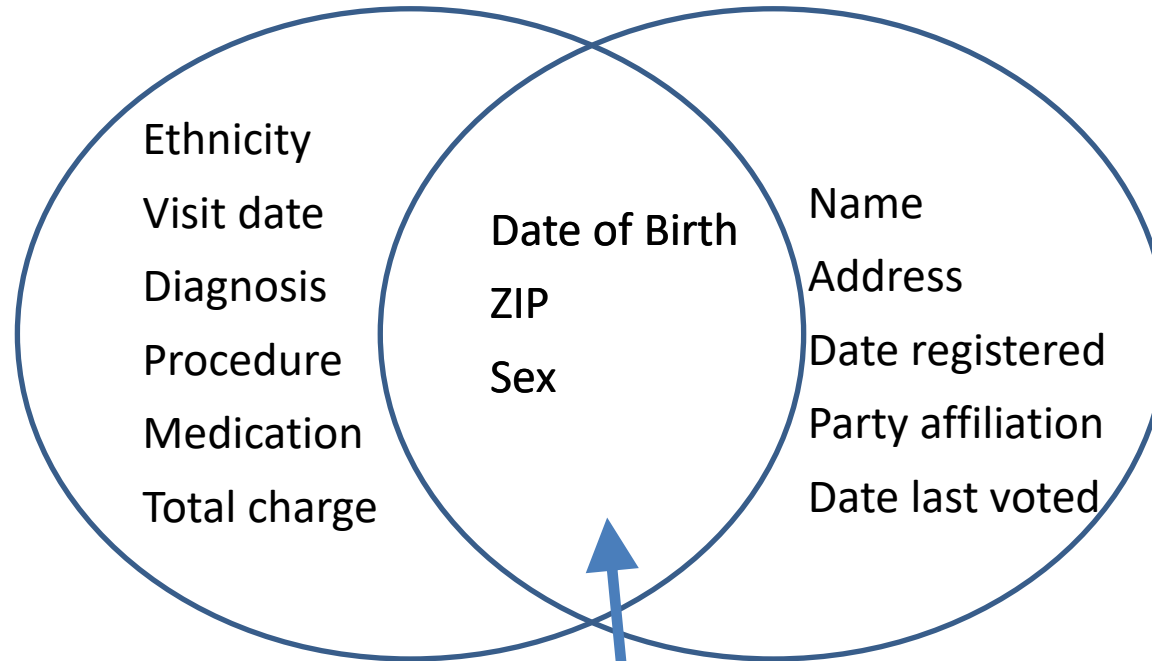
# Linkage attack (Sweeney '97)

- A unique record fully de-anonymize the record

- (DoB, ZIP, Sex) uniquely identifies 87% of US population

- Re-identified medical records of William Weld
(MA governor at the time)

- In Cambridge voters list
  - Six people shared his DoB
  - Three of which were men
  - He was the only one in his ZIP code

- Significant impact on privacy policymaking
and the health privacy legislation HIPAA
(Health Insurance Portability
and Accountability Act)

# AOL search history release (2006)

- In Aug 4th, 2006 AOL released users search requests to the public
- 20 million queries by 650,000 users over 3 months
- Goal: provide real query data by real users
- IP address replaced by random numbers
- In Aug 7th, 2006 AOL deleted the data

# AOL search history release (2006)

| | | | | |
|---|---|---|---|---|
| 4417749 best dog for older owner | 3/6/2006 | 11:48:24 | 1 | http://www.canismajor.com |
| 4417749 best dog for older owner | 3/6/2006 | 11:48:24 | 5 | http://dogs.about.com |
| 4417749 landscapers in lilburn ga. | 3/6/2006 | 18:37:26 | | |
| 4417749 effects of nicotine | 3/7/2006 | 19:17:19 | 6 | http://www.nida.nih.gov |
| 4417749 best retirement in the world | 3/9/2006 | 21:47:26 | 4 | http://www.escapeartist.com |
| 4417749 best retirement place in usa | 3/9/2006 | 21:49:37 | 10 | http://www.clubmarena.com |
| 4417749 best retirement place in usa | 3/9/2006 | 21:49:37 | 9 | http://www.committment.com |
| 4417749 bi polar and heredity | 3/13/2006 | 20:57:11 | | |
| 4417749 adventure for the older american | 3/17/2006 | 21:35:48 | | |
| 4417749 nicotine effects on the body | 3/26/2006 | 10:31:15 | 3 | http://www.geocities.com |
| 4417749 nicotine effects on the body | 3/26/2006 | 10:31:15 | 2 | http://health.howstuffworks.com |
| 4417749 wrinkling of the skin | 3/26/2006 | 10:38:23 | | |
| 4417749 mini strokes | 3/26/2006 | 14:56:56 | 1 | http://www.ninds.nih.gov |
| 4417749 panic disorders | 3/26/2006 | 14:58:25 | | |
| 4417749 jarrett t. arnold eugene oregon | 3/23/2006 | 21:48:01 | 2 | http://www2.eugeneweekly.com |
| 4417749 jarrett t. arnold eugene oregon | 3/23/2006 | 21:48:01 | 3 | http://www2.eugeneweekly.com |
| 4417749 plastic surgeons in gwinnett county | 3/28/2006 | 15:04:23 | 1 | http://www.wedalert.com |
| 4417749 plastic surgeons in gwinnett county | 3/28/2006 | 15:04:23 | 4 | http://www.implantinfo.com |
| 4417749 plastic surgeons in gwinnett county | 3/28/2006 | 15:31:00 | | |
| 441774960 single men | 3/29/2006 | 20:11:52 | 6 | http://www.adultlovecompass.com |
| 441774960 single men | 3/29/2006 | 20:14:14 | | |
| 4417749 clothes for 60 plus age | 4/19/2006 | 12:44:03 | | |
| 4417749 clothes for age 60 | 4/19/2006 | 12:44:41 | 10 | http://www.news.cornell.edu |
| 4417749 clothes for age 60 | 4/19/2006 | 12:45:41 | | |
| 4417749 lactose intolerant | 4/21/2006 | 20:53:51 | 2 | http://digestive.niddk.nih.gov |
| 4417749 lactose intolerant | 4/21/2006 | 20:53:51 | 10 | http://www.netdoctor.co.uk |
| 4417749 dog who urinate on everything | 4/28/2006 | 13:24:07 | 6 | http://www.dogdaysusa.com |
| 4417749 fingers going numb | 5/2/2006 | 17:35:47 | | |

# AOL search history release (2006)

**The New York Times**

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

By Michael Barbaro and Tom Zeller Jr.

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.

And search by search, click by click, the identity of AOL user No. 4417749 became easier to discern. There are queries for "landscapers in Lilburn, Ga," several people with the last name Arnold and "homes sold in shadow lake subdivision gwinnett county georgia."



Thelma Arnold, 62
Widow
Lives in Lilburn, GA

**Data itself leaks PII**

# Netflix Prize (2006)

- Netflix recommends movies to its subscribers
- In 2006 offered $1,000,000 for 10% improvement in its algorithm
- Published training data:
  - More than 100 million ratings from over 480,000 randomly chosen anonymous users on nearly 18,000 movie titles
  - All PII have been removed, all customer id replaced by random numbers
- Prize won by Bellkore's Pragmatic Chaos team, 2009

# Netflix Prize (2006)

- Anonymized data included: rating (1-5 stars), date, watch/didn't watch
- 213 dated ratings per used, on average
- Narayanan and Shmatikov re-identified the data

# Netflix Prize (2006)

- A source of auxiliary information: **IMDb**
  - Individuals may rate movies
  - Many use their real identify (not anonymous)
  - Visible data includes ratings, dates, comments

## IMDb Datasets

Subsets of IMDb data are available for access to customers for personal and non-commercial use. You can hold local copies of this data, and it is subject to our terms and conditions. Please refer to the Non-Commercial Licensing and copyright/license and verify compliance.

**Data Location**

The dataset files can be accessed and downloaded from https://datasets.imdbws.com/. The data is refreshed daily.

**IMDb Dataset Details**

Each dataset is contained in a gzipped, tab-separated-values (TSV) formatted file in the UTF-8 character set. The first line in each file contains headers that describe what is in each column. A '\N' is used to denote that a particular field is missing or null for that title/name. The available datasets are as follows:

# Netflix Prize (2006)

- **Sparse data cannot be anonymized!**
- Considering just watch/didn't watch for 90% of the records there isn't a **single** other record which is more than 30% similar
- Focus on movies that are not in top 10,000
- The whole point of privacy is that my record is similar to other records
- Here, to make two records "close" the data is destroyed



https://arxiv.org/abs/cs/0610105

# Netflix Prize (2006)

Results of the attack

- With 8 movie ratings and dates that may have a 3-days error, 96% of Netflix clients whose data was released can be uniquely identified in the dataset

- For 89%, 2 ratings and dates are enough to reduce the set of plausible records to 8 out of almost 500,000

Consequences

- Learn about movies that IMDb users didn't want to tell the world: sexual orientation, religious beliefs, political attitude, etc.

- In 2009 four Netflix users filled a lawsuit against Netflix

- In 2010 Netflix cancelled the second prize competition

# Privacy is more than re-identification

Medical encounter data

- Ambulance collects an elderly neighbor

- Daily medical encounter data shows that every elderly admitted patient was diagnosed with tachycardia, influenza, broken arm, panic attack

- Learn the neighbor suffers from one of these 4 complaints

- Next day, can rule out influenza, broken arm

- Re-identification fails to capture privacy risks!

# NYC Taxi and Limo Commission (2014)

- TLC is the regulator for establishing public transport policy setting and enforcing the fare rate in taxis, etc.

- Published statistics about taxi rides

- 2014 Whong filled a FOILed request (Freedom Of Information Law)

- Got 2 datasets (90 GB of data) trips and fares

https://chriswhong.com/open-data/foil_nyc_taxi/

# NYC Taxi and Limo Commission (2014)

| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | medallion | hack_license | vendor_id | pickup_datetime | payment_type | fare_amoun | surcharge | mta_tax | tip_amount | tolls_amoun | total_amount |
| 2 | 89D227B655E5C82AECF13C3I | BA96DE419E711691B944 | CMT | 1/1/13 15:11 | CSH | 6.5 | 0 | 0.5 | 0 | 0 | 7 |
| 3 | 0BD7C8F5BA12B88E0B67BED | 9FD8F69F0804BDB5549F | CMT | 1/6/13 0:18 | CSH | 6 | 0.5 | 0.5 | 0 | 0 | 7 |
| 4 | 0BD7C8F5BA12B88E0B67BED | 9FD8F69F0804BDB5549F | CMT | 1/5/13 18:49 | CSH | 5.5 | 1 | 0.5 | 0 | 0 | 7 |
| 5 | DFD2202EE08F7A8DC9A57B0 | 51EE87E3205C985EF843: | CMT | 1/7/13 23:54 | CSH | 5 | 0.5 | 0.5 | 0 | 0 | 6 |
| 6 | DFD2202EE08F7A8DC9A57B0 | 51EE87E3205C985EF843: | CMT | 1/7/13 23:25 | CSH | 9.5 | 0.5 | 0.5 | 0 | 0 | 10.5 |
| 7 | 20D9ECB2CA0767CF7A015641 | 598CCE5B9C1918568DEE | CMT | 1/7/13 15:27 | CSH | 9.5 | 0 | 0.5 | 0 | 0 | 10 |
| 8 | 496644932DF3932605C22C79 | 513189AD756FF14FE670 | CMT | 1/8/13 11:01 | CSH | 6 | 0 | 0.5 | 0 | 0 | 6.5 |
| 9 | 0B57B9633A2FECD3D3B1944 | CCD4367B417ED6634D9I | CMT | 1/7/13 12:39 | CSH | 34 | 0 | 0.5 | 0 | 4.8 | 39.3 |
| 10 | 2C0E91FF20A856C891483ED6 | 1DA2F6543A62B8ED934: | CMT | 1/7/13 18:15 | CSH | 5.5 | 1 | 0.5 | 0 | 0 | 7 |

# NYC Taxi and Limo Commission (2014)

```
6B111958A39B24140C973B262EA9FEA5,D3B035A03C8A34DA17488129DA581EE7,VTS,5,,2013-12-03
15:46:00,2013-12-03 16:47:00,1,3660,22.71,-73.813927,40.698135,-74.093307,40.829346

medallion, hack_license, vendor_id, rate_code, store_and_fwd_flag, pickup_datetime,
dropoff_datetime, passenger_count, trip_time_in_secs, trip_distance, pickup_longitude,
pickup_latitude, dropoff_longitude, dropoff_latitude
```

- MD5 values of taxi number and driver license
- After a taxi ride one can learn information about the driver
- If someone is taking a taxi you can see where they're going
- Are they good tippers

https://chriswhong.com/open-data/foil_nyc_taxi/

# Identifiers vs. Sensitive attributes

- Key attributes: name, address, etc. (uniquely identifying)
- Quasi-identifiers: ZIP, DoB, etc.
- Sensitive attributes: medical records, etc.

| Key Attribute | | Quasi-identifier | | Sensitive attribute |
|---|---|---|---|---|
| Name | DOB | Gender | Zipcode | Disease |
| Andre | 1/21/76 | Male | 53715 | Heart Disease |
| Beth | 4/13/86 | Female | 53715 | Hepatitis |
| Carol | 2/28/76 | Male | 53703 | Brochitis |
| Dan | 1/21/76 | Male | 53703 | Broken Arm |
| Ellen | 4/13/86 | Female | 53706 | Flu |
| Eric | 2/28/76 | Female | 53706 | Hang Nail |

29

# k-Anonymity (Sweeney and Samarati 98)

- The information for each person contained in the released table cannot be distinguished from at least $k - 1$ individuals whose information also appears in the release
- Any quasi-identifier present in the released table must appear in at least $k$ records
- Simple and syntactic property of the dataset
- Very popular technique

# k-Anonymity (Sweeney and Samarati 98)

| | Race | Birth | Gender | ZIP | Problem |
|---|---|---|---|---|---|
| t1 | Black | 1965 | m | 0214* | short breath |
| t2 | Black | 1965 | m | 0214* | chest pain |
| t3 | Black | 1965 | f | 0213* | hypertension |
| t4 | Black | 1965 | f | 0213* | hypertension |
| t5 | Black | 1964 | f | 0213* | obesity |
| t6 | Black | 1964 | f | 0213* | chest pain |
| t7 | White | 1964 | m | 0213* | chest pain |
| t8 | White | 1964 | m | 0213* | obesity |
| t9 | White | 1964 | m | 0213* | short breath |
| t10 | White | 1967 | m | 0213* | chest pain |
| t11 | White | 1967 | m | 0213* | chest pain |

Figure 2 Example of *k*-anonymity, where *k*=2 and **Ql**={*Race, Birth, Gender, ZIP*}

https://epic.org/privacy/reidentification/Sweeney_Article.pdf

# k-Anonymity (Sweeney and Samarati 98)

**Released table**

| | Race | Birth | Gender | ZIP | Problem |
|---|---|---|---|---|---|
| t1 | Black | 1965 | m | 0214* | short breath |
| t2 | Black | 1965 | m | 0214* | chest pain |
| t3 | Black | 1965 | f | 0213* | hypertension |
| t4 | Black | 1965 | f | 0213* | hypertension |
| t5 | Black | 1964 | f | 0213* | obesity |
| t6 | Black | 1964 | f | 0213* | chest pain |
| t7 | White | 1964 | m | 0213* | chest pain |
| t8 | White | 1964 | m | 0213* | obesity |
| t9 | White | 1964 | m | 0213* | short breath |
| t10 | White | 1967 | m | 0213* | chest pain |
| t11 | White | 1967 | m | 0213* | chest pain |

**External data source**

| Name | Birth | Gender | ZIP | Race |
|---|---|---|---|---|
| Andre | 1964 | m | 02135 | White |
| Beth | 1964 | f | 55410 | Black |
| Carol | 1964 | f | 90210 | White |
| Dan | 1967 | m | 02174 | White |
| Ellen | 1968 | f | 02237 | White |

# k-Anonymity (Sweeney and Samarati 98)

**Microdata**

| QID | | | SA |
|---|---|---|---|
| Zipcode | Age | Sex | Disease |
| 47677 | 29 | F | Ovarian Cancer |
| 47602 | 22 | F | Ovarian Cancer |
| 47678 | 27 | M | Prostate Cancer |
| 47905 | 43 | M | Flu |
| 47909 | 52 | F | Heart Disease |
| 47906 | 47 | M | Heart Disease |

**Generalized table**

| QID | | | SA |
|---|---|---|---|
| Zipcode | Age | Sex | Disease |
| 476** | 2* | * | Ovarian Cancer |
| 476** | 2* | * | Ovarian Cancer |
| 476** | 2* | * | Prostate Cancer |
| 4790* | [43,52] | * | Flu |
| 4790* | [43,52] | * | Heart Disease |
| 4790* | [43,52] | * | Heart Disease |

- Released table is 3-anonymous
- Alice's quasi-identifier (47677, 29, F) does not reveal her disease

# k-Anonymity (Sweeney and Samarati 98)

- Unsorted matching attack
- Records appear in the same order as in the original table
- Solution: randomize order before releasing

| Race | ZIP |
|------|-----|
| Asian | 02138 |
| Asian | 02139 |
| Asian | 02141 |
| Asian | 02142 |
| Black | 02138 |
| Black | 02139 |
| Black | 02141 |
| Black | 02142 |
| White | 02138 |
| White | 02139 |
| White | 02141 |
| White | 02142 |

PT

| Race | ZIP |
|------|-----|
| Person | 02138 |
| Records | 02139 |
| Person | 02141 |
| Person | 02142 |
| Person | 02138 |
| Person | 02139 |
| Person | 02141 |
| Person | 02142 |
| Person | 02138 |
| Person | 02139 |
| Person | 02141 |
| Person | 02142 |

GT1

# Quiz: what does k-Anonymity provide

- Membership discloser:
  attacker cannot tell that a given person is in the dataset

- Sensitive attribute discloser:
  attacker cannot tell that a given person has a certain sensitive attribute

- Identity discloser:
  attacker cannot tell which record corresponds to which person

# Quiz: what does k-Anonymity provide

- Membership discloser:
  attacker cannot tell that a given person is in the dataset

- Sensitive attribute discloser:
  attacker cannot tell that a given person has a certain sensitive attribute

- Identity discloser:
  attacker cannot tell which record corresponds to which person

This interpretation is correct,
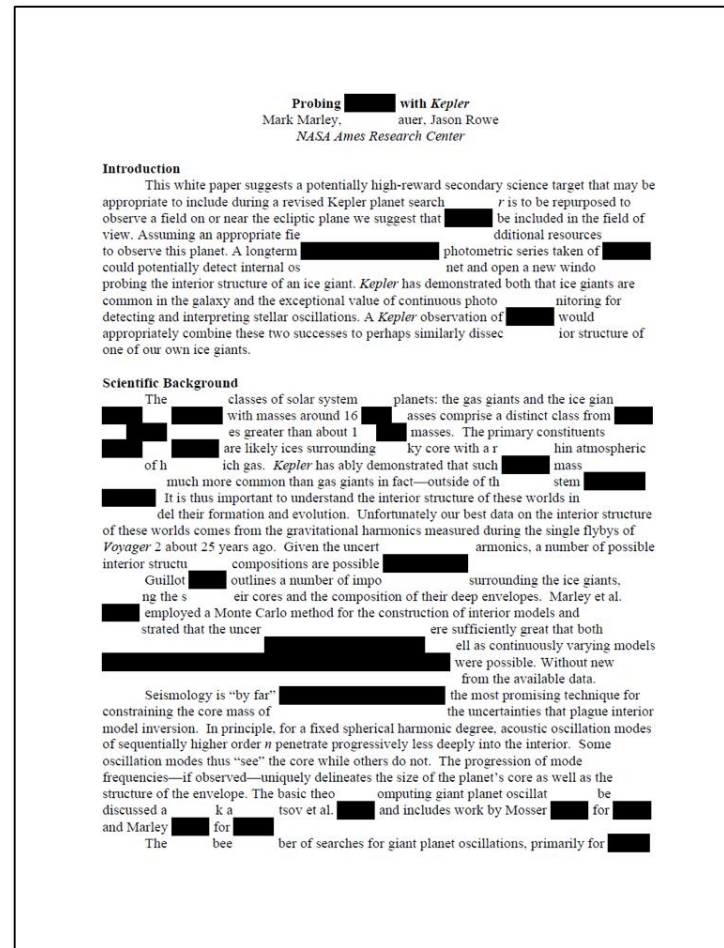 assuming the attacker does not know anything other than quasi-identifiers

# A chain of measures and counter measures

- $k$-anonymity  [Sweeney and Samarati 98]

- Attacks against $k$-anonymity [Machanavajjhala et al. 06]
  Proposed $L$-diversity

- Attacks against $L$-diversity [Xiao and Tao 07]
  Proposed $M$-invariance

- Proposed $T$-closeness [Li et al. 07]

- Attacks against all the above [Ganta, Kasiviswanathan, Smith 08]

# Popular idea #2

- Information not explicitly given cannot be harmful
- E.g., redaction

**Probing** ▮ with *Kepler*

Mark Marley, ▮auer, Jason Rowe

*NASA Ames Research Center*

**Introduction**

This white paper suggests a potentially high-reward secondary science target that may be appropriate to include during a revised Kepler planet search ▮ *r* is to be repurposed to observe a field on or near the ecliptic plane we suggest that ▮ be included in the field of view. Assuming an appropriate fie ▮ dditional resources to observe this planet. A longterm ▮ photometric series taken of ▮ could potentially detect internal os ▮ net and open a new windo ▮ probing the interior structure of an ice giant. *Kepler* has demonstrated both that ice giants are common in the galaxy and the exceptional value of continuous photo ▮ nitoring for detecting and interpreting stellar oscillations. A *Kepler* observation of ▮ would appropriately combine these two successes to perhaps similarly dissec ▮ ior structure of one of our own ice giants.

**Scientific Background**

The ▮ classes of solar system ▮ planets: the gas giants and the ice gian ▮ with masses around 16 ▮ asses comprise a distinct class from ▮ es greater than about 1 ▮ masses. The primary constituents ▮ are likely ices surrounding ▮ ky core with a r ▮ hin atmospheric of h ▮ ich gas. *Kepler* has ably demonstrated that such ▮ mass ▮ much more common than gas giants in fact—outside of th ▮ stem ▮ It is thus important to understand the interior structure of these worlds in ▮ del their formation and evolution. Unfortunately our best data on the interior structure of these worlds comes from the gravitational harmonics measured during the single flybys of *Voyager* 2 about 25 years ago. Given the uncert ▮ armonics, a number of possible interior structu ▮ compositions are possible ▮

Guillot ▮ outlines a number of impo ▮ surrounding the ice giants, ng the s ▮ eir cores and the composition of their deep envelopes. Marley et al. ▮ employed a Monte Carlo method for the construction of interior models and strated that the uncer ▮ ere sufficiently great that both ▮ ell as continuously varying models ▮ were possible. Without new ▮ from the available data.

Seismology is "by far" ▮ the most promising technique for constraining the core mass of ▮ the uncertainties that plague interior model inversion. In principle, for a fixed spherical harmonic degree, acoustic oscillation modes of sequentially higher order $n$ penetrate progressively less deeply into the interior. Some oscillation modes thus "see" the core while others do not. The progression of mode frequencies—if observed—uniquely delineates the size of the planet's core as well as the structure of the envelope. The basic theo ▮ omputing giant planet oscillat ▮ be discussed a ▮ k a ▮ tsov et al. ▮ and includes work by Mosser ▮ for ▮ and Marley ▮ for ▮ The ▮ bee ▮ ber of searches for giant planet oscillations, primarily for ▮

# Popular idea #2

- The President's Daily Brief (PDB) is a top-secret document given each morning to the US president

- August 6th, 2001 George W. Bush received a PDB Bin Laden and El Qaeda are planning to strike in the US

- Declassified and released to the 9/11 Commission in 2004



Declassified and Approved for Release, 10 April 2004

**Bin Ladin Determined To Strike in US**

Clandestine, foreign government, and media reports indicate Bin Ladin since 1997 has wanted to conduct terrorist attacks in the US. Bin Ladin implied in US television interviews in 1997 and 1998 that his followers would follow the example of World Trade Center bomber Ramzi Yousef and "bring the fighting to America."

After US missile strikes on his base in Afghanistan in 1998, Bin Ladin told followers he wanted to retaliate in Washington, according to a ████████ service.

An Egyptian Islamic Jihad (EIJ) operative told an ████ service at the same time that Bin Ladin was planning to exploit the operative's access to the US to mount a terrorist strike.

The millennium plotting in Canada in 1999 may have been part of Bin Ladin's first serious attempt to implement a terrorist strike in the US. Convicted plotter Ahmed Ressam has told the FBI that he conceived the idea to attack Los Angeles International Airport himself, but that Bin Ladin lieutenant Abu Zubaydah encouraged him and helped facilitate the operation. Ressam also said that in 1998 Abu Zubaydah was planning his own US attack.

Ressam says Bin Ladin was aware of the Los Angeles operation.

Although Bin Ladin has not succeeded, his attacks against the US Embassies in Kenya and Tanzania in 1998 demonstrate that he prepares operations years in advance and is not deterred by setbacks. Bin Ladin associates surveilled our Embassies in Nairobi and Dar es Salaam as early as 1993, and some members of the Nairobi cell planning the bombings were arrested and deported in 1997.

Al-Qa'ida members—including some who are US citizens—have resided in or traveled to the US for years, and the group apparently maintains a support structure that could aid attacks. Two al-Qa'ida members found guilty in the conspiracy to bomb our Embassies in East Africa were US citizens, and a senior EIJ member lived in California in the mid-1990s.

A clandestine source said in 1998 that a Bin Ladin cell in New York was recruiting Muslim-American youth for attacks.

We have not been able to corroborate some of the more sensational threat reporting, such as that from a ████████ service in 1998 saying that Bin Ladin wanted to hijack a US aircraft to gain the release of "Blind Shaykh" 'Umar 'Abd al-Rahman and other US-held extremists.

continued

For the President Only
6 August 2001
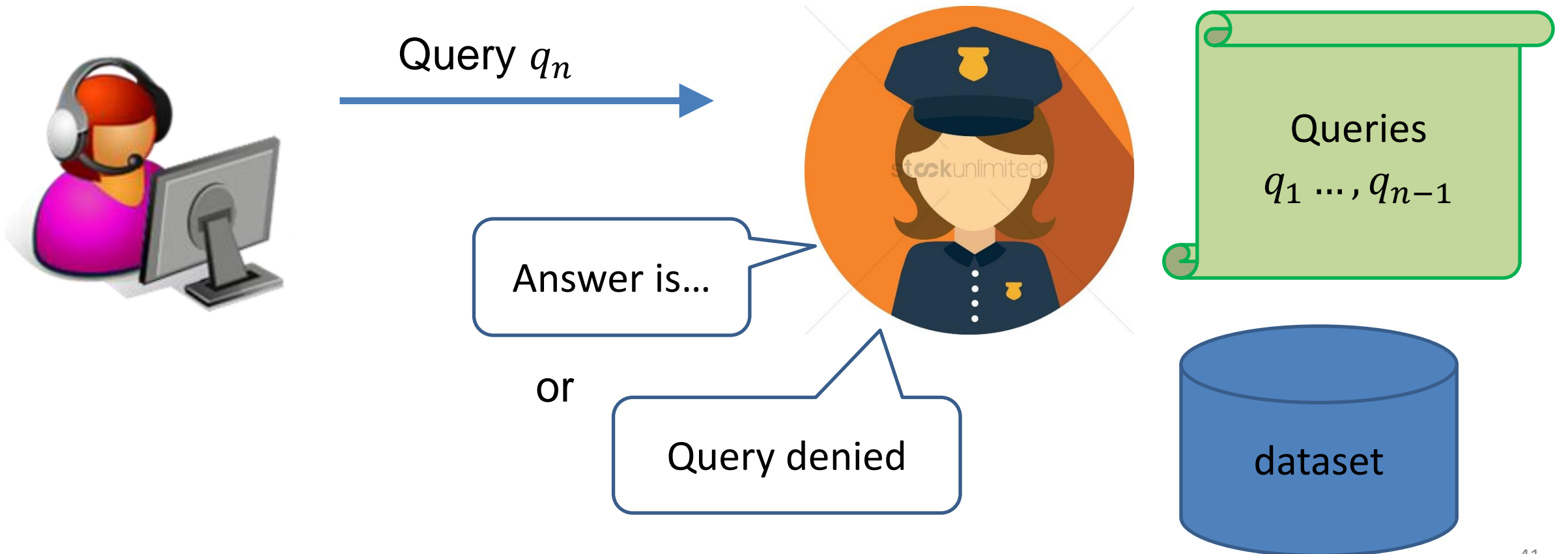
Declassified and Approved for Release, 10 April 2004

# Popular idea #2

- Naccache and Whelan analyzed the geometry of the font
- 1530 plausible words
- The "an" reduced to 7 candidates: Ukrainian, uninvited, unofficial, incursive, Egyptian, indebted and Ugandan
- Egyptian is the only one who made sense in the context

An Egyptian Islamic Jihad (EIJ) operative told an ██████████ service at the same time that Bin Ladin was planning to exploit the operative's access to the US to mount a terrorist strike.

# Query auditing

- Refuse to answer queries that would compromise privacy



Query $q_n$

Answer is…

or

Query denied

Queries
$q_1 \ldots, q_{n-1}$

dataset

# Example: sum/max auditing

- Sensitive info: $d_i$ (real)

$$q_1 = \text{sum}(d_1, d_2, d_3)$$

$$q_2 = \text{max}(d_1, d_2, d_3)$$

$$\text{sum}(d_1, d_2, d_3) = 15$$

query denied

dataset

# Example: sum/max auditing

- Sensitive info: $d_i$ (real)

$\max(d_1, d_2, d_3) \geq 5$

No denial if
$\max(d_1, d_2, d_3) > 5$

$\max(d_1, d_2, d_3) = 5$

$d_1 = d_2 = d_3 = 5$

$\mathrm{sum}(d_1, d_2, d_3) = 15$

query denied

dataset

# Popular idea #3: add noise

- Mask numbers by adding a random number between $[-a, a]$
  - Privacy $2a@100\%$ confidence, Privacy $a@50\%$ confidence, ...
- The larger the interval the better the privacy
- Example:
  - For each person mask the age by adding a random number between $[-100, 100]$
  - Gives privacy 200@100% confidence
  - But, masked age -99 $\Rightarrow$ a baby of age 0 or 1

# So far

- Many ideas fall short of providing data privacy
- Auxiliary information
- Data itself may leak information
- Sparse dataset cannot be anonymized
- Privacy is more than re-identifying

# Outline

- Popular ideas that do not work
  + privacy horror stories

- An approach that works

# What went wrong?

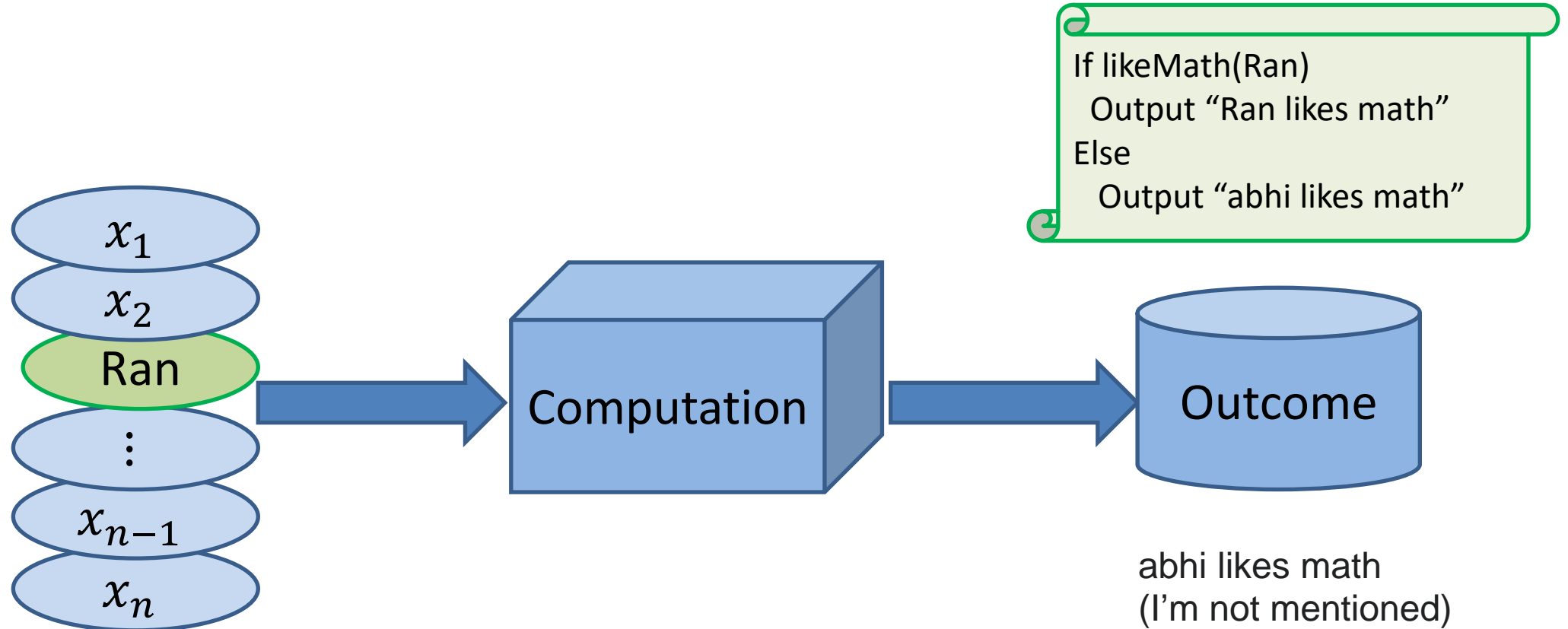**Privacy is NOT a property of the outcome but of the computation!!!**

# What went wrong?



Output "Ran likes math"

$x_1$

$x_2$

Ran

$\vdots$

$x_{n-1}$

$x_n$

Computation

Outcome

Ran likes math

Is my privacy breached? Yes / No / Cannot tell

# What went wrong?



If likeMath(Ran)
  Output "Ran likes math"
Else
  Output "abhi likes math"

abhi likes math
(I'm not mentioned)

Is my privacy breached? Yes / No / Cannot tell

# Recall sematic security

Real world

Ideal world

$\text{Enc}(m)$
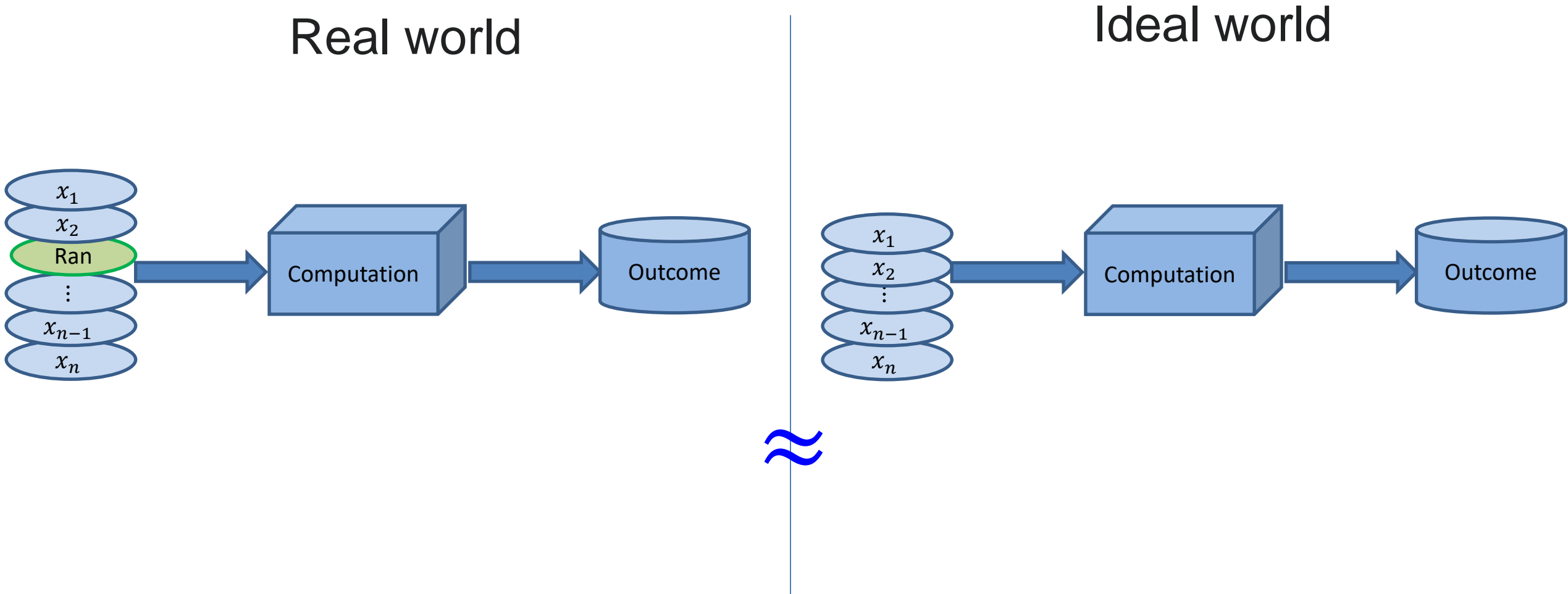
$\approx$

An encryption scheme is semantically secure if whatever can be learned given the ciphertext can be learned without the ciphertext

# Privacy analogue



A computation is "private" if whatever can be learned with my record in the DB can be learned without my record

# Differential Privacy
# [Dwork, McSherry, Nissim, Smith 2006]

A mechanism / algorithm / computation $M$ has **$\varepsilon$-differential privacy** if for any pair of neighboring databased $D_1, D_2$ (differing by 1 record) and for any $S \subseteq \text{Range}(M)$

$$\Pr[M(D_1) \in S] \leq e^{\varepsilon} \cdot \Pr[M(D_2) \in S]$$

# Differential Privacy

Adopted by:

- US census Bureau
- Google
- Apple
- YouTube
- Many more