

W4111

# Introduction to Databases

## Spring 2019

Computer Science Department  
Columbia University

# Data

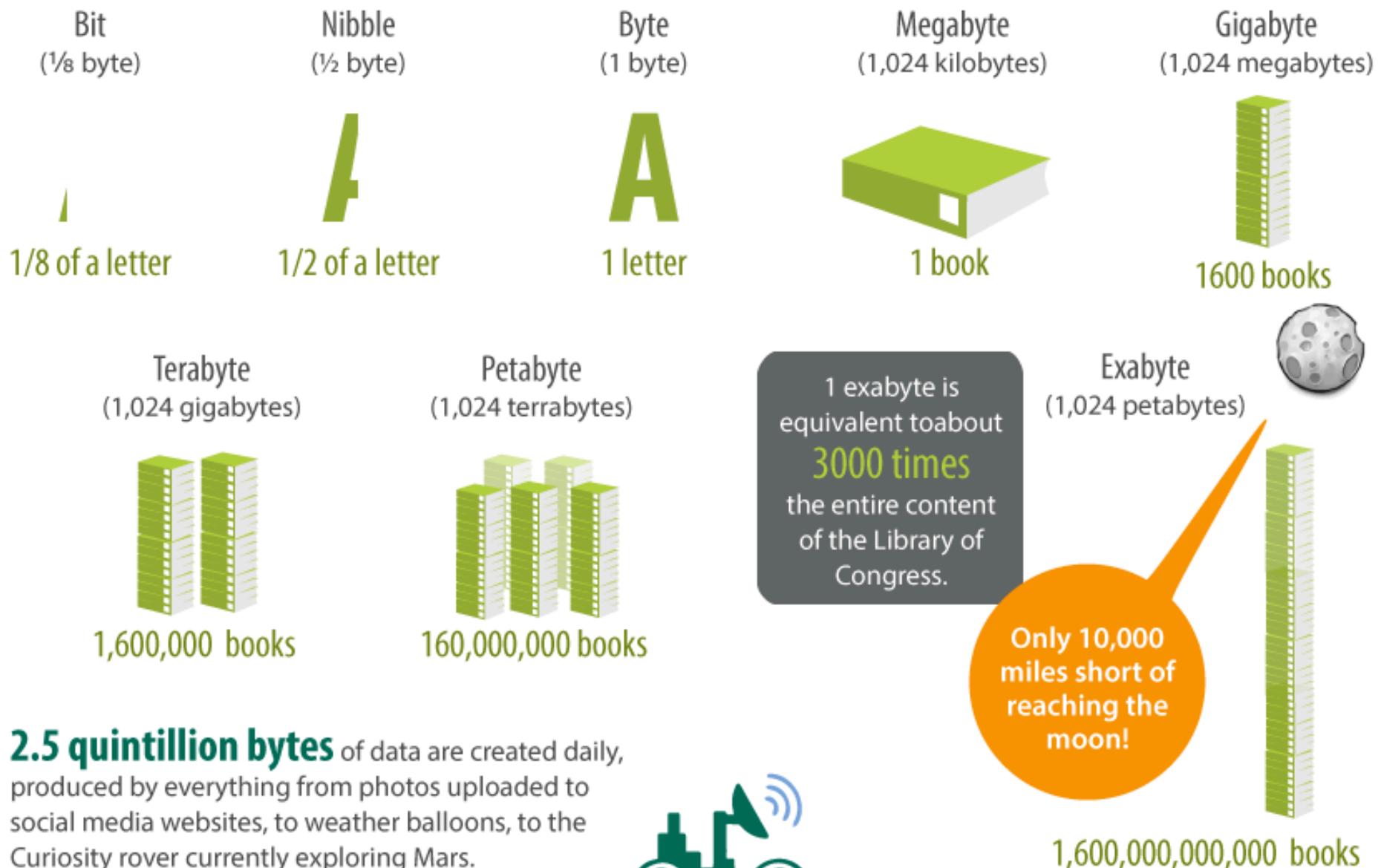
**Data**  
**is for serious business**

**Data**  
is at the center of most things.

**Data**  
is at the center of *everything*



# Data Sizes



**2.5 quintillion bytes** of data are created daily, produced by everything from photos uploaded to social media websites, to weather balloons, to the Curiosity rover currently exploring Mars.

Bigger Than Big Data

## ► THE PAST

Digital storage grew annually by **23%** between 1986 and 2007.

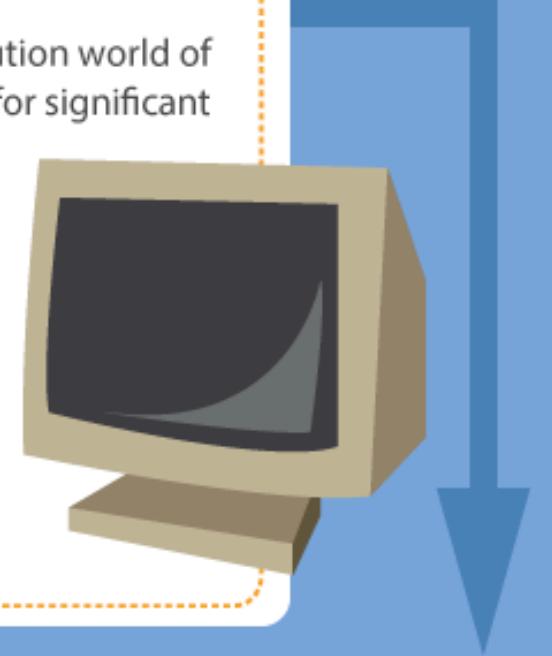
Most data was stored on **videotapes** such as VHS cassettes in the pre-digital revolution world of the late 1980s, Vinyl LP records, audio cassette tapes, and photography accounted for significant portions as well.

Paper-based storage represented **33%** of all data storage on its own in 1986.

**25%** of all data stored in the world in 2000 was stored digitally.

2002 is the first year that digital storage capacity overtook analog capacity.

**94%** of all data was stored in digital format by 2007.



## ► PRESENT

Today, more than **2.5 exabytes** (2.5 billion gigabytes) of data is generated every single day. This is expected to continue growing at a significant rate with mobile devices accounting for much of this data.

Some experts have estimated that **90%** of all of the data the world today was produced within the last two years.

# How did we get here?

# Data was Manual

67

June 11<sup>1928</sup>

Geo. A. Kelly  
June 16 Mrs. Chas. Long Jr  
June 16 Nellora Wright  
June 14 Charity A. Bonds  
" " Mr. H. A. Carpenter  
" " Mr. & Mrs. Carpenter  
July 10 James Ostrom Troop 251  
July 10 F. W. Gemmings  
July 10 Millicent Gemmings  
Walt Klein  
July 11 Mrs. Rawo. & Daughter  
" " Mrs. Ralph Pease  
" " Mrs. A. H. Favout  
" " Mrs. J. A. Miller  
Mrs. C. J. Morris  
Mrs. O. J. Vista  
Mary S. Gossypium  
Mrs. L. Holden Hoffman  
Mrs. Key & Young  
Mrs. A. L. Whitney  
Miss. Mrs. T. G. Gadda  
Mrs. J. L. Roberts

Phoenix, Arizona.  
Phoenix Arizona  
Phoenix Arizona  
Prescott - Arizona.  
Sage Granulated  
Prescott

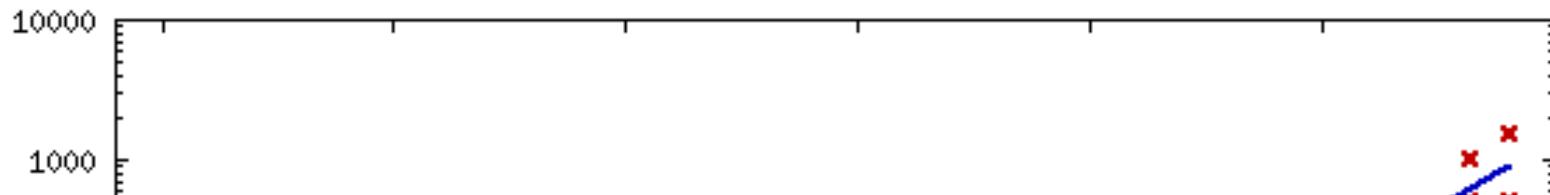
8. M. Vernon St. Prescott  
Dewey Arizona  
Deer Valley, Arizona  
Ph. 219 - N.Y.  
Sage Granulated Calif.  
Prescott  
"

Arizona  
"

# Data was *Expensive*



# Data is Cheap



**U32 Shadow™ 1TB External USB 3.0 Portable Hard Drive**  
by Oyen Digital

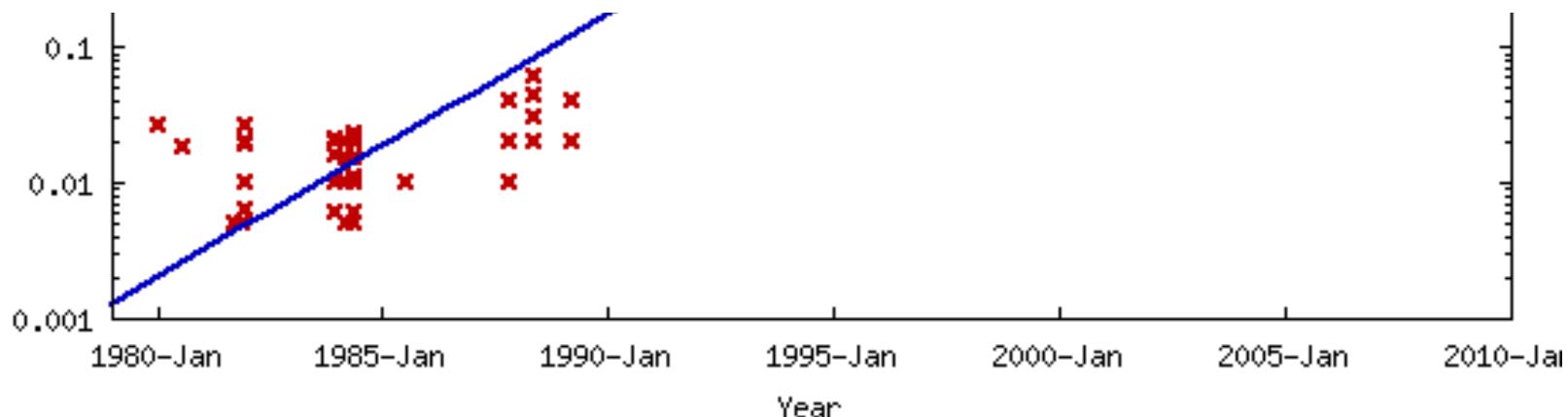
**\$84.95 ✓Prime**

Only 4 left in stock - order soon.

More Buying Choices

**\$84.95 new** (2 offers)

**\$74.33 used** (1 offer)



# Data is Automated

## Physical devices



# *Data is Automated*

Physical devices

Software logs

# Data is *Ubiquitous*

Physical devices

Software logs

Phones



# Data is *Ubiquitous*

Physical devices

Software logs

Phones

GPS/Cars



1-888-411-2188

STORES

sleep  number.

- Information You provide to Us on a financing application
- Information You provide to Us to allow Your participation in a partner's loyalty program

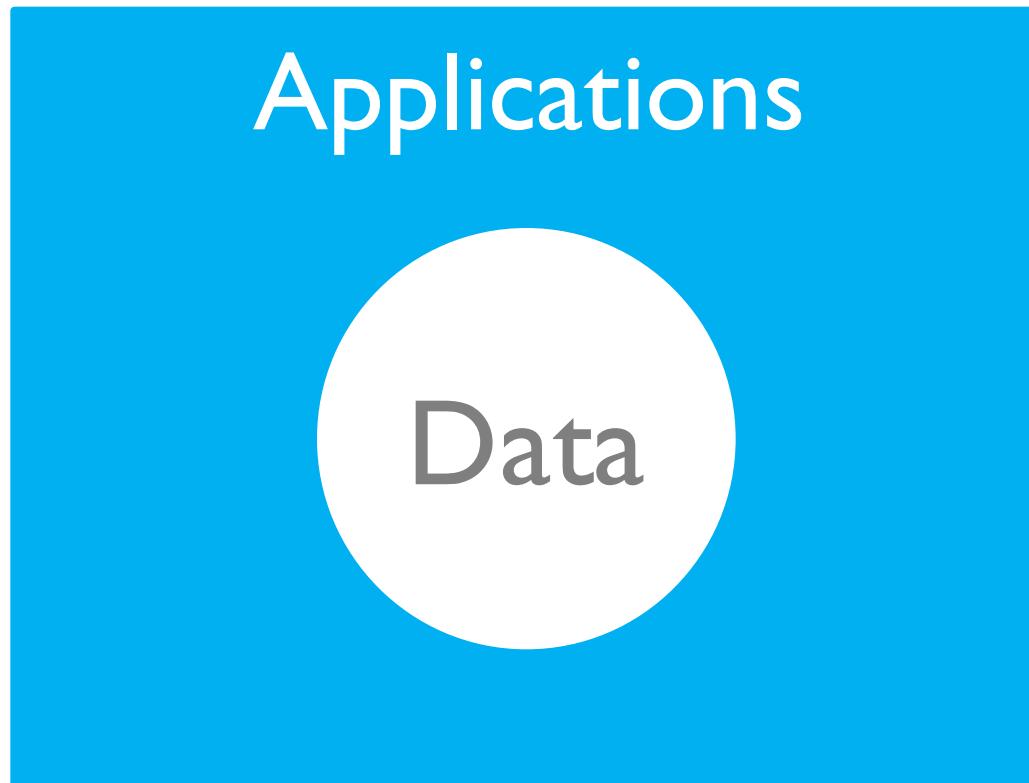
Once You create a User Profile, We also may collect Personal Information, which may include, among other types of information:

- Revised or updated User Profile information
- Biometric and sleep-related data about how You, a Child, and any person that uses the Bed slept, such as that person's movement, positions, respiration, and heart rate while sleeping
-  Audio in Your room to detect snoring and similar sleep conditions





# What Applications?



# What are we doing with data?

Health



# What are we doing with data?

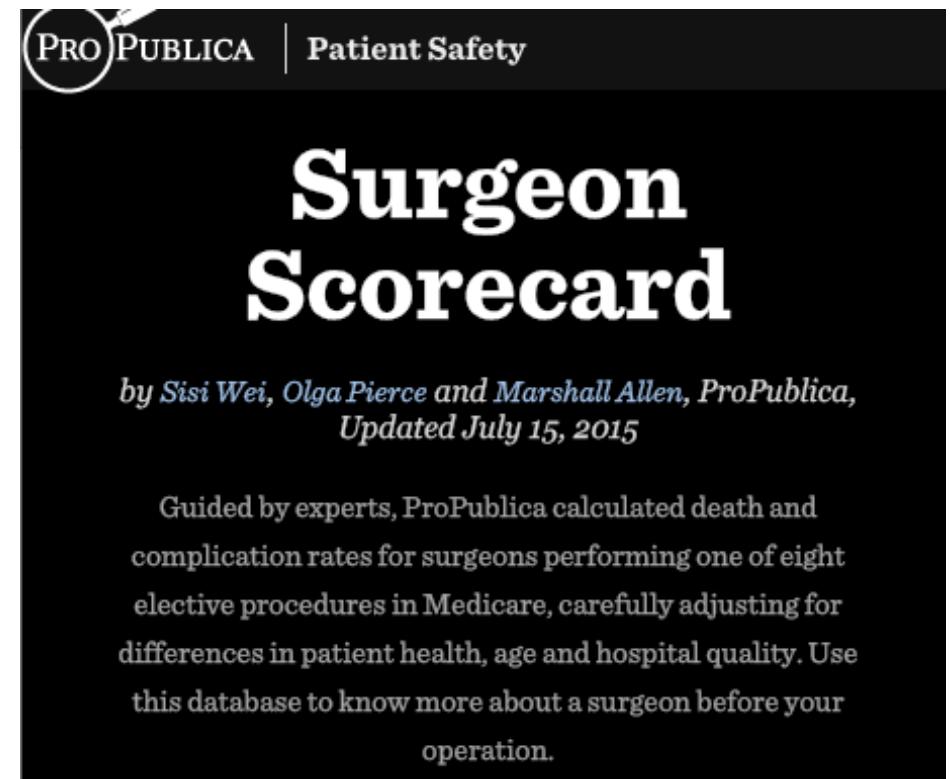
Health



# What are we doing with data?

Health

Investigative Journalism



The image shows a screenshot of the ProPublica Surgeon Scorecard website. At the top left is the ProPublica logo, which includes a circular icon with the letters 'PRO' and the word 'PUBLICA'. To the right of the logo is the text 'Patient Safety'. Below the header, the title 'Surgeon Scorecard' is displayed in large, bold, white font. Underneath the title, the authors are credited as 'by Sisi Wei, Olga Pierce and Marshall Allen, ProPublica, Updated July 15, 2015'. A descriptive paragraph follows, explaining that the data was calculated by experts using Medicare records to determine death and complication rates for eight elective procedures, adjusting for patient health, age, and hospital quality. It encourages users to use the database to learn about surgeons before their operations.

PRO PUBLICA | Patient Safety

# Surgeon Scorecard

by Sisi Wei, Olga Pierce and Marshall Allen, ProPublica, Updated July 15, 2015

Guided by experts, ProPublica calculated death and complication rates for surgeons performing one of eight elective procedures in Medicare, carefully adjusting for differences in patient health, age and hospital quality. Use this database to know more about a surgeon before your operation.

# What are we doing with data?

Health

Investigative Journalism

Recommendations



# What are we doing with data?



MACHINE BIAS



## Besieged Facebook Says New Ad Limits Aren't Response to Lawsuits



The social network is removing 5,000 options that regulators say enable advertisers to discriminate.



by Ariana Tobin and [Jeremy B. Merrill](#), Aug. 23, 12:48 p.m. EDT



Facebook's move to eliminate 5,000 options that enable advertisers on its platform to limit their audiences is unrelated to lawsuits accusing it of fostering housing and

FOLLOW PROPUBLICA

Twitter

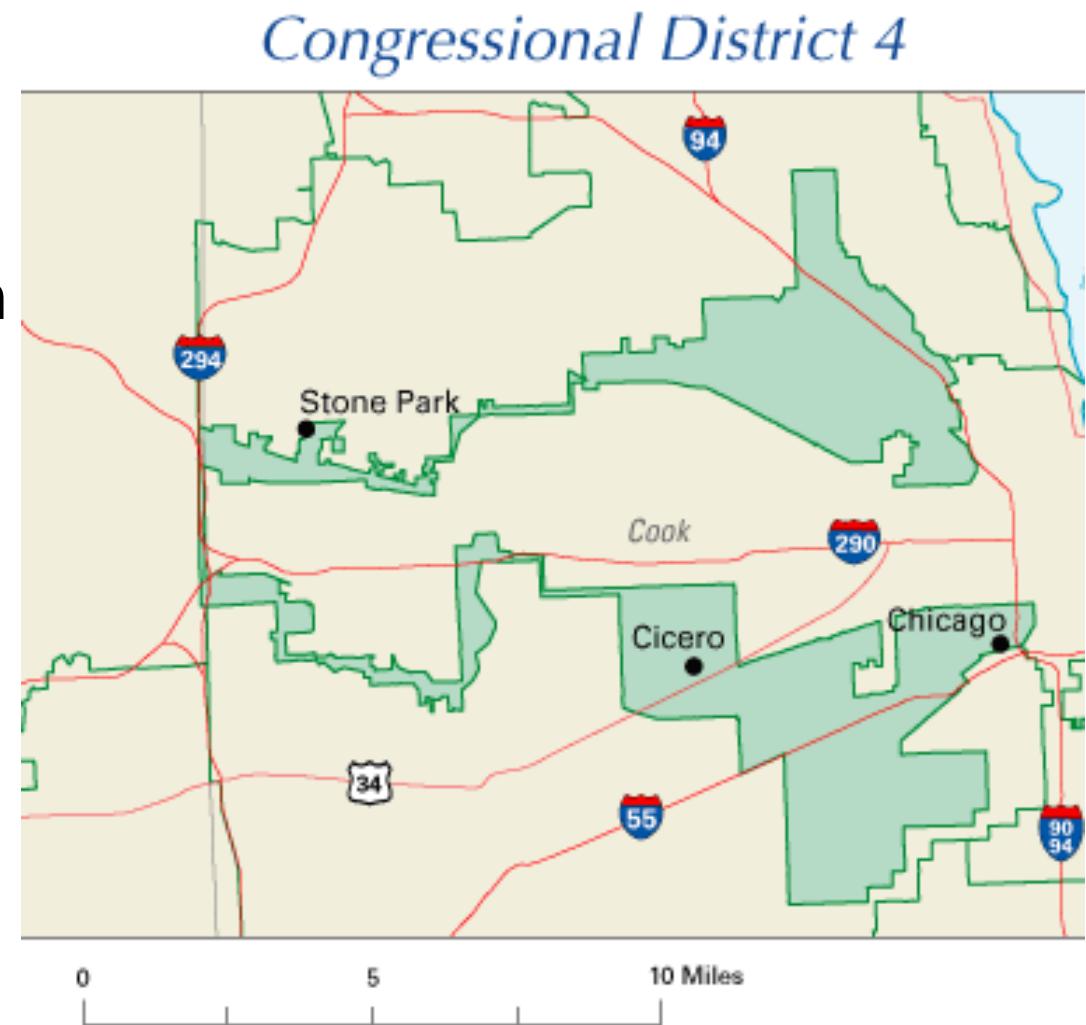
# What are we doing with data?

Health

Investigative Journalism

Recommendations

Politics



# What are we doing with data?

Health

Investigative Journalism

Recommendations

Politics

Make a  
contribution

Subscribe

Find a job

Sign in / Register

Search ▾

News

Opinion

Sport

Culture

Lifestyle

More ▾

US edition ▾

The  
Guardian



The  
Cambridge  
Analytica  
Files

A year-long investigation into  
Facebook, data, and influencing  
elections in the digital age

Key stories

Hide



TIME

Subscribe

2012 ELECTION

**Inside the Secret World of the  
Data Crunchers Who Helped  
Obama Win**

Data-driven decisionmaking played a huge role in creating a second term for the 44th President and will be one of the more closely studied elements of the 2012 cycle

# What are we doing with data?

Health

Investigative Journalism



set station news arts & life music programs

shop

parallels MANY STORIES, ONE WORLD



4:21

+ QUEUE

DOWNLOAD

EMBED

Facial Recognition In China Is Big Business As Local Governments Boost Surveillance

April 3, 2018 · 10:40 AM ET



ROB SCHMITZ



Every day, the NSA intercepts and stores **1.7 billion** emails, phone calls, texts and other electronic

# What are we doing with data?

Health

Investigative Journalism

Recommendations

Politics

Surveillance

Identity



30 APR 2012 RESEARCH & IDEAS

## India's Ambitious National Identification Program

Comments 30 Email Print Download Share [f Recommend](#) Share 92

The Unique Identification Authority of India has been charged with implementing a nationwide program to register and assign a unique 12-digit ID to every Indian resident—some 1.2 billion people—by 2020. In a new case, Professor Tarun Khanna and HBS India Research Center Executive Director Anjali Raina discuss the complexities of this massive data management project.

**“YOU ARE BASICALLY DENIED ALMOST EVERYTHING IF YOU CAN'T PROVE WHO YOU ARE.”**

# What data?

Applications



Data

# What data?

Fake data



# What data?

# Fake data

# Biased data



DATA    TUTORIALS    FINDINGS    PUBLICATIONS    NEWS

AK									ME		
WA	ID	MT	ND	MN		MI		NY	MA	RH	
OR	UT	WY	SD	IA	WI	IN	OH	PA	NJ	CT	
CA	NV	CO	NE	MO	IL	KY	WV	VA	MD	DE	
	AZ	NM	KS	AR	TN	NC	SC	DC			
			OK	LA	MS	AL	GA				

# THE STANFORD OPEN POLICING PROJECT

On a typical day in the United States, police officers make more than 50,000 traffic stops. Our team is gathering, analyzing, and releasing records from millions of traffic stops by law enforcement agencies across the country. Our goal is to help researchers,



# EU may fine political groups misusing personal data to skew elections

It's hoping to prevent another Cambridge Analytica scandal.

W  
Yale



Jon Fingas, @jonfingas  
08.27.18 in Internet

4  
Comments

275  
Shares



another record fine from the EU



Chris Smith @chris\_writes

August 17th, 2018 at 12:32 AM

Share

Tweet

Google has been dealt two huge blows in Europe in recent years, where antitrust investigations have ruled the company abused its position in search as well as in the mobile market. The company received two record fines as a result, which added up to [more than €6.74 billion \(\\$7.66 billion\)](#). On top of that, a third investigation is in the works and could bring over additional fines.

# What data?

Fake data

Biased data

Personal data

Mixed data

Reservation

About

Menu

Reviews

SIDES

Fruit Plate	\$7	S
Patatas Bravas, Spicy-Tangy Sauce and Rosemary Aioli	\$9	H

Powered by  singleplatform from Constant Contact | Owner Verified



SUNLIGHT  
FOUNDATION

LOGIN

search 

Follow Us



BLOG

TOOLS

APIS

POLICY

ISSUES

PRESS

ABOUT

CONTACT

DONATE

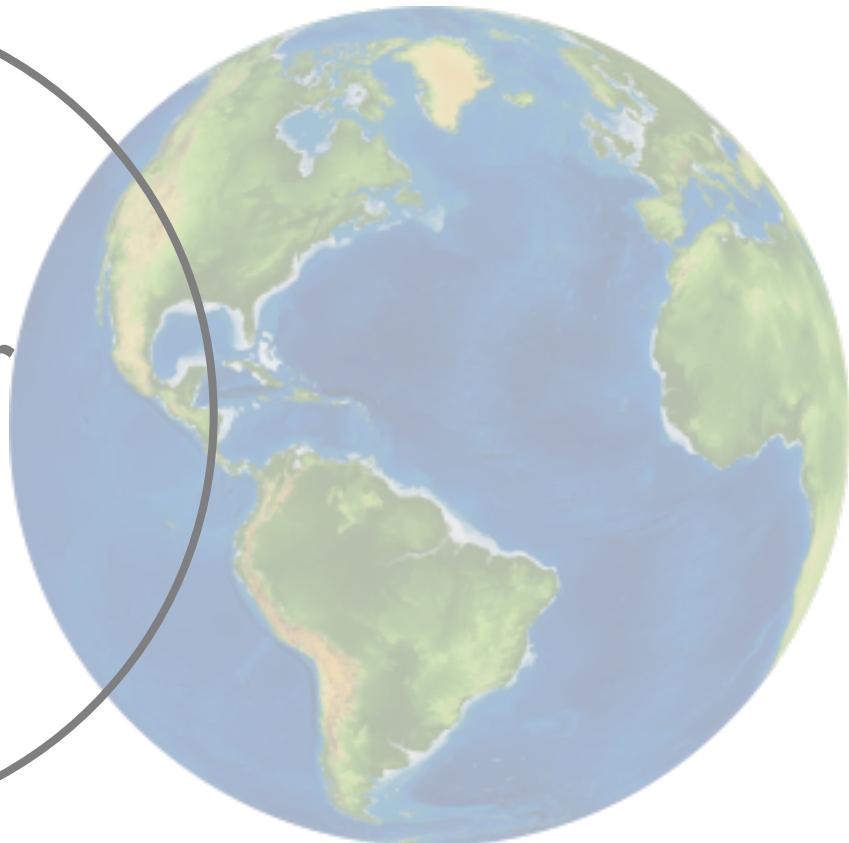
JOIN

Making government & politics more accountable & transparent.

Computer  
Science



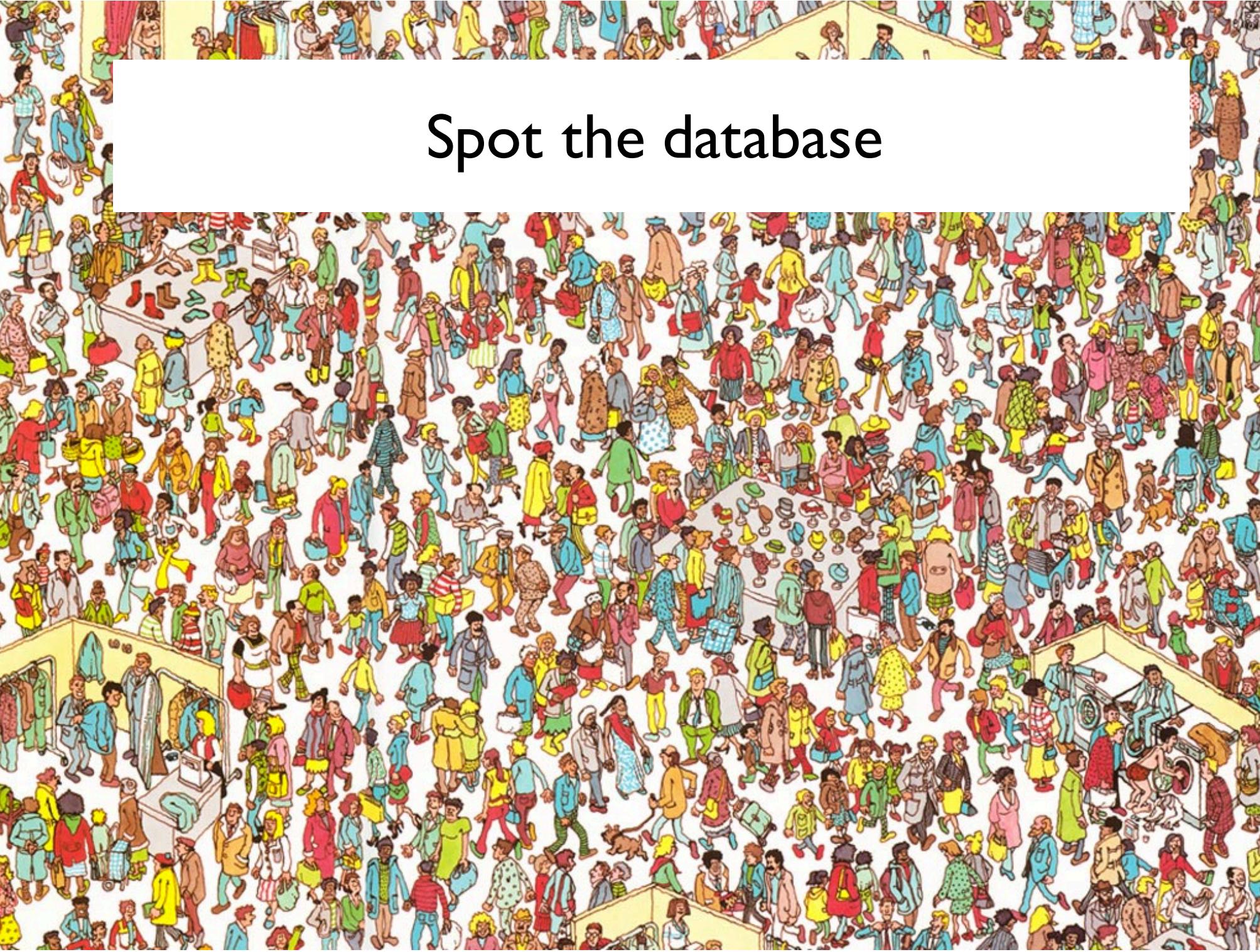
Computer  
Science





Data will be crucial to  
how we live  
as individuals and as a society

Data ~~will be~~ **is** crucial to  
how we live  
as individuals and as a society



# Spot the database



•••• AT&T 3:00 PM

Contacts +

Search

A

Apple Inc.

C

Call Recorder

F

Julia Fillory

Mike Fillory me

G

Justin Gilmore

Thomas Gilmore

Willa Good

H

Barry T. Hubbard

M

Favorites

Recent

Contacts

Keypad

Voiceicemail

A  
B  
C  
D  
E  
F  
G  
H  
I  
J  
K  
L  
M  
N  
O  
P  
Q  
R  
S  
T  
U  
V  
W  
X  
Y  
Z  
#



Home Notifications Messages

Search Twitter

What's happening?

sirrice retweeted

**IOC MEDIA** @iocmedia · Aug 2

Congratulations to the World Flying Disc Federation (WFDF), which was granted full IOC recognition at the #128IOCSession today!

519 317

**Eliran Sapir** @eliransapir · Jul 30

c-span.org/video/?327380-...

**Fred Werner** @SustainableFred · Jul 28

@berkeleyside the sun put on a show over Berkeley BEFORE sunset today



Trends · Change

**#GOPDebate**

Five Things To Watch For, While Watching The GOP Debate

64.7K Tweets about this trend

**#Ashes2015**

Australia 60 all out: Stats and facts that will leave you bamboozled

122K Tweets about this trend

2012-01-04 00:01:23,180 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: Receiving  
010

2012-01-04 00:01:23,184 INFO org.apache.hadoop.hdfs.server.datanode.DataNode.clienttrace  
cliID: DFSClient\_-603743753, offset: 0, srvID: DS-292194659-127.0.1.1-50010-13247633001

2012-01-04 00:01:23,185 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: PacketRespon

2012-01-04 00:01:23,291 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: Receiving  
10 Is this a Database?

2012-01-04 00:01:23,293 INFO org.apache.hadoop.hdfs.server.datanode.DataNode.clienttrace  
cliID: DFSClient\_-603743753, offset: 0, srvID: DS-292194659-127.0.1.1-50010-132476330017

2012-01-04 00:01:23,293 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: PacketRespon

2012-01-04 00:01:23,324 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: Receiving  
010

2012-01-04 00:01:23,326 INFO org.apache.hadoop.hdfs.server.datanode.DataNode.clienttrace  
cliID: DFSClient\_-603743753, offset: 0, srvID: DS-292194659-127.0.1.1-50010-1324763300176

2012-01-04 00:01:23,327 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: PacketRespon

2012-01-04 00:01:23,409 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: Receiving  
10

2012-01-04 00:01:23,411 INFO org.apache.hadoop.hdfs.server.datanode.DataNode.clienttrace  
, cliID: DFSClient\_-603743753, offset: 0, srvID: DS-292194659-127.0.1.1-50010-1324763300

2012-01-04 00:01:23,411 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: PacketRespon

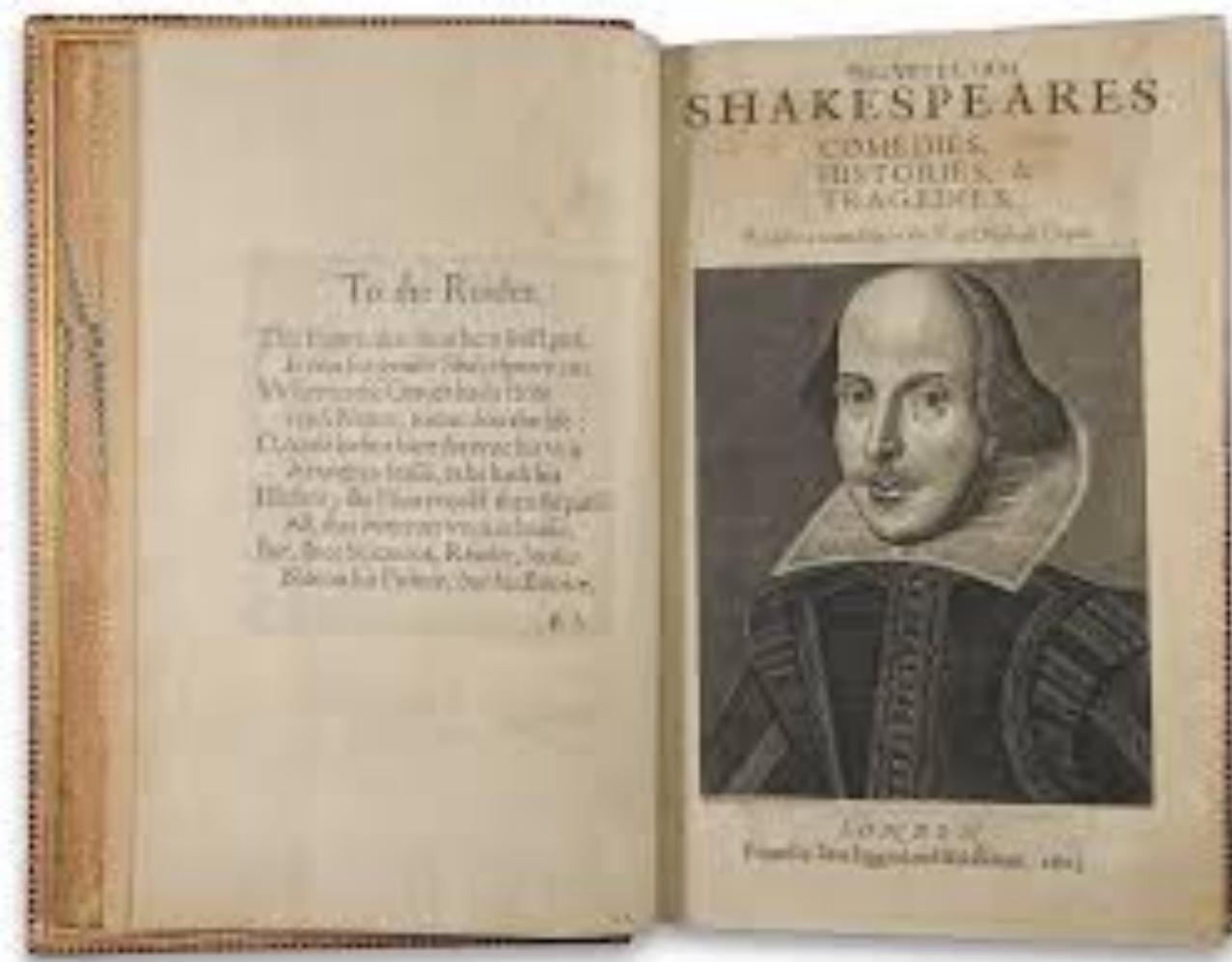
2012-01-04 00:01:23,433 INFO org.apache.hadoop.hdfs.server.datanode.DataNode.clienttrace  
cliID: DFSClient\_-2054881890, offset: 0, srvID: DS-292194659-127.0.1.1-50010-1324763300

2012-01-04 00:01:23,494 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: Receiving  
10

2012-01-04 00:01:23,498 INFO org.apache.hadoop.hdfs.server.datanode.DataNode.clienttrace  
, cliID: DFSClient\_-2054881890, offset: 0, srvID: DS-292194659-127.0.1.1-50010-1324763300

2012-01-04 00:01:23,498 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: PacketRespon

2012-01-04 00:01:23,523 INFO org.apache.hadoop.hdfs.server.datanode.DataNode: Receiving  
010



### To the Reader.

The humor of the better trifles,  
As fayre for fayre, like a glasse in a bosome,  
Wherinnesse Countes hauie to see  
Thei'res faces, fayre, as theye be,  
Cleane helpe haue for meuchas to a  
Prestyng fayre, as he had for  
Hilfe, to haue fayre, when he gaue  
All that he had, to such as he  
Fayre, fayre, fayre, fayre, fayre,  
Fayre, fayre, fayre, fayre, fayre,

6.5

J. O. P. B. 2. 2.  
Facsimile of the First Folio of 1623.

# What is a Database?

Structured data

# What is a Database?

Lots of  
Structured data

# Database Management System (DBMS)

A system to **store, manage** and **access** databases

# Database Management System (DBMS)

System to **safely** and **reliably** store **lots** of **persistent** structured data and is **convenient** for **multiple** users to **efficiently** access and modify.

# Is a script a DBMS?

Javascript/Python Script

Data stored in variables (RAM)

Very fast access

# Is Excel a DBMS?

Microsoft office security

Visually access/modify/compute over data cells

Click save to store persistently

# Is the file system a DBMS?

Manages files that are persistently stored on disk

Open/read/seek/write access to files

Access via file names

Access control via permissions

# Is the file system a DBMS?

You and a friend edit the same text file

Save at the same time

What happens?

1. Your changes survive
2. Friend's changes survive
3. Both changes survive
4. No changes survive
5.  $\neg \backslash (\cup) \backslash$

# Is the file system a DBMS?

You edit a text file

Computer crashes

What happens?

1. All changes survive
2. No changes survive
3. Changes from last save survive
4.  $\neg \backslash (\tau) \backslash$

# Is the file system a DBMS?

A screenshot of a Microsoft PowerPoint slide. The title "Is the file system a DBMS?" is displayed in large black font at the top. Below the title, the slide content area contains the text "COMS W4111" and "Introduction to Databases". The Microsoft ribbon is visible at the top, showing tabs for Animations, Slide Show, Review, Insert, Format, and Slide Show. A red box highlights the "Recovered File 2" tab in the ribbon. A yellow callout bubble is positioned over the "Insert" tab, containing the text "Insert an item such as a text box, slide number, hyperlink, or symbol".

Recovered File 2

Animations    Slide Show    Review    Insert    Format    Slide Show

Paragraph    Insert    Format    Slide Show

Insert an item such as a text box, slide number, hyperlink, or symbol

COMS W4111

Introduction to Databases

# Want Guarantees from DBMS

You want to write a hot new app on a DBMS.  
What do you *not* want to worry about?

Failures disk, machine, human, corruption, deity

Lots of users

Ad-hoc data access

Data formats csv? tsv? custom format?

# Database Management System (DBMS)

System to **safely** and **reliably** store **lots** of **persistent** structured data and is **convenient** for **multiple** users to **efficiently** access and modify.

# Database Management System (DBMS)

<b>Safe</b>	Consistent and correct data after failures
<b>Reliable</b>	99.99+% Uptime
<b>Lots</b>	>>RAM (terabytes)
<b>Persistent</b>	Lives longer than DBMS application
<b>Convenient</b>	Physical Independence. Declarative.
<b>Multiple Users</b>	Concurrent access. Access control.
<b>Efficient</b>	<i>Fast:</i> 100k+ queries / sec

# DBMSes in the Wild

## Classic Disk-based Relational

\$\$: Oracle, IBM, Microsoft, Teradata, EMC, etc

Free: MySQL, PostgreSQL, SQLite

## New Relational

In-Memory, Column-store, Streaming

## Non-traditional

Search (Google, Bing, Lucene), Scientific, Geo, Graph

## NoSQL

Big Data: Hadoop, Spark, etc

Key-value: Mongo, BerkeleyDB, Cassandra, etc

## DBMS-as-a-Service

MS Azure, Google BigQuery, Amazon Redshift/RDS ...

# Encompasses most of CS

OS	DBMS directly manages hardware
Languages	SQL is a domain specific language
Theory	Algorithms, models, NP-complete
AI/ML	Knowledge Discovery, KDD
Logic	Relational Algebra = 1 <sup>st</sup> order logic

## Scalable Computer Science

# Good time to learn!

Cloud programmer

Data science

Data engineer

Machine learning engineer



DATA

## Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

# 2 Key Concepts

Data Independence

Declarative Languages

Serve to insulate application programmers  
from the system implementation

# Data Independence

<b>External Schema</b>	Describe how users see data
<b>Conceptual Schema</b>	Describes logical structure
<b>Physical Schema</b>	Describes files and indexes

External Schema

Conceptual Schema

Physical Schema

“Data”

# Example App: Guuber

Users(**uid int**, name str, age int)

Drivers(**did int**, name str)

Rides(**uid int**, **did int**, distance float, drive\_time float)



# Data Independence

UID	Name	Age
0	Eugene	17
1	Luis	20
2	Ken	30

0,Eugene,17  
1,Luis,20  
2,Ken,30  
CSV File

What is the number of adults?

# Data Independence

UID	Name	Age
0	Eugene	17
1	Luis	20
2	Ken	30

0,Eugene,17  
1,Luis,20  
2,Ken,30  
CSV File

```
n = 0
for line in csv_file:
    attributes = line.split(",")
    if attributes[2] >= 18:
        n += 1
```

# Data Independence

UID	Name	Age
0	Eugene	17
1	Luis	20
2	Ken	30

0,1,2  
Eugene,Luis,Ken  
17,20,30  
CSV File

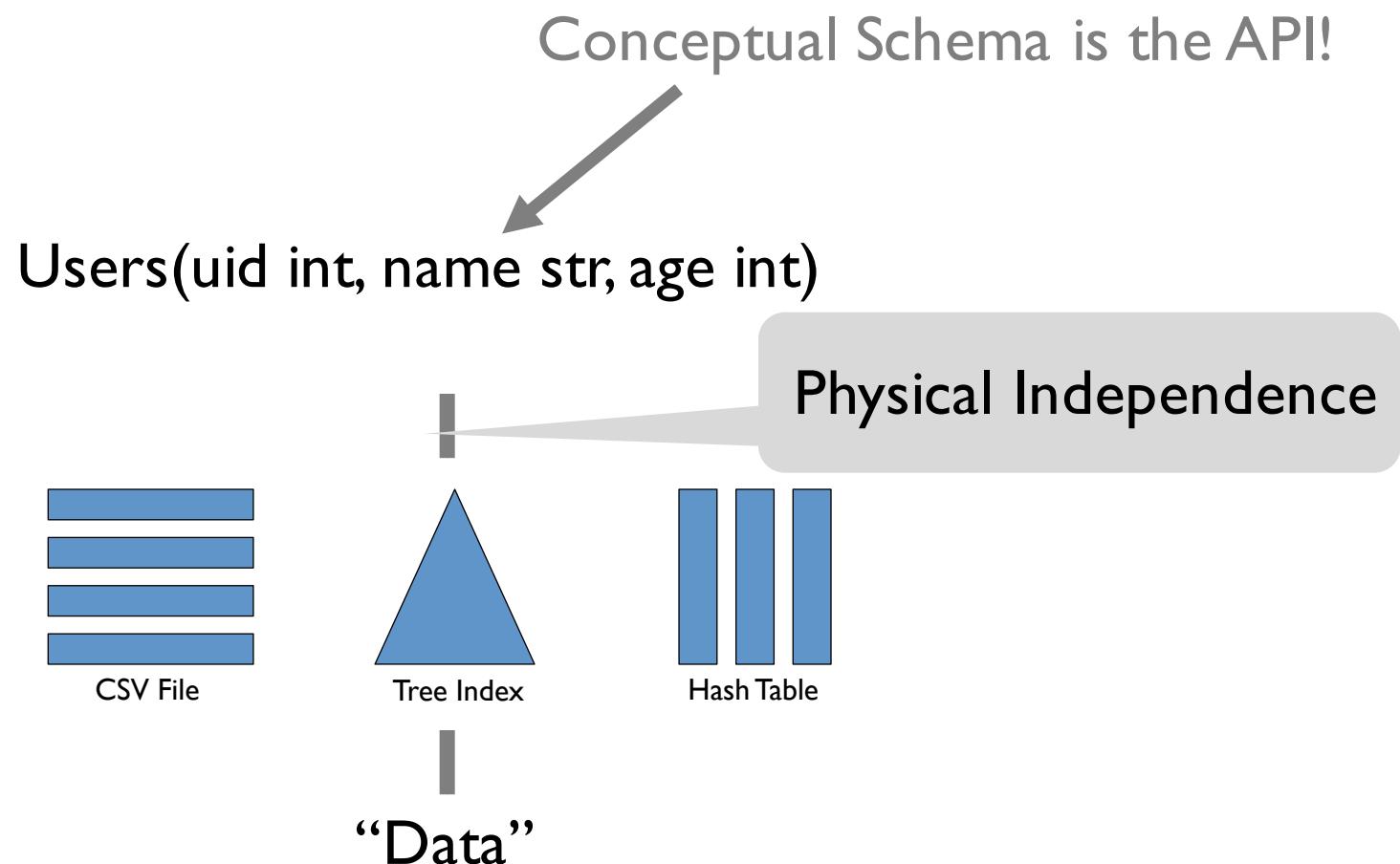
~~n = 0  
For line in csv\_file:  
    attributes = line.split(",")  
    if attributes[2] >= 18:  
        n += 1~~

# Data Independence

**Conceptual Schema**

Describes logical structure

**Physical Schema**  
Describes files and indexes



# Data Independence

Users(uid int, name str, age int)

Drivers(did int, name str)

Rides(uid int, did int, distance float, drive\_time float)

“Welcome back Mr. Wu”

# Data Independence

Users(uid int, **fname str, lname str**, age int)

Drivers(did int, name str)

Rides(uid int, did int, distance float, drive\_time float)

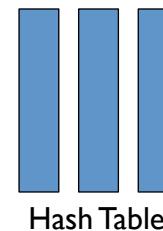
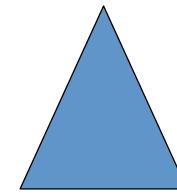
“Welcome back Mr. Wu”

# Data Independence

**Conceptual Schema**  
Describes logical structure

`Users(uid int, name str, age int)`

**Physical Schema**  
Describes files and indexes



**Physical Independence**

“Data”

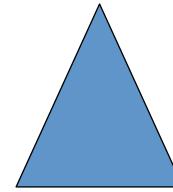
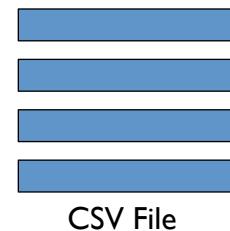
# Data Independence

**Conceptual Schema**  
Describes logical structure

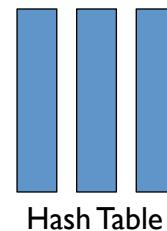
`Users(uid int, fname str, lname str, age int)`

Physical Independence

**Physical Schema**  
Describes files and indexes



Tree Index



Hash Table

“Data”

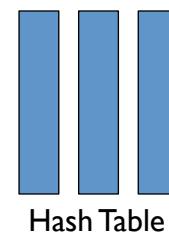
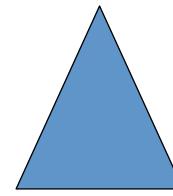
# Data Independence

External Schema	View 1	View 2	View 3
Describe how users see data			
Conceptual Schema			
Describes logical structure			

Users(uid int, **fname str, lname str, age int**)

Logical Independence

Physical Independence



“Data”

# Data Independence

## Physical Independence

Protection from changes in physical structure of data

## Logical Independence

Protection from changes in logical structure of data

**One of most important properties of a DBMS**

# Declarative

**What you want,**

“Make me a sandwich”

Buy from pb&j store

Make BLT

½ Tuna

Veggie

**not how to do it.**

“Take two slices of wheat bread out of the 2<sup>nd</sup> shelf, put them next to each other...”

What if on 1<sup>st</sup> shelf?  
Out of wheat bread?  
No counter space?

# Declarative

“I want all highly rated fast drivers”

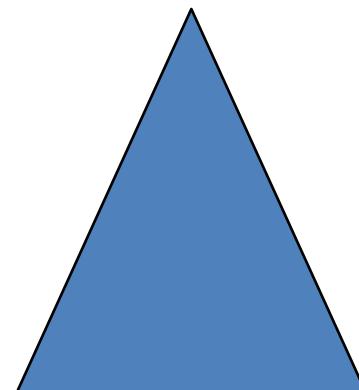
---

DBMS

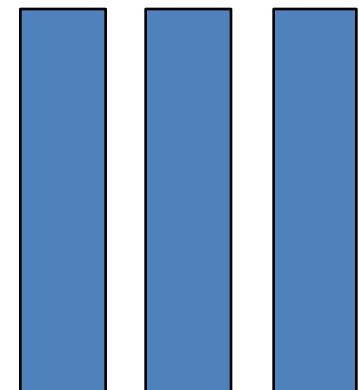
---



CSV File



Tree Index



Hash Table

# Declarative

“I want all highly rated fast drivers”

---

DBMS

---

Node

Node

Node

# Declarative

“I want all highly rated fast drivers”

---

DBMS

---

Node

# Classic Components in Databases

Concurrency Control

Transactions

Atomicity

Recovery and Logs

# Transaction: Execution of a DB Program

Def: *atomic* sequence of DBMS actions

```
Begin;  
<read beth's account>  
<deduct from beth's account>  
<increase eugene's account>  
Commit; (or Abort;)
```

# Transaction: Execution of a DB Program

Def: *atomic* sequence of DBMS actions

Each fully executed transaction must leave DB in  
*consistent state* if DB is consistent before transaction

- Users specify simple *integrity constraints* on data, and DBMS enforces the constraints.
- DBMS does not understand semantics of its data  
e.g., doesn't know how bank interest is computed
- User's responsibility to ensure transaction (run alone) preserves consistency

# **Concurrency Control**

**Concurrently running multiple user programs needed for good performance**

Disk accesses are frequent & slow. Keep CPU working on several user programs while waiting.

**Concurrency can cause inconsistencies**

- e.g., check cleared while account balance being computed.
- *Really* hard to program against

**DBMS ensures such problems don't arise**

- programmers can pretend to use a single-user system.

# Scheduling Concurrent Transactions

Transactions  $T_1, \dots, T_n$  are run concurrently  
Equivalent to a *serial* ordering (as if no concurrency)

**Locks:**  $T_i$  requests and waits for lock before read/write.

e.g.,  $T_i$  locks the database, updates, then releases

e.g.,  $T_i$  locks the table, updates, then releases

e.g.,  $T_i$  locks rows, updates, then releases

Will talk about how this works later in course.

# Atomicity

Def: Xact fully completes, or never happened  
even after failures e.g., crashes

Record all actions Xact did during execution in a log

1. **Write ahead logging**: before making any change, ensure the change is safely recorded in log
2. After failure, read log and undo any incomplete Xacts

# The Log

A log record contains enough info to undo actions:

- Transaction id

- T<sub>i</sub> writes an object: old and new values

- Log record *must* be safely stored before the changed data

- T<sub>i</sub> commits/aborts: store commit/abort action

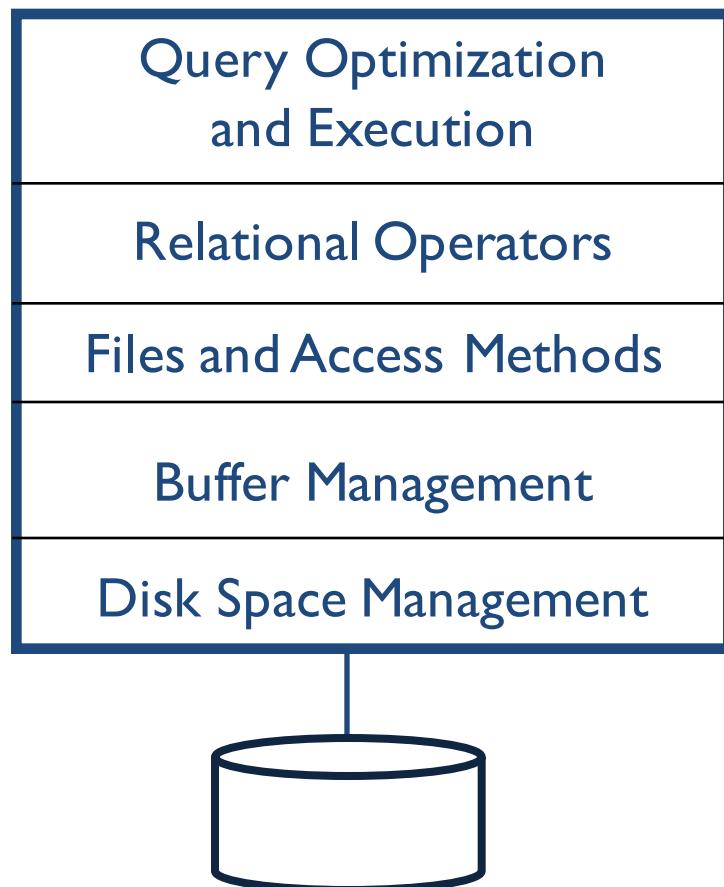
All logging, recovery and concurrency control activities hidden away from user.

# Classic Structure of a DBMS

Typical layered architecture  
DBMS, not OS, manages  
memory and disk

Doesn't show concurrency  
control & recovery components

These layers  
must consider  
concurrency  
control and  
recovery



# Database Courses at Columbia

# **COMS W4111 - Intro to Databases**

Prerequisites: CS3137 or CS3134; fluency in Python

Intro to DBMSes

Data Models Entity-relation, Relational, ...

Relational Algebra

SQL

Applications + SQL cursors, APIs, embedded ...

Normalization

Peek at DBMS internals:

Storage and indexing

Query optimization

Transaction Processing

# **COMS W4112-Database Sys. Impl.**

Prerequisites: CS3137 or CS3134; fluency in Python

Storage Methods and Indexing

Query Processing and Optimization for 1NF Relations,  
including external sorting

Materialized Views and Use in Query Optimization

Query Processing and Optimization for ORDBMSs

Transaction Processing and Recovery

Parallel & Distributed DBMSes: Query Proc. and  
Optimization

Parallel and Distributed Databases: Transaction Processing

Performance Considerations Beyond I/Os

# **COMS E6111-Advanced Databases**

Prerequisites: CS4111; fluency in Java or Python

Information Retrieval

Web Search

Distributed Information Retrieval and Web Search

Data Mining

Data Warehousing, OLAP, Decision Support

Information Extraction

Scalable Visualization and Interaction

Supporting data analysis

Exploration, explanation and exhibition techniques

# **COMS E6xxx-DB Research Seminars**

Prerequisites: CS4111; fluency in Java or Python

**6113 Database Research Topics**

**6998.002 Interactive Data  
Exploration Systems**

**6998.005 Database Topics in  
Research & Practice**

# Administrivia

# Next Up

HW0 is out.

Due by Friday 1/25 10AM sharp.

No Late Submission Accepted  
0 in class if not submitted in time

# Your Instructor: **Eugene Wu**

B.S. @U.C. Berkeley

Ph.D. @MIT

PostDoc @U.C. Berkeley

Assistant Professor since Fall 2015

Databases, visualization, data analysis  
data cleaning, crowdsourcing.

# Your Instructor: **Eugene Wu**

## Contact

[www.eugenewu.net](http://www.eugenewu.net)

[ewu@cs.columbia.edu](mailto:ewu@cs.columbia.edu)

421 Mudd

## Office hours

Thurs 3-4PM

By appointment by email

# Class Resources

Class web page  
[w4111.github.io](https://w4111.github.io)

Discussion board  
piazza (linked from website, public)

Announcements from class staff:  
Website

# Your TAs

Amita Shukla

Yiru Chen

Zhicheng Wu

Ziao Wang

All TA office hours in CS TA Room  
TA office hours will be set next week  
(see class web page)

# Class Information: Prerequisites

COMS W3134 - *Data Structures in Java* or  
COMS W3137 - *Data Structures and Algorithms*  
(equivalent courses taken elsewhere are acceptable as well)

Fluency in **Python**

# Class Information: Lectures

Tuesdays and Thursdays

8:30-10AM

209 Havemeyer

(here)

# An aside: Success

What does succeeding in this course mean?

Timescales

How to encourage a collaborative environment?

What discourages it?

Assessment

# Grading Information

Midterm I: 25%

Midterm 2: 40% cumulative

HW: 15% (4 HWs equally weighed)

Project I: 15%

Project 2: 5%

Extra credit: scribe notes + advanced assignments

Median grade: B or slightly higher.

Alternative or make-up exams will not be given.

All homework assignments are equally weighted.

Project I has higher weight than Project 2.

# Exam Dates

Midterm I: 3/7, in class

Midterm 2: 4/30 last day of class, in class, cumulative

Makeup exams are not scheduled

# Homework

Homeworks usually due at 10AM of due date.

Assignment will specify submission instructions.

No extensions or exceptions.

Three grace late days for hws throughout the semester.

After using all grace days, 25% grade deduction per late day.

Check full details on web site.

# Projects (more details soon)

Two projects.

Teams of two

Run on cloud infrastructure

Get CS account if your team doesn't have a computer

Language is Python

Project 1

Model and build your own database web application

Explore “traditional” relational database features.

Project 2

TBD

# Projects (cont.)

3 grace late days total for project parts I and 2.

No extensions or exceptions for project part 3 submission.

After using all grace days, 25% grade deduction per late day.

Check full details on web site.

# Extra Credit

Added *after* the curve

Scribe notes: 0-5% extra credit

Advanced Assignments:

- ~same value as HW
- Goes into more depth
- Hack on the DataBass system

# Collaboration Policy

Read Syllabus on course site for allowed conduct

CS Dept academic honesty policies

<http://www.cs.columbia.edu/education/honesty>

We will not tolerate *any* cheating

# Collaboration Policy

Discussing lectures and course material strongly encouraged

Homework and exams are *individual*. No exceptions  
Any libraries or code however minor must be disclosed.

Projects are done in *teams*; no collaboration between teams.

Contact the instructor right away if you have any questions  
or are falling behind.

# Textbook

Raghu Ramakrishnan, Johannes Gehrke: *Database Management Systems*, 3<sup>rd</sup> edition, McGraw-Hill, 2002

*Available from*

*Bookculture bookstore 536 W. 112th St.*

*Online retailers*

*Upperclass-persons*

*On reserve in Engineering Library*

# Contests and Rewards

## Project I contest

Four best projects chosen as contest winners.

Winners get:

10% boost in Project I grade.

If time allows: demo your project in class.

# On-going Feedback

Please provide feedback throughout the course.

- What is useful or confusing in lecture
- Thoughts about software stack
- Thoughts about assignments

Email me, come to office hours, talk to staff or:

# On-going Feedback

Use form on website

The image shows a feedback form titled "Feedback form" with a light gray background. At the top, it says "Please share your comments and suggestions for the course!" Below that, there is a note in red text: "\* Required". A section labeled "Feedback \*" asks for comments about what worked or was confusing/difficult. Another section labeled "Improvements" asks for suggestions to improve things.

**Feedback \***  
Share what worked or what was confusing/difficult

**Improvements**  
What change would you suggest to improve things?

Slides borrow material from  
Prof. Gravano

Prof. Hellerstein & Franklin@Cal

Prof. Madden & Stonebraker@MIT

(and by transitivity Raghu Ramakrishnan and Johannes Gehrke)

w4111.github.io  
ewu@cs.columbia.edu

**DO HOMEWORK 0!**