

W4121  
Computer Systems for Data Science  
Spring 2018

Roxana Geambasu, Sambit Sahu, Eugene Wu

<https://w4121.github.io/>

1

Data

2

Data

is for serious business

3

Data

is at the center of most things.

4

Data

is at the center of *everything*

5

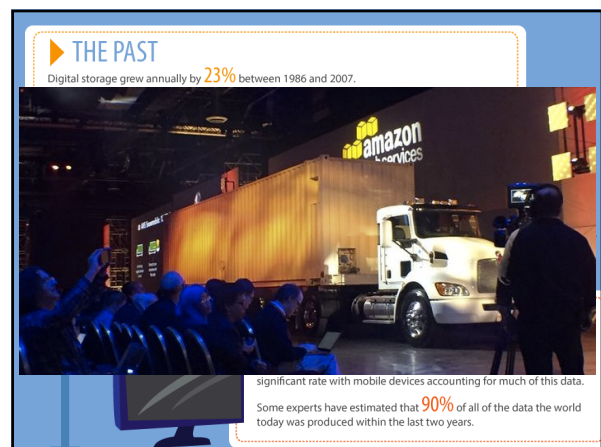
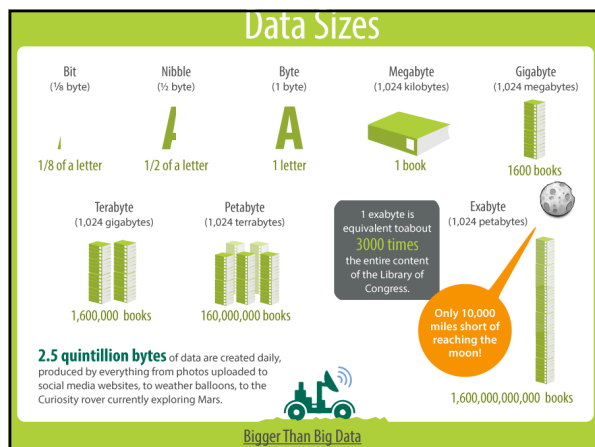
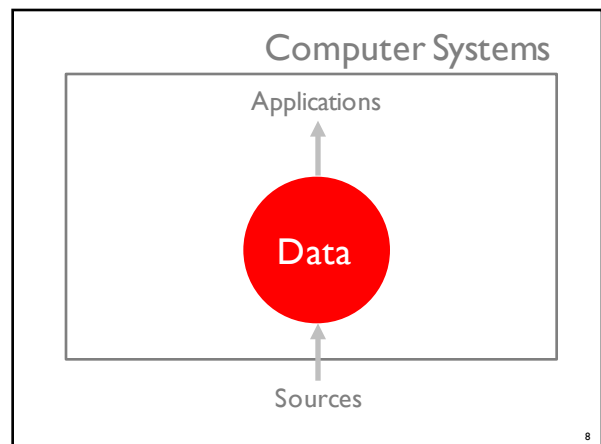
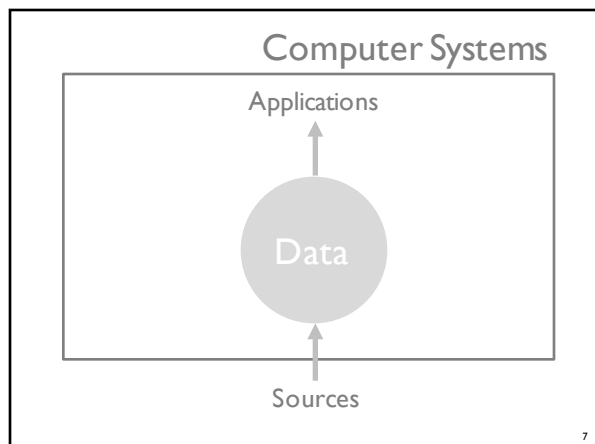
Computer Systems

How is it used?

Data

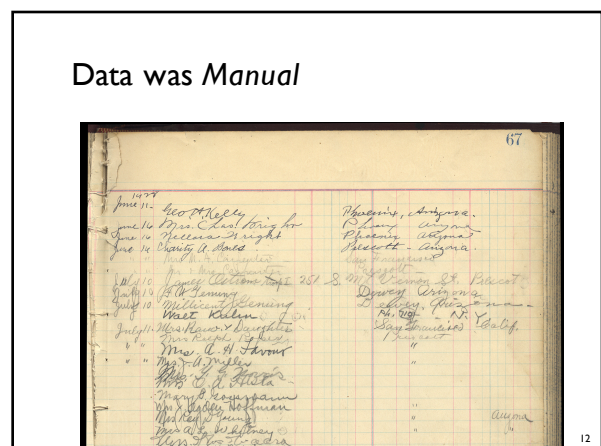
Where does it come from?

6



### How did we get here?

11



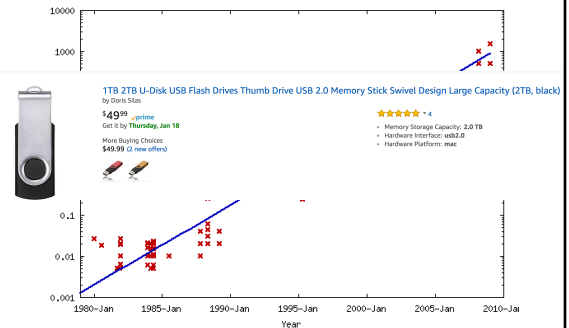


## Data was *Expensive*



13

## Data is *Cheap*



14

## Data is *Automated*

Physical devices



15

## Data is *Automated*

Physical devices  
Software logs

16

## Data is *Ubiquitous*

Physical devices  
Software logs  
Phones



17

## Data is *Ubiquitous*

Physical devices  
Software logs  
Phones  
GPS/Cars



18

## Data is Everywhere

Physical devices  
Software logs  
Phones  
GPS/Cars  
Internet of Things



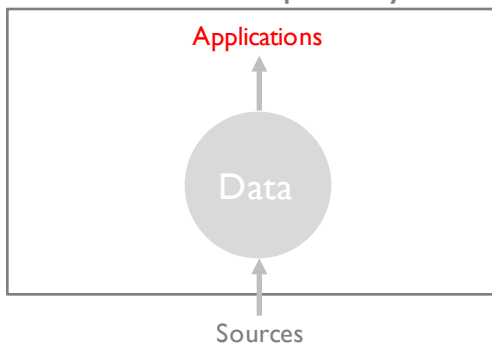
19

## Data is Temporal

```
cmd: dyn-128-59-158-173 /root/.ssh/ (ms) [root@ms] # 38
Jan 18 22:32:01 dyn-160-39-38-195 WindowServer[185] -darnings: send_datagram_available_ping: pid 271 failed to act on a ping
disappeared before timing out.
... Last message received 5 times ...
Jan 18 22:32:13 dyn-160-39-38-195 Papers[64391] -dnotice: Downloading patches:
Jan 18 22:32:13 dyn-160-39-38-195 Papers[64391] -darnings: Finished sync: 2/28/16, 10:33 PM
Jan 18 22:32:13 dyn-160-39-38-195 Papers[64391] -darnings: SYNCING: We will sync again in 120 seconds
Jan 18 22:32:14 dyn-160-39-38-195 Papers[64391] -darnings: [WARNING] DropboxSD: error making request to /2/Files/crunch.
... CWD) Cannot create folder 'CompressedCheckpoints' because a file or folder already exists at path '/papers/Library/and
CompressedCheckpoints'
Jan 18 22:32:22 dyn-160-39-38-195 WindowServer[185] -darnings: send_datagram_available_ping: pid 271 failed to act on a ping
disappeared before timing out.
Jan 18 22:32:22 dyn-160-39-38-195 kernel[0] -dnotice: Google Chrome He[64477] triggered unmet of range 0x7ffff8000000-0x7f
000000 of D1D shared region in VM map 0x42f17d93b8b0d4e9. While not abnormal for debuggers, this increases system memory fo
ce until the target exits.
Jan 18 22:32:24 dyn-160-39-38-195 kernel[0] -dnotice: Google Chrome He[64478] triggered unmet of range 0x7ffff8000000-0x7f
000000 of D1D shared region in VM map 0x42f17d93c0450839. While not abnormal for debuggers, this increases system memory fo
ce until the target exits.
Jan 18 22:32:35 dyn-160-39-38-195 WindowServer[185] -darnings: send_datagram_available_ping: pid 271 failed to act on a ping
disappeared before timing out.
Jan 18 22:32:38 dyn-160-39-38-195 kernel[0] -dnotice: Google Chrome He[64482] triggered unmet of range 0x7ffff8000000-0x7f
000000 of D1D shared region in VM map 0x42f17d93b36c2c39. While not abnormal for debuggers, this increases system memory fo
ce until the target exits.
Jan 18 22:32:43 dyn-160-39-38-195 kernel[0] -dnotice: Google Chrome He[64484] triggered unmet of range 0x7ffff8000000-0x7f
000000 of D1D shared region in VM map 0x42f17d93b36c2c39. While not abnormal for debuggers, this increases system memory fo
ce until the target exits.
Jan 18 22:32:45 dyn-160-39-38-195 kernel[0] -dnotice: Google Chrome He[64486] triggered unmet of range 0x7ffff8000000-0x7f
000000 of D1D shared region in VM map 0x42f17d93b36c2c39. While not abnormal for debuggers, this increases system memory fo
```

20

## Computer Systems



21

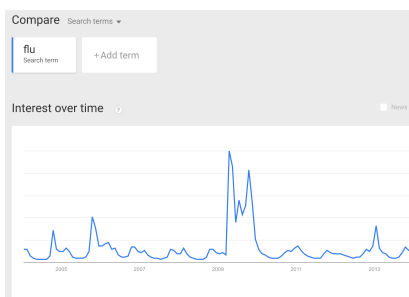
## Data Science Applications

Health



## Data Science Applications

Health



<https://www.google.org/flutrends/>

23

## Data Science Applications

Health

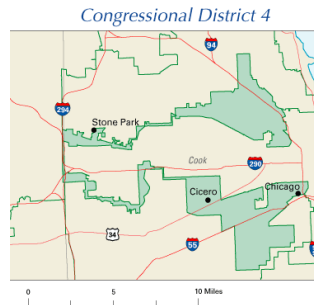
### Thank you for stopping by.

Google Flu Trends and Google Dengue Trends are **no longer publishing** current estimates of Flu and Dengue fever based on search patterns. The historic estimates produced by Google Flu Trends and Google Dengue Trends are available below. It is still early days for

<https://www.google.org/flutrends/>

## What are we doing with data?

Health  
Politics

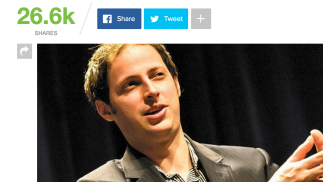


25

## What are we doing with data?

Health  
Politics

Triumph of the Nerds: Nate Silver Wins in 50 States



26

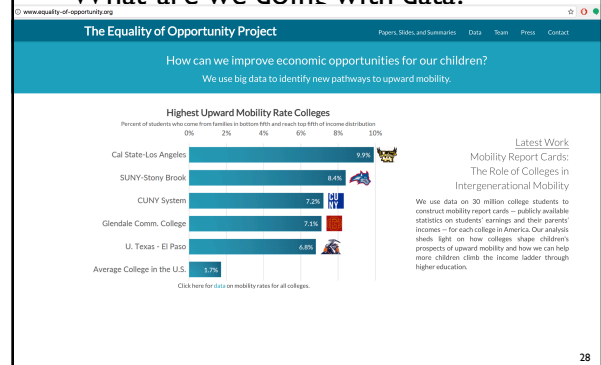
## What are we doing with data?

Health  
Politics



27

## What are we doing with data?



28

## What are we doing with data?

Health  
Politics  
Investigative Journalism  
Society



<https://linear.gorcs.com/2017/12/23/hamath-palihapitiya-on-facebook-quest-games-quest/>

29

## What are we doing with data?

### Epidemiological modeling of online social network dynamics

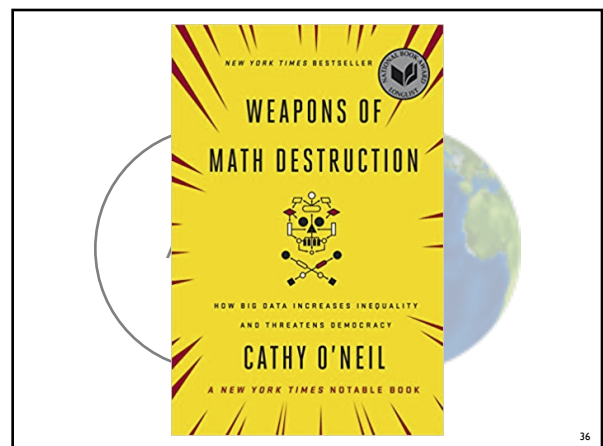
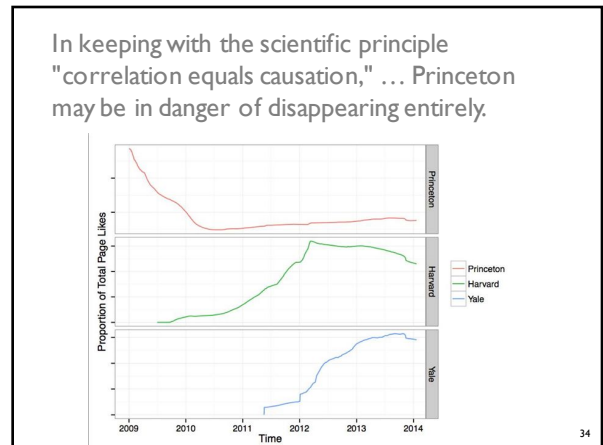
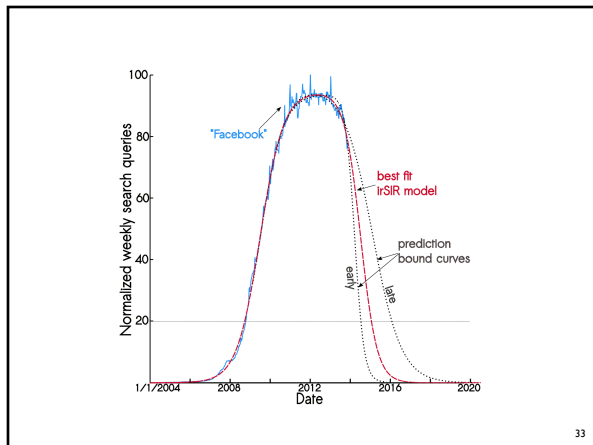
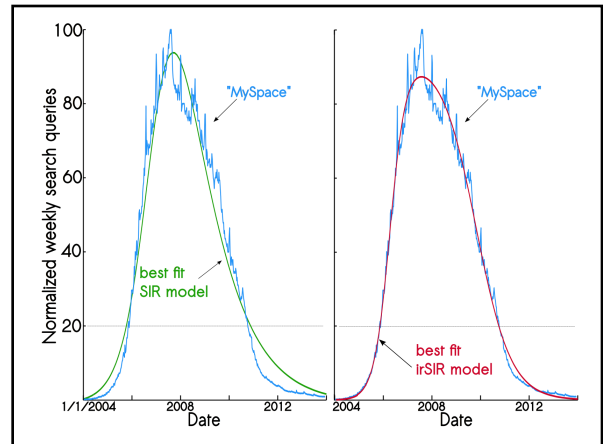
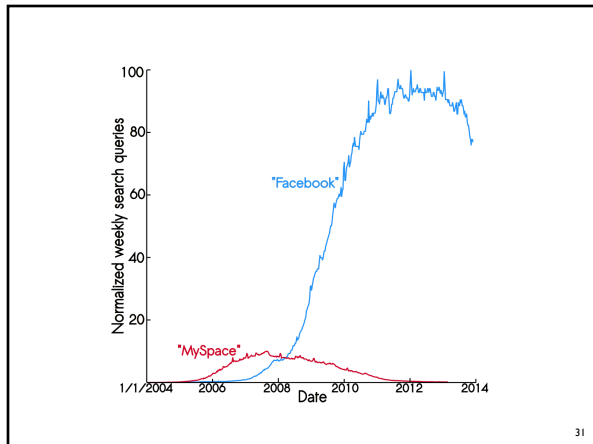
John Cammarella<sup>1</sup>, Joshua A. Spechler<sup>1\*</sup>  
<sup>1</sup> Department of Mechanical and Aerospace Engineering, Princeton University, Princeton, NJ, USA  
 \* E-mail: Corresponding spechler@princeton.edu

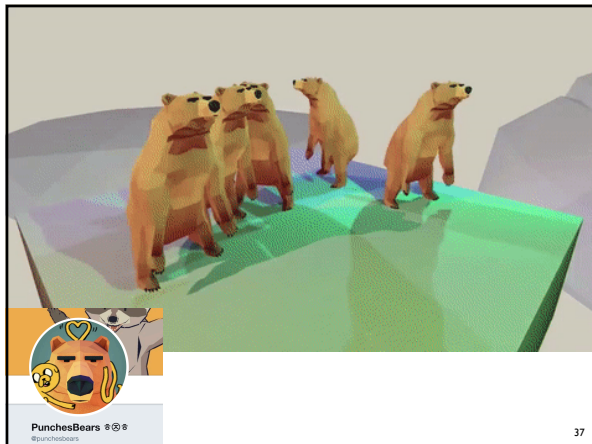
### Abstract

The last decade has seen the rise of immense online social networks (OSNs) such as MySpace and Facebook. In this paper we use epidemiological models to explain user adoption and abandonment of OSNs, where adoption is analogous to infection and abandonment is analogous to recovery. We modify the traditional SIR model of disease spread by incorporating infectious recovery dynamics such that contact between a recovered and infected member of the population is required for recovery. The proposed infectious recovery SIR model (iSIR model) is validated using publicly available Google search query data for "MySpace" as a case study of an OSN that has exhibited both adoption and abandonment phases. The iSIR model is then applied to search query data for "Facebook," which is just beginning to show the onset of an abandonment phase. Extrapolating the best fit model into the future predicts a rapid decline in Facebook activity in the next few years.

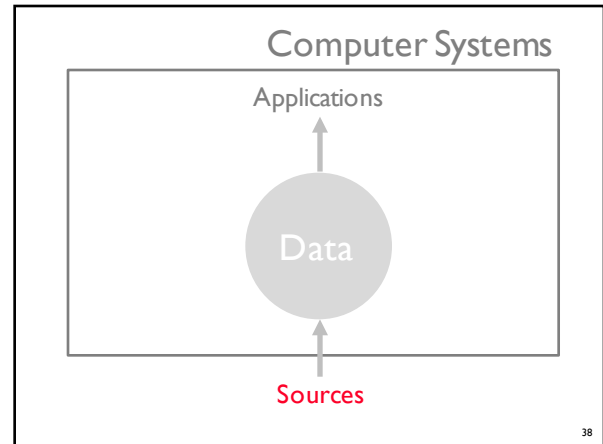
Extrapolating the best fit model into the future predicts a rapid decline in Facebook activity in the next few years

30





37



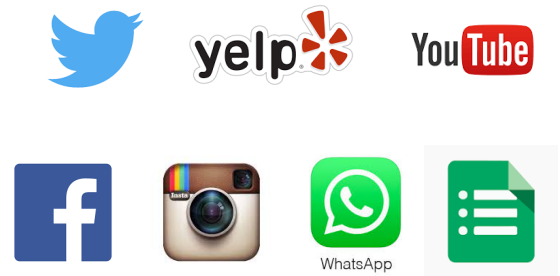
38

## Autogenerated – record every...

- Mouse click
- Car drive
- Ad impression
- Webpage visit
- Billing transaction
- Network message
- Error
- Video stream

39

## User generated



40

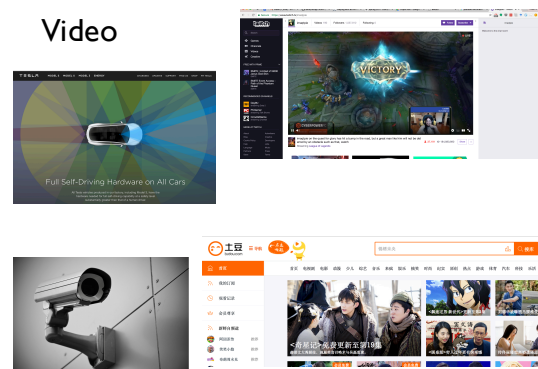
Then there's the Q50, a smart watch for children. Marketed as a way to help parents easily communicate with and keep track of their kids, bugs in the watch would allow hackers to "intercept all communications, remotely listen to the child's surroundings and spoof the child's location," according to [a report](#) by Top10VPN, a consumer research company this month.

And the BB-8 droid, which was released with "The Last Jedi" this month, also had an insecure

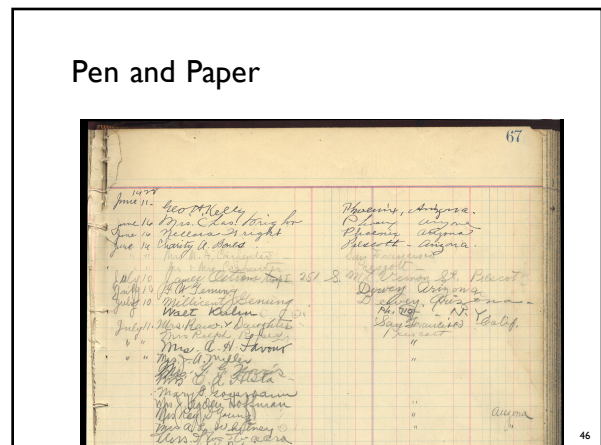
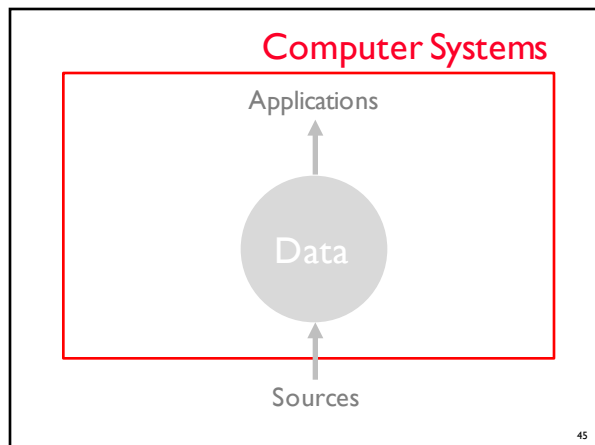
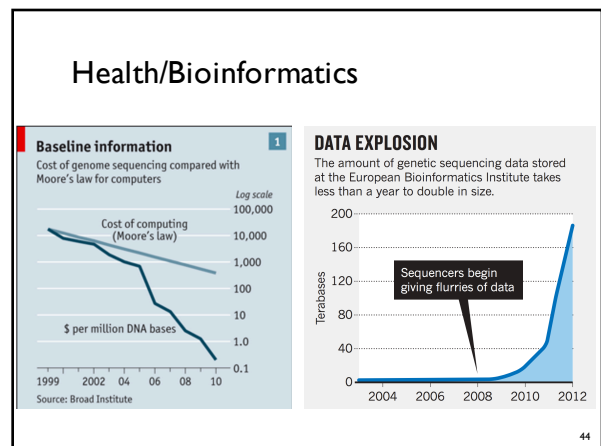
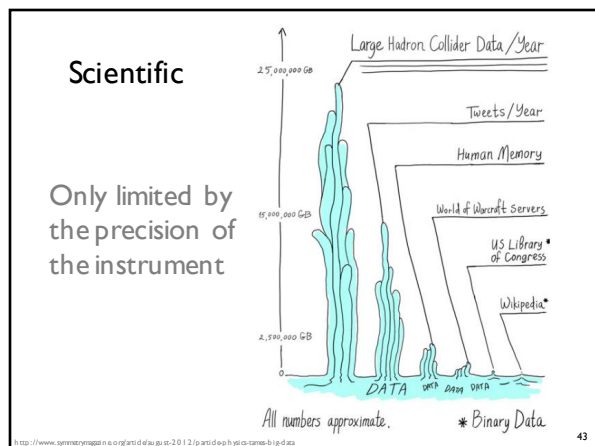
<https://twitter.com/internetofshitfangmen>

41

## Video



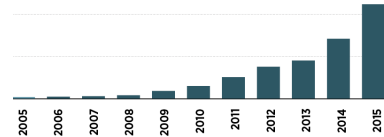
42







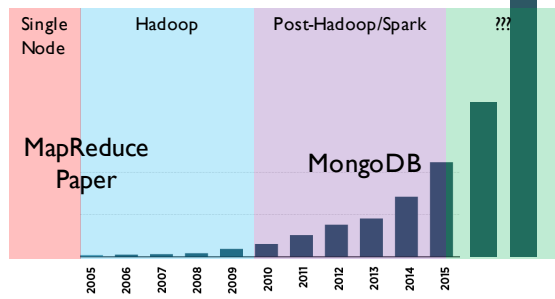
## Data Volume Over Time



Note: Post-2013 figures are predicted. Source: UNECE

50

## Big Data Systems Over Time



Note: Post-2013 figures are predicted. Source: UNECE

51

## Post-Spark: ML Systems?

### Lines of code in google's ML system

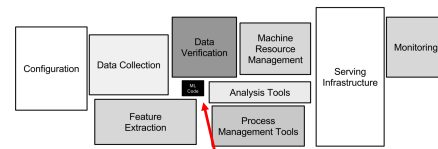


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

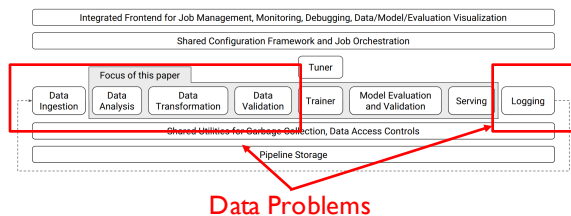
**ML Training**

Hidden technical debt in machine learning systems

52

## Post-Spark: ML Systems?

### TFX: Google's TensorFlow ML Platform



**Data Problems**

TFX

53

## Post-Spark: ML Systems?

### Massive data management to support ML

#### Many data problems

collection, cleaning, merging validation, analysis, monitoring processing finding, versioning sharing

54

## Course Goal

Understand fundamental principles behind large-scale data science systems

Data Processing Techniques for “big” data

Experience with modern data science tools (Spark)

55

data

57



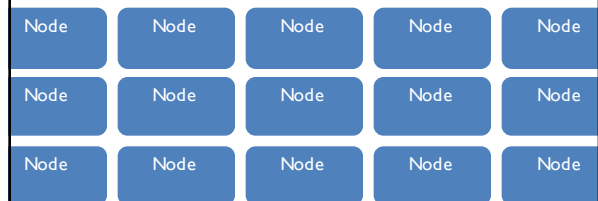
58



59



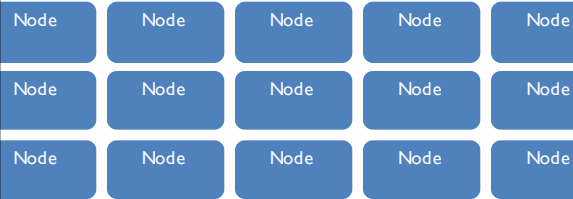
60





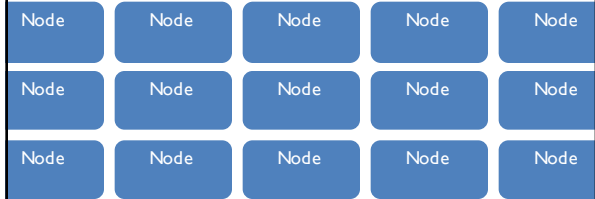
## Fundamental Issues

Applies to *any* multi-node system



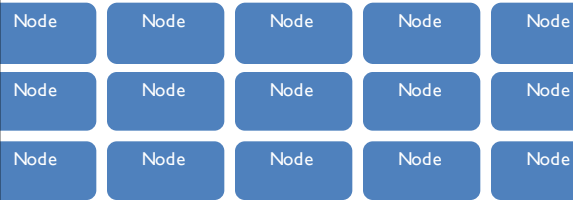
## Data Processing Issues

How to lay out, index, structure data to answer queries quickly, correctly



## Spark/Application Issues

How to use modern distributed computing system?



## Big Data Systems in the Wild

Spark  
Google Cloudflow  
Azure Cloud  
AWS/Redshift  
Tensorflow  
Cludera  
...

65

## Course Structure Overview

Three key modules and focus areas:

1. Data modeling and visualization (Wu)
  - \* Various data models and storage
  - \* Graph processing and big-data visualization
2. Storage at Scale (Ganbasu)
  - \* Challenges and core techniques for scalability and fault tolerance
  - \* Distributed transaction on sharded databases
  - \* Replication architectures and protocols
  - \* Design and implementation of Spanner, Google's geo-distributed, transactional store
3. Processing at Scale (Sahu)
  - \* Batch processing with Map Reduce and higher-level programming construct
  - \* Real-time responsive analytics with Spark and Spark Streams

Designing Machine Learning Systems with Big Data

66

## Course Administrative Details

Course materials

1. Primarily lecture notes. Additional reference readings will be provided as needed based on the lecture topics including research papers.

All course related submissions will be done using courseworks.

Important deadlines and communications will be done using Courseworks Announcement.

3-4 TAs will be available to assist in the course. We will announce their contact emails.

Good programming background in one of the languages Python/Java

[w4121.github.io](https://github.com/w4121)

67

## Grading and Project

### Grading

1. 60% Homework
2. 30% Tests/Quizzes
3. 10% Participation (ask/answer questions)

### Optional Project

1. 0-40% Extra credit (does not affect curve)

### A+s

1. Hand selected by instructors for exceptional work

68

## Logistics

Register with piazza

We will not answer direct emails

69

## 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  
Cheating = Failing grade

70

## Module I: Data Modeling Topics

Data models

Data cleaning

Data wrangling, Entity Resolution, Explanation

Large scale analytics

Visualizations and scaling them

71

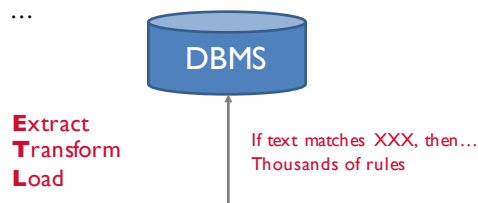
## How does data get into a DBMS?

Entity resolution

Data extraction

Missing data

...



72

## How does data get into a DBMS?

Entity resolution

Data extraction

Missing data

...

Extract  
Transform  
Load





The image is a screenshot of a web browser displaying a New York Times article. At the top, the New York Times logo is on the left, and a search icon is on the right. Below the logo, the word "Technology" is written. To the right of "Technology", the words "SUBSCRIBE" and "LOG IN" are visible. The article's title, "Medicaid's Data Gets an Internet-Era Makeover", is prominently displayed in a large, black, serif font. Below the title are three circular social media sharing icons: Twitter (a blue bird), Facebook (a blue 'f'), and a generic share icon (three dots). The byline "By STEVE LOHR" is followed by the date "JANUARY 9, 2017". The first paragraph of the article begins with "Jini Kim's relationship with [Medicaid](#) is business and personal." The second paragraph starts with "Her San Francisco start-up, [Nuna](#), while working with the federal government, has built a cloud-computing database of the nation's 74 million Medicaid patients and their treatment." The third paragraph begins with "[Medicaid](#), which provides health care to low-income people, is administered state by state. Extracting, cleaning and curating the information from so many [disparate](#) and dated computer systems was an extraordinary achievement, health and technology specialists say. This new collection of data could inform the coming [debate on Medicaid spending](#)." The fourth paragraph starts with "Andrew M. Slavitt, acting director of the Centers for Medicare and Medicaid Services, described the cloud database as 'near historic.'" Largely because Medicaid information resides in so many state-level computing silos, Mr. Slavitt explained, "we've never had a systemwide view across the program."

# How does data get into a DBMS?

The diagram illustrates the process of interactive data cleaning. A blue cylinder labeled "DBMS" is at the top. Below it, a black icon of a person represents the user. A curved arrow points from the user to the DBMS, and a straight arrow points back from the DBMS to the user, indicating a two-way interaction. The text "Interactive data cleaning" is written in red to the left of the arrows.

Interactive data cleaning

75

# How does data get into a DBMS?

Text → data records

Building sensor IDs – no consistency, arbitrary

BLDA | C600A\_ART

BLDC | C2\_\_\_\_TMR

**Automated Metadata Construction To Support Portable Building Applications**

Arka Bhattacharya, David Culler Electrical Engineering and Computer Sciences, UC Berkeley arka.culler@eecs.berkeley.edu	Dezhi Hong, Kamin Whitehouse University of Virginia dh5gm.whitehouse@virginia.edu
Jorge Ortiz IBM Research jortiz@us.ibm.com	Eugene Wu Computer Science, Columbia University ewu@cs.columbia.edu

**ABSTRACT**

Active sensor networks are an active topic for emerging systems, it has long been a core challenge in the sensor network. Often, a critical step in the

76

[illegible]

# Large scale analytics

Data volumes too large to even scan once

## How to deal with this?

- Spend more time
- Concurrency
- Reduce data size
- Read less data
- Do less work
- Waste less time doing work

78

# Large scale analytics

- Columnar databases
- In-memory databases
- Intermediate results
- Graph “databases”
- Sketching and sampling

## Visualization

How to think about and approach visualization

Modern visualization tools

How to scale visualizations

80

## Module 2: Storage at Scale

- Two key reasons for distributed systems:
  - Scaling: system capacity grows proportionally with # of machines.
  - Fault tolerance: being able to continue operation despite failures, which can happen constantly in a large system.
- But achieving scale and fault tolerance (at scale) is hard.
  - Consistency, coherence, semantics are one challenge.
  - Fault tolerance requires coordination, which limits scalability.
- The second model will teach key techniques and protocols for scaling and fault tolerance, with a particular focus on one system: Google's Spanner storage system.

81

## Module 3: Processing at Scale

Computation on huge amount of data is not a luxury – it is a necessity!

Imagine Facebook logs for logins. FB wants to compute how many people are logging in from which continents for each hour.

How to compute?

82

## What's the big deal?

- How big is the data?
- Huge data – the data file does not fit into single server's disks...how do you compute if data does not even fit into server's storage?
- Data is on multiple servers – on a cluster of servers. So how do you compute and where do you compute what?
- How do we compute the final results?
- Who takes care of some machine or computing failure?
- How do you automate such computations spread across machines?

83

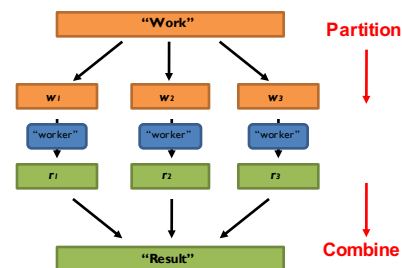
## Computing Models for Big Data

We will learn two computing paradigm for big data on a cluster of machines

- Batch processing with Map Reduce
  - Idea is to divide and conquer the task – compute partial results on smaller chunks of data and then merge the partial results to compute final result
  - Move computing task to where data is
- Real-time processing with Spark
  - Map Reduce is great but too slow due to lot of disk based operations
  - Spark computes with in-memory data

84

## Divide and Conquer



85

## So what is Hadoop/Map Reduce

Hadoop/Map Reduce is a computing system on a cluster of machines that provide at the minimum the following

- Storage across a cluster of machines (HDFS)
- A computation model to divide-conquer a task (map-reduce)
- A runtime to enable map-reduce style of computation

86

## Why MapReduce not efficient for iterative computations?

- MapReduce is an excellent computing model that scales for log processing type of computations described earlier.
- What about iterative models that use the same data again and again?
  - Every operation is to read and write to disk. So every iteration requires reading and writing to disk. Too many disk based operations for iterative computing.
  - Many machine learning based computations are iterative in nature.
- So what is the solution? Can the data be somehow kept in memory until all the operations on it completes...
- **Spark Model:** Resilient Distributed Datasets (RDD)
  - Recent computing model that is 100x faster and more suitable for iterative and real-time analytics
  - We will learn how to write real-time analytics using Spark and Spark Streams.

87