# CC5212-1
## Procesamiento Masivo de Datos
## Otoño 2016

## Lecture 1: Introduction

Aidan Hogan

aidhog@gmail.com

# THE VALUE OF DATA

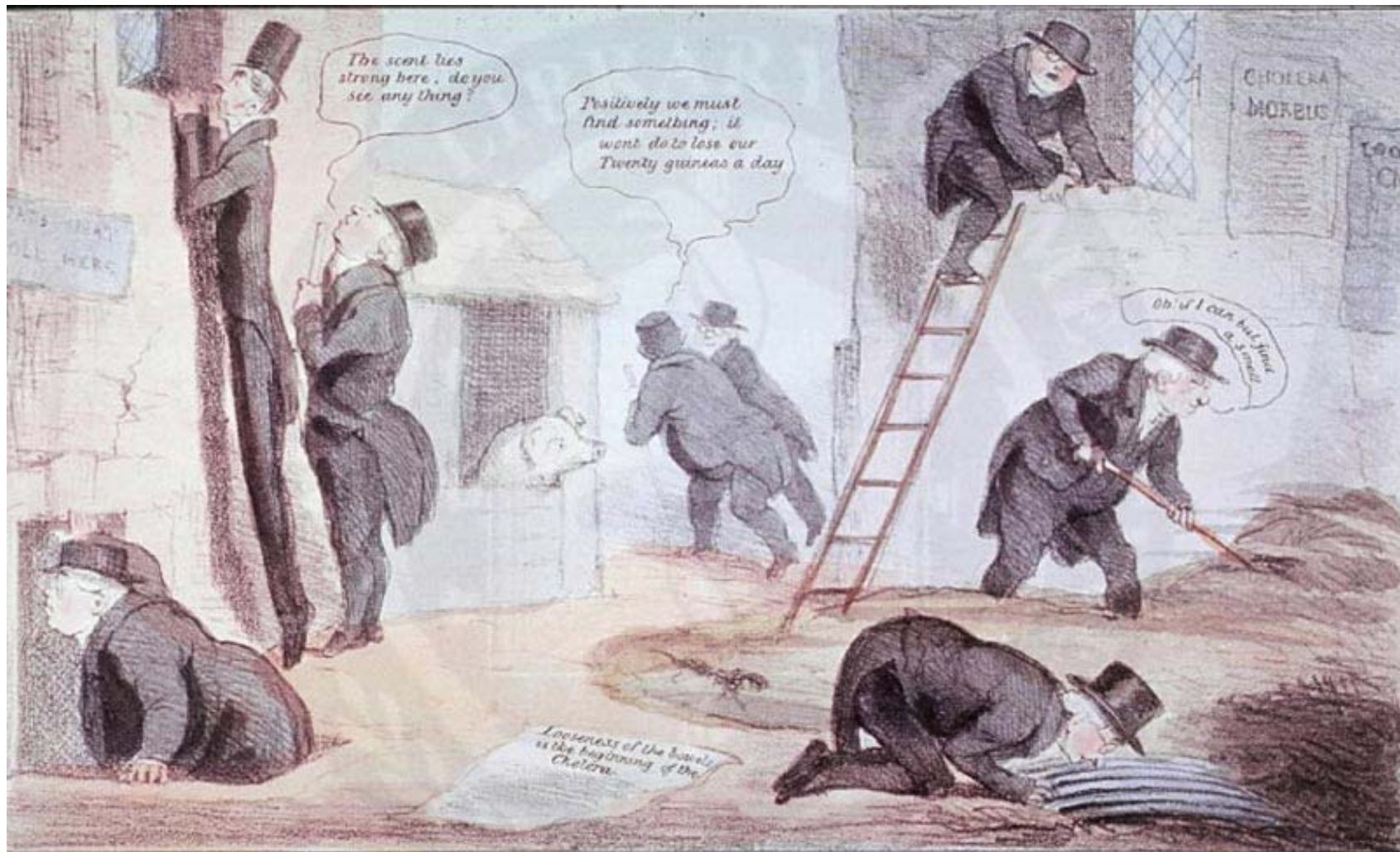# Soho, London, 1854





BLUE STAGE OF THE SPASMODIC CHOLERA.
Sketch of a Girl who died of Cholera, in Sunderland, November, 1831.

# The mystery of cholera

# The hunt for the invisible cholera

# Cholera: Galen's miasma theory

# John Snow: 1813–1858

# The Survey of Soho
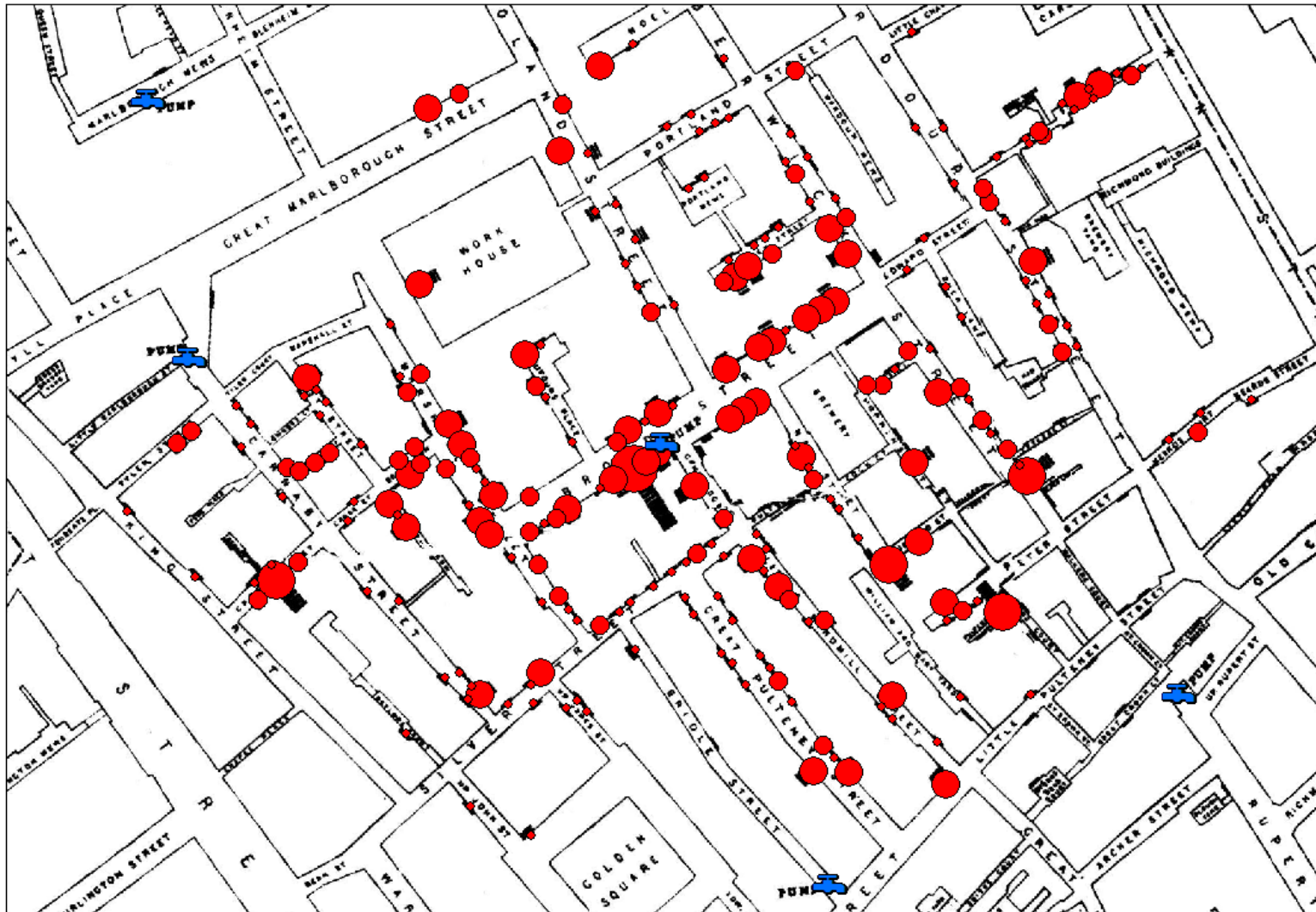
# Data collection

TABLE VI.

*The Mortality from Cholera in 1854, in Thirty-one Sub-Districts, as compared with Calculations founded on the Results shewn in Table v.*

| Registration Districts. | Registration Sub-Districts. | Population in 1851. | Estimated population supplied with water as under. | | | Deaths from cholera in 1854. | | Calculated mortality in the population, supplied with water as under. | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Southwark and Vauxhall Co. | Lambeth Co. | Both Companies together. | Total deaths. | Deaths per 10,000 living. | Southwark and Vauxhall Co. at 160 per 10,000. | Lambeth Co. at 27 per 10,000. | The two Companies. | Calculated deaths per 10,000 supplied by the two Companies. |
| St. Saviour, Southw. - | 1. Christchurch - - - - | 16,022 | 2,015 | 13,234 | 16,149 | 113 | 71 | 46 | 36 | 82 | 57 |
| | 2. St. Saviour - - - - | 19,709 | 16,337 | 898 | 17,235 | 378 | 192 | 261 | 2 | 263 | 153 |
| St. Olave - - - - | 1. St. Olave - - - - - | 8,015 | 8,745 | 0 | 8,745 | 161 | 201 | 140 | 0 | 140 | 160 |
| | 2. St. John, Horselydown | 11,360 | 9,360 | 0 | 9,360 | 152 | 134 | 150 | 0 | 150 | 160 |
| Bermondsey - - - | 1. St. James - - - - - | 18,899 | 23,173 | 693 | 23,866 | 362 | 192 | 370 | 2 | 372 | 156 |
| | 2. St. Mary Magdalen - | 13,934 | 17,258 | 0 | 17,258 | 247 | 177 | 276 | 0 | 276 | 160 |
| | 3. Leather Market - - | 15,295 | 14,003 | 1,092 | 15,095 | 237 | 155 | 224 | 3 | 227 | 150 |
| St. George, Southw. - | 1. Kent Road - - - - | 18,126 | 12,630 | 3,997 | 16,627 | 177 | 98 | 202 | 11 | 213 | 134 |
| | 2. Borough Road - - - | 15,862 | 8,937 | 6,672 | 15,609 | 271 | 171 | 143 | 18 | 161 | 104 |
| | 3. London Road - - - | 17,836 | 2,872 | 11,497 | 14,369 | 95 | 53 | 46 | 31 | 79 | 55 |
| Newington - - - - | 1. Trinity - - - - - | 20,922 | 10,132 | 8,370 | 18,502 | 211 | 101 | 162 | 22 | 184 | 99 |
| | 2. St. Peter, Walworth - | 29,861 | 14,274 | 10,724 | 24,998 | 391 | 131 | 228 | 29 | 257 | 103 |
| | 3. St. Mary - - - - - | 14,033 | 2,983 | 5,484 | 8,467 | 92 | 66 | 48 | 15 | 63 | 74 |

# What the data showed …



Dessin satirique (1866)

# What the data showed …

# 616 deaths, 8 days later ...



The Red Granite kerbstone marks the site of the historic **BROAD STREET PUM** associated with Dr. John Snow's discovery in 1854 that Cholera is conveyed by water

# Cholera notice ca. 1866



**CHOLERA**
**AND WATER.**

**BOARD OF WORKS**

FOR THE LIMEHOUSE DISTRICT,
Comprising Limehouse, Ratcliff, Shadwell,
and Wapping.

The INHABITANTS of the District within
which CHOLERA IS PREVAILING, are
earnestly advised

**NOT TO DRINK ANY WATER**
**WHICH HAS NOT**
**PREVIOUSLY BEEN BOILED.**

Fresh Water ought to be Boiled every
Morning for the day's use, and what
remains of it ought to be thrown away
at night. The Water ought not to stand
where any kind of dirt can get into it,
and great care ought to be given to see
that Water Butts and Cisterns are free
from dirt.

BY ORDER,

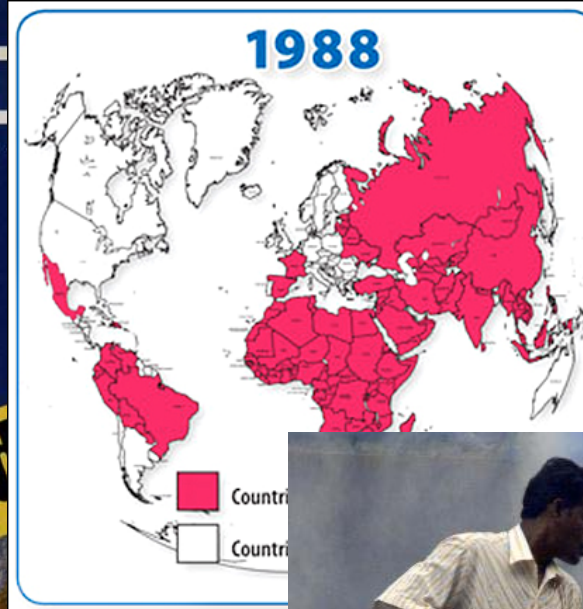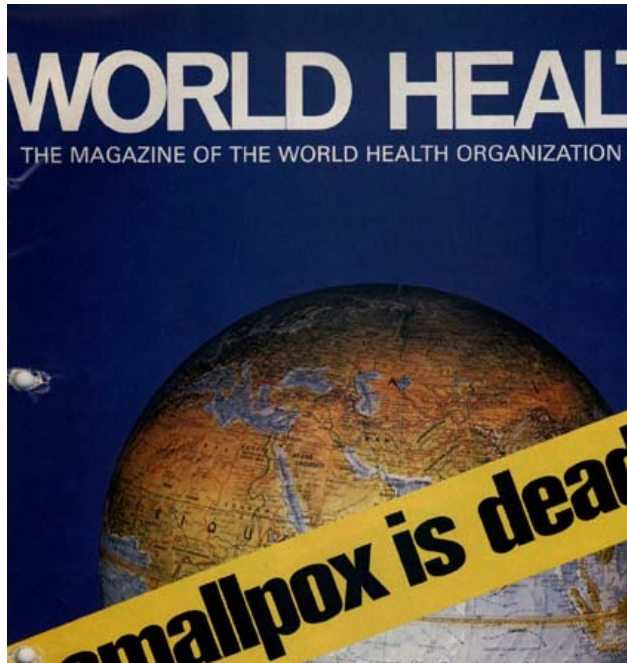**THOS. W. RATCLIFF,**
CLERK OF THE BOARD.



John Snow

# Thirty years before discovery of *V. cholerae*

# John Snow: Father of Epidemiology

# Epidemiology's Success Stories

# Value of data: Not just epidemiology

# (Paper) notebooks no longer good enough

# THE GROWTH OF DATA

# "Big Data"



**Wikipedia**
≈ 5.9 TB of data
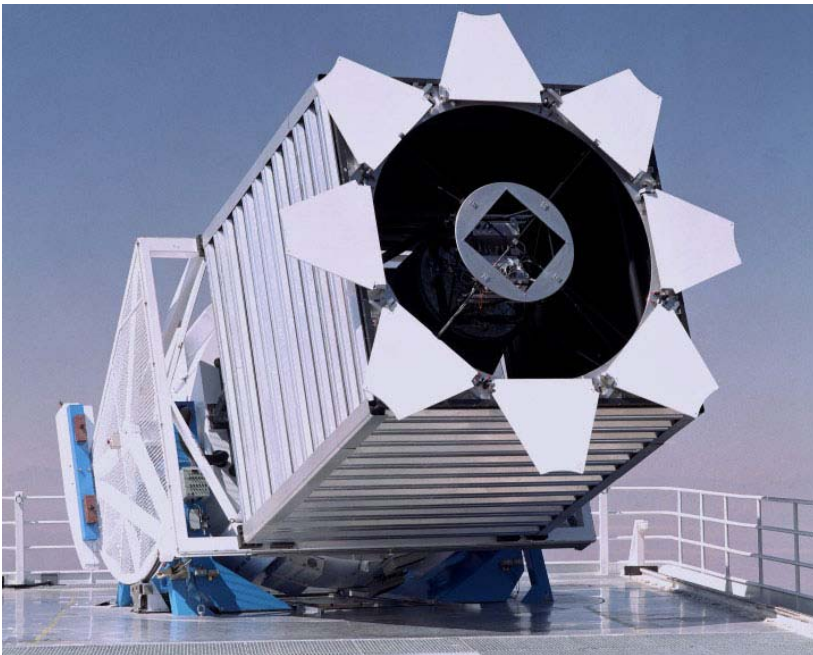    (*Jan. 2010 Dump*)

1 Wiki = 1 Wikipedia

# "Big Data"



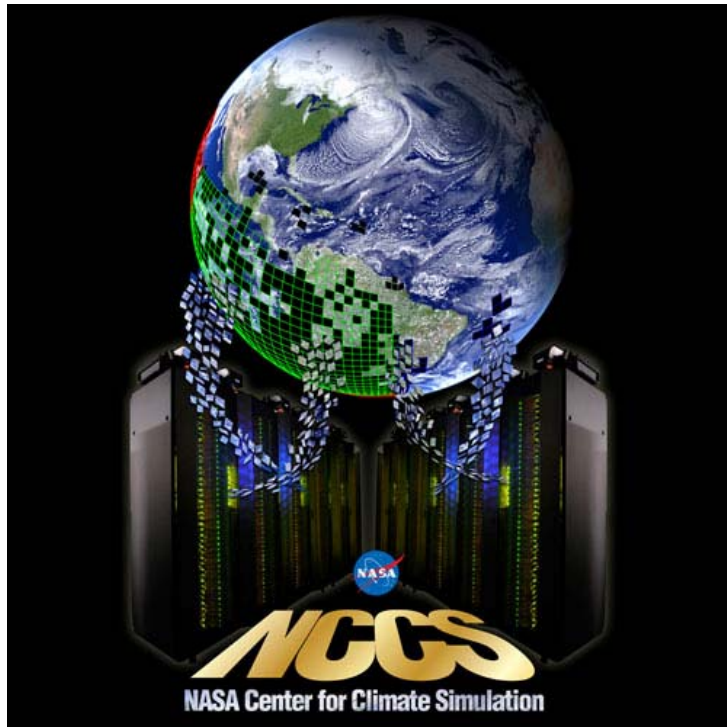**US Library of Congress**
≈ 235 TB archived
≈ 40 Wiki

# "Big Data"



**Sloan Digital Sky Survey**
≈ 200 GB/day
≈ 73 TB/year
≈ 12 Wiki/year

# "Big Data"



**NASA Center for Climate Simulation**
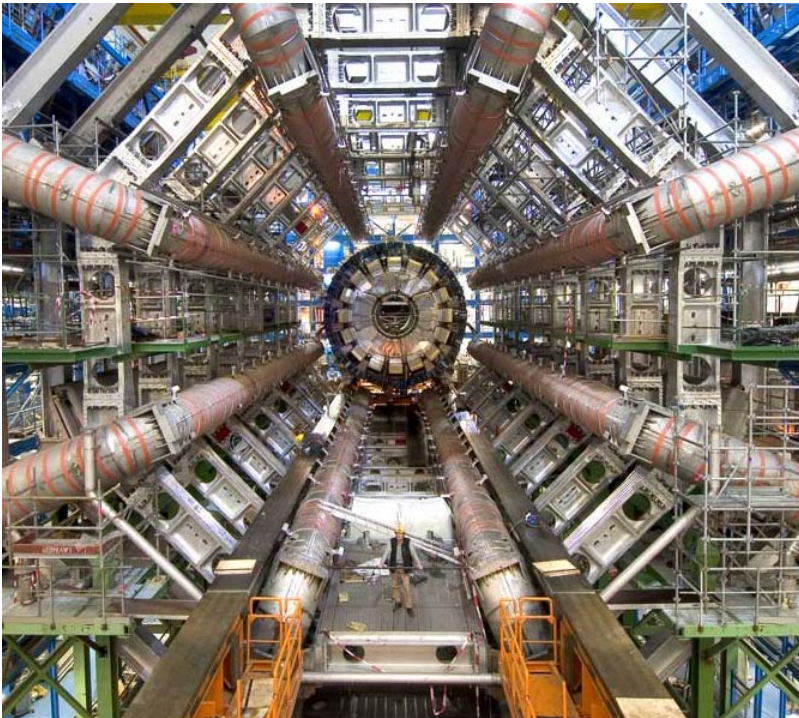≈ 32 PB archived
≈ 5,614 Wiki

# "Big Data"



**Facebook**
≈ 100 TB/day added
≈ 17 Wiki/day
≈ 6,186 Wiki/year
   (*as of Mar. 2010*)

# "Big Data"



**Large Hadron Collider**
≈ 15 PB/year
≈ 2,542 Wikipedias/year

# "Big Data"



**Google**
≈ 20 PB/day <u>processed</u>
≈ 3,389 Wiki/day
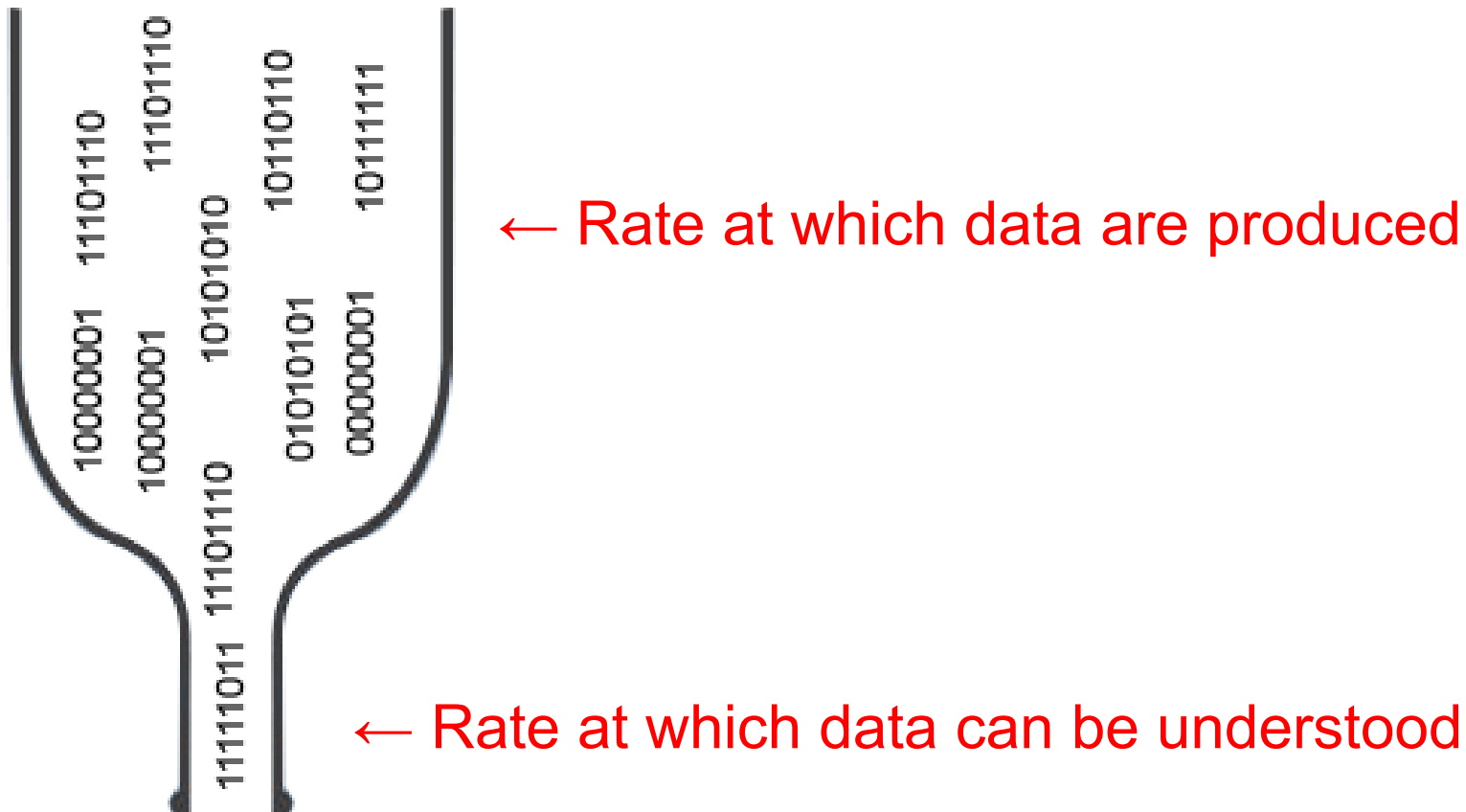≈ 7,300,000 Wiki/year
      *(Jan. 2010)*

# "Big Data"



**Internet (2016)**
≈ 1.3 ZB/year
≈ 220,338,983 Wiki/year
*(2016 IP traffic; Cisco est.)*

# Data: A Modern-day Bottleneck?

← Rate at which data are produced
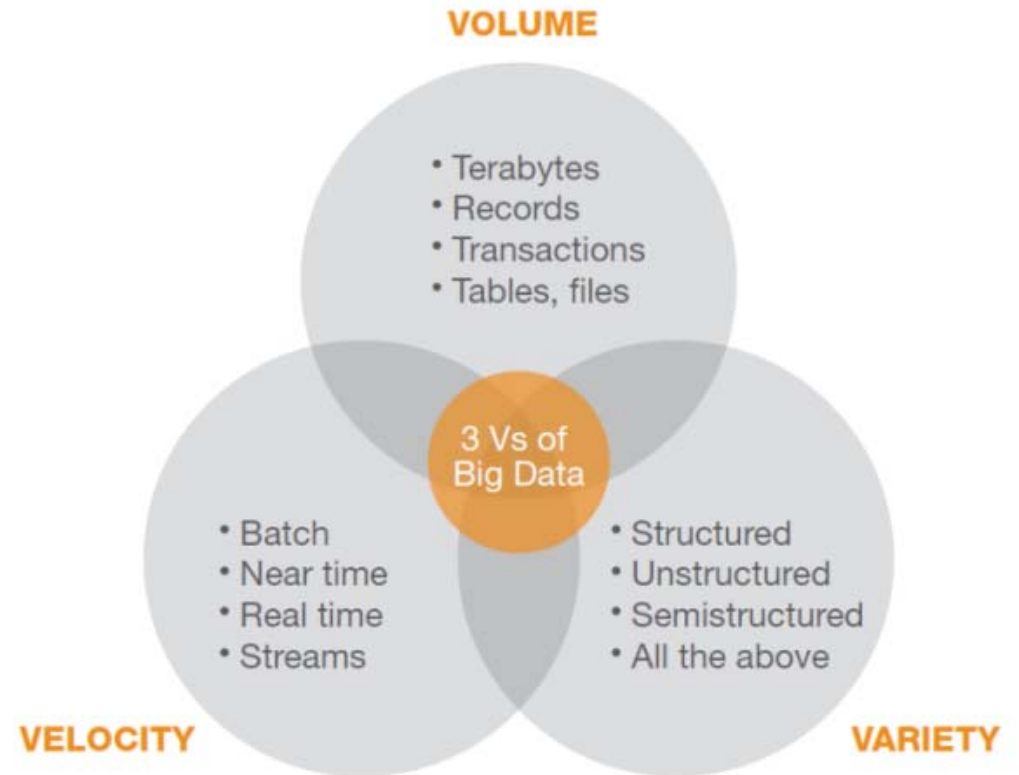
← Rate at which data can be understood

# BIG DATA

# "Big Data"

- A buzz-word: no *precise* definition?
- Data that are too big to process by "conventional means"
- A call for Computer Scientists to produce new techniques to crunch even more data

- Storage, processing, querying, analytics, data mining, applications, visualisations …

# How many V's in "Big Data"?



**VOLUME**
- Terabytes
- Records
- Transactions
- Tables, files

3 Vs of Big Data

**VELOCITY**
- Batch
- Near time
- Real time
- Streams

**VARIETY**
- Structured
- Unstructured
- Semistructured
- All the above

- Three 'V's:
  - Volume (large amounts of data)
  - Velocity (rapidly changing data)
  - Variety (different data sources and formats)
- Maybe more (Value, Veracity)

# "BIG DATA" IN ACTION …

# Social Media

# What's happening here? (Trendsmap)

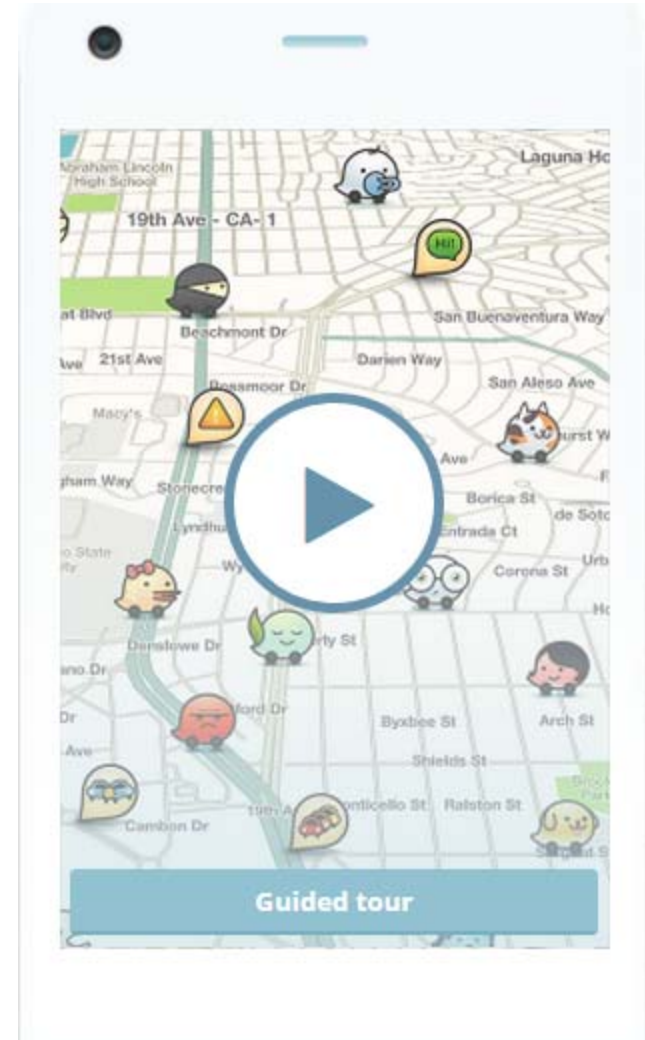*"What are the hot topics of discussion in an area"*

- Analyse tags of geographical tweets

# What's the fastest route? (Waze)
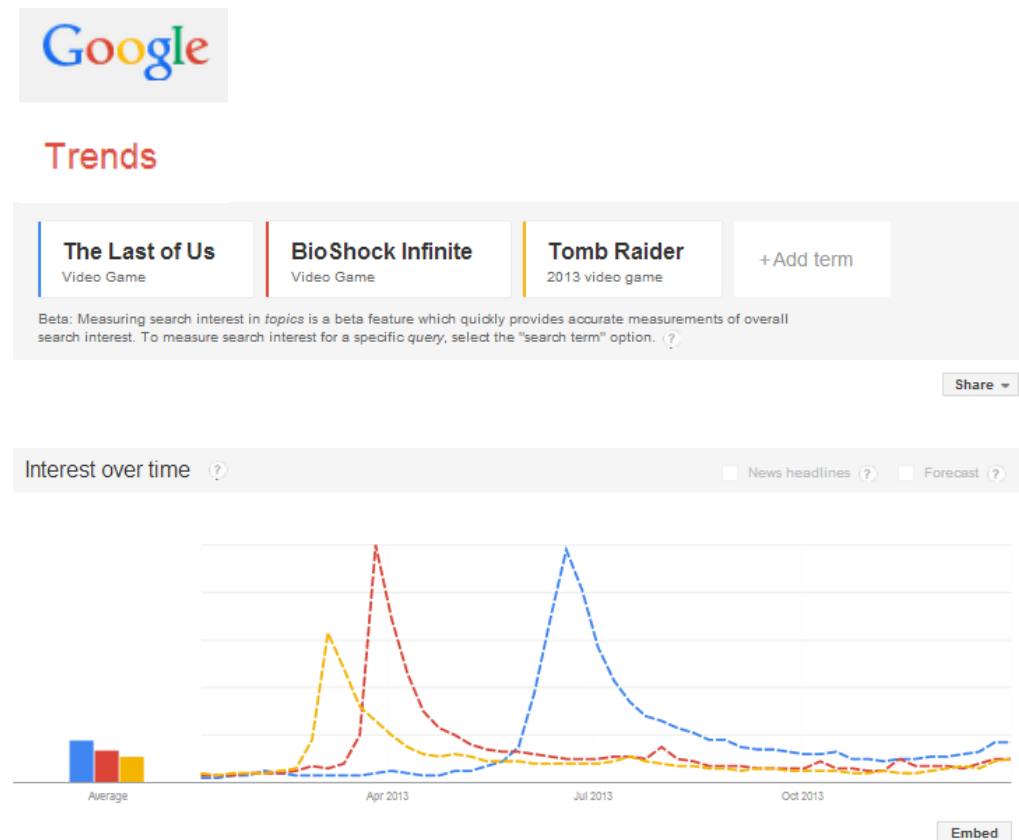
*"What's the fastest route to get home right now?"*

- Processes real journeys to build background knowledge

- "Participatory Sensing"

# Christmas Predictions for Stores

*"What will be the hot items to stock up on this Christmas? We don't want to sell out!"*

- Analyse product hype on Twitter, Search Engines and Social Networks
- Analyse transaction histories

# Get Elected President (Narwhal)

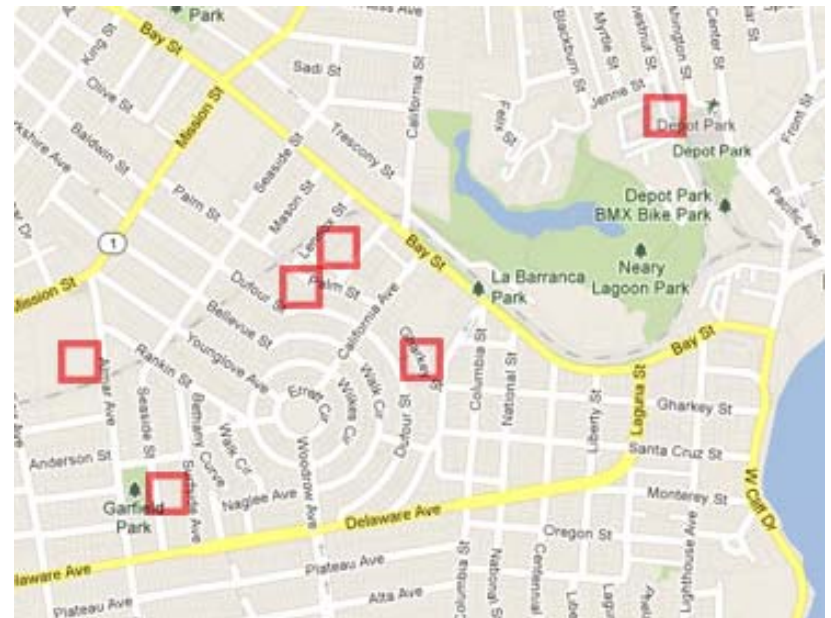*"Who are the undecided voters and how can I convince them to vote for me?"*

- User profiles built and integrated from online sources
- Targeted emails sent to voters based on profile

# Predicting Pre-crime (PredPol)

*"What areas of the city are most need of police patrol at 13:55 on Mondays?"*

- PredPol system used by Santa Cruz (US) police to target patrols
- Predictions based on analysis of 8 years of historical crime data
- Minority Report!

# IBM Watson: Jeopardy Winner

*"William Wilkinson's "An Account of the Principalities of Wallachia and Moldavia" inspired this author's most famous novel."*

- Indexed 200 million pages of structured and unstructured content
- An ensemble of 100 techniques simulating AI-like behaviour



Check it out on YouTube!

# "BIG DATA" NEEDS
# "MASSIVE DATA PROCESSING" …

# Every Application is Different …

- **Data** can be
  - Structured data (JSON, XML, CSV, Relational Databases, HTML form data)
  - Unstructured data (text document, comments, tweets)
  - And everything in-between!
  - **Often a mix!**

# Every Application is Different …

- **Processing** can involve:
  - Natural Language Processing (sentiment analysis, topic extraction, entity recognition, etc.)
  - Machine Learning and Statistics (pattern recognition, classification, event detection, regression analysis, etc.)
  - Even inference! (Datalog, constraint checking, etc.)
  - And everything in-between!
  - **Often a mix!**

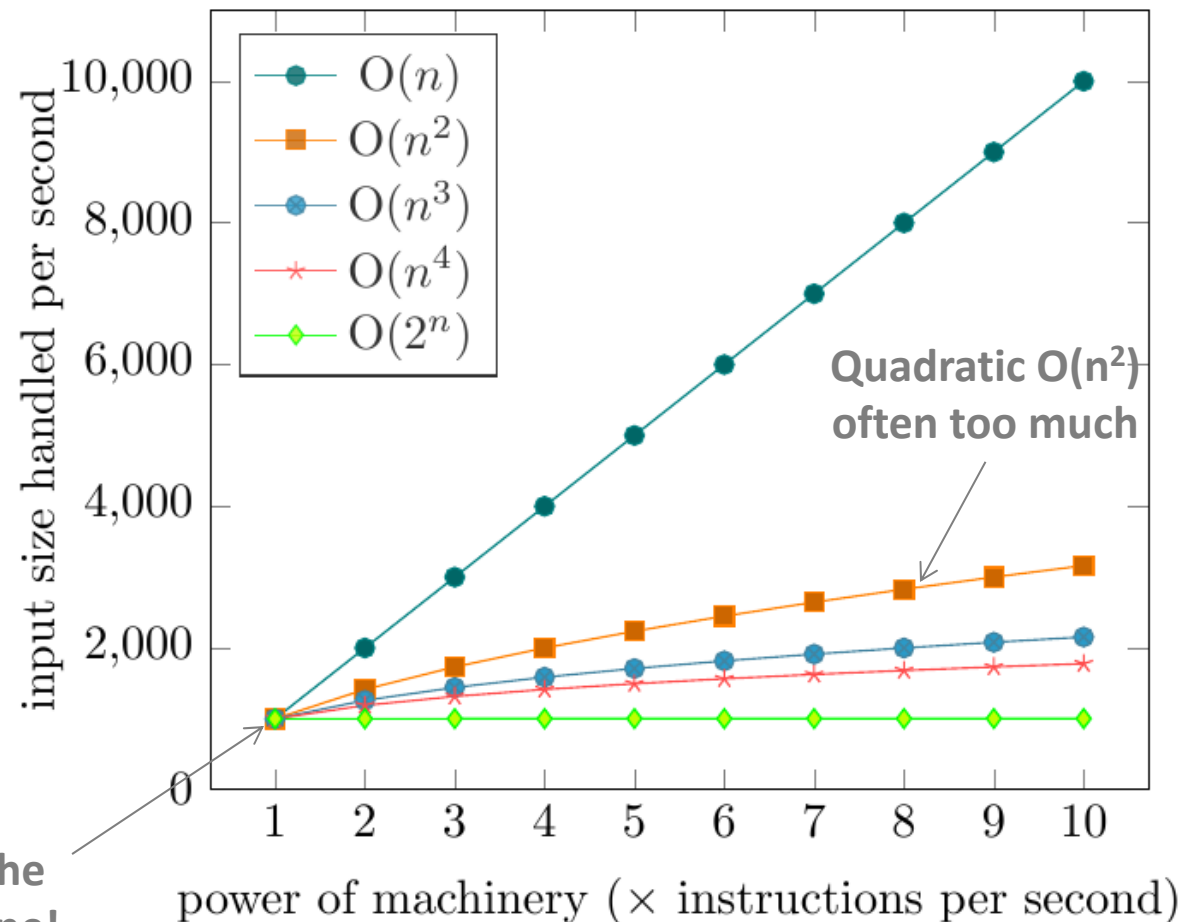# Scale is a Common Factor ...

- Cannot run expensive algorithms

*I have an algorithm.*

*I have a machine that can process 1,000 input items in an hour.*

*If I buy a machine that is <u>n</u> times as powerful, how many input items can I process in an hour?*

*Depends on algorithm complexity of course!*

Note: Not the same machine!



Legend: $O(n)$, $O(n^2)$, $O(n^3)$, $O(n^4)$, $O(2^n)$

Quadratic $O(n^2)$ often too much

y-axis: input size handled per second — 0, 2,000, 4,000, 6,000, 8,000, 10,000

x-axis: power of machinery ($\times$ instructions per second) — 1 2 3 4 5 6 7 8 9 10

# Scale is a Common Factor ...

- One machine that's *n* times as powerful?

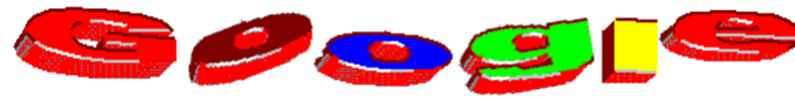**VS.**

- *n* machines that are equally as powerful?

# Scale is a Common Factor …

- Data-intensive (our focus!)
  - Inexpensive algorithms / Large inputs
  - e.g., Google, Facebook, Twitter

- Compute-intensive (not our focus!)
  - More expensive algorithms / Smaller inputs
  - e.g., climate simulations, chess games, combinatorials

- No black and white!

# "MASSIVE DATA PROCESSING" NEEDS "DISTRIBUTED COMPUTING" …

# Distributed Computing

- ## Need more than one machine!

- Google ca. 1998:

# Distributed Computing

- ## Need more than one machine!

- Google ca. 2014:

# Data Transport Costs

- Need to divide tasks over many machines
  - Machines need to communicate
    - … but not too much!
  - Data transport costs (*simplified*):

| Main Memory | Solid-state Disk | Hard-disk | Network |
|:---:|:---:|:---:|:---:|

**Need to minimise network costs!**

# Data Placement

- Need to think carefully about where to put what data!

*I have four machines to run my website. I have 10 million users.*

*Each user has personal profile data, photos, friends and games.*

*How should I split the data up over the machines?*

*Depends on application of course!*
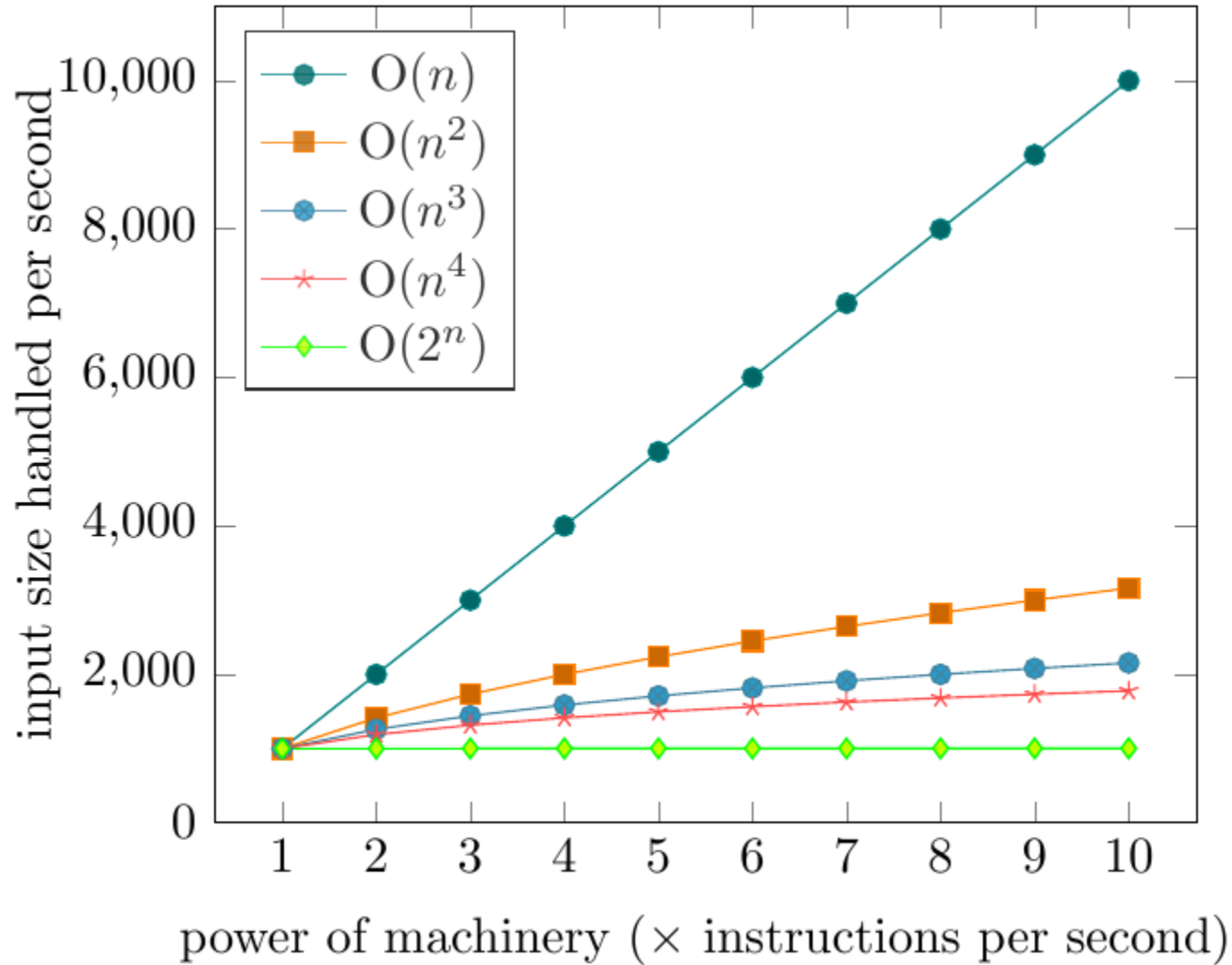
*(But good design principles apply universally!)*

# Network/Node Failures

- Need to think about failures!



**Lot of machines: likely one will break!**

# Network/Node Failures

- Need to think (<span style="color:red">even more!</span>) carefully about where to put what data!

*I have four machines to run my website. I have 10 million users.*

*Each user has a personal profile, photos, friends and apps.*

*How should I split the data up over the machines?*

<span style="color:red">*Depends on application of course!*</span>

*(But good design principles apply universally!)*

# Human Distributed Computation



Similar Principles!

# "DISTRIBUTED COMPUTING" LIMITS & CHALLENGES ...

# Distribution Not Always Applicable!

# Distributed Development Difficult

- Distributed systems can be complex
- Tasks take a long time!
  - Bugs may not become apparent for hours
  - Lots of data = lots of counter-examples
  - **Need to balance load!**
- Multiple machines to take care of
  - Data in different locations
  - Logs and messages in different places
  - **Need to handle failures!**

# Frameworks/Abstractions can Help

- ## For Distrib. Processing



- ## For Distrib. Storage

# HOW DOES TWITTER WORK?

Based on 2013 slides by Twitter lead architect: Raffi Krikorian

"Twitter Timelines at Scale"

# Big Data at Twitter

- 150 million active worldwide users
- 400 million tweets per day
  - mean: 4,600 tweets per second
  - max: 143,199 tweets per second

- 300 thousand queries/sec for user timelines

- 6 thousand queries/sec for custom search

*What should be the priority for optimisation?*

# Supporting timelines:write

- 300 thousand queries per second

# High-fanout



@ladygaga ✔
31 million followers

@katyperry ✔
28 million followers

@justinbieber ✔
28 million followers

@barackobama ✔
23 million followers

# Supporting timelines: read

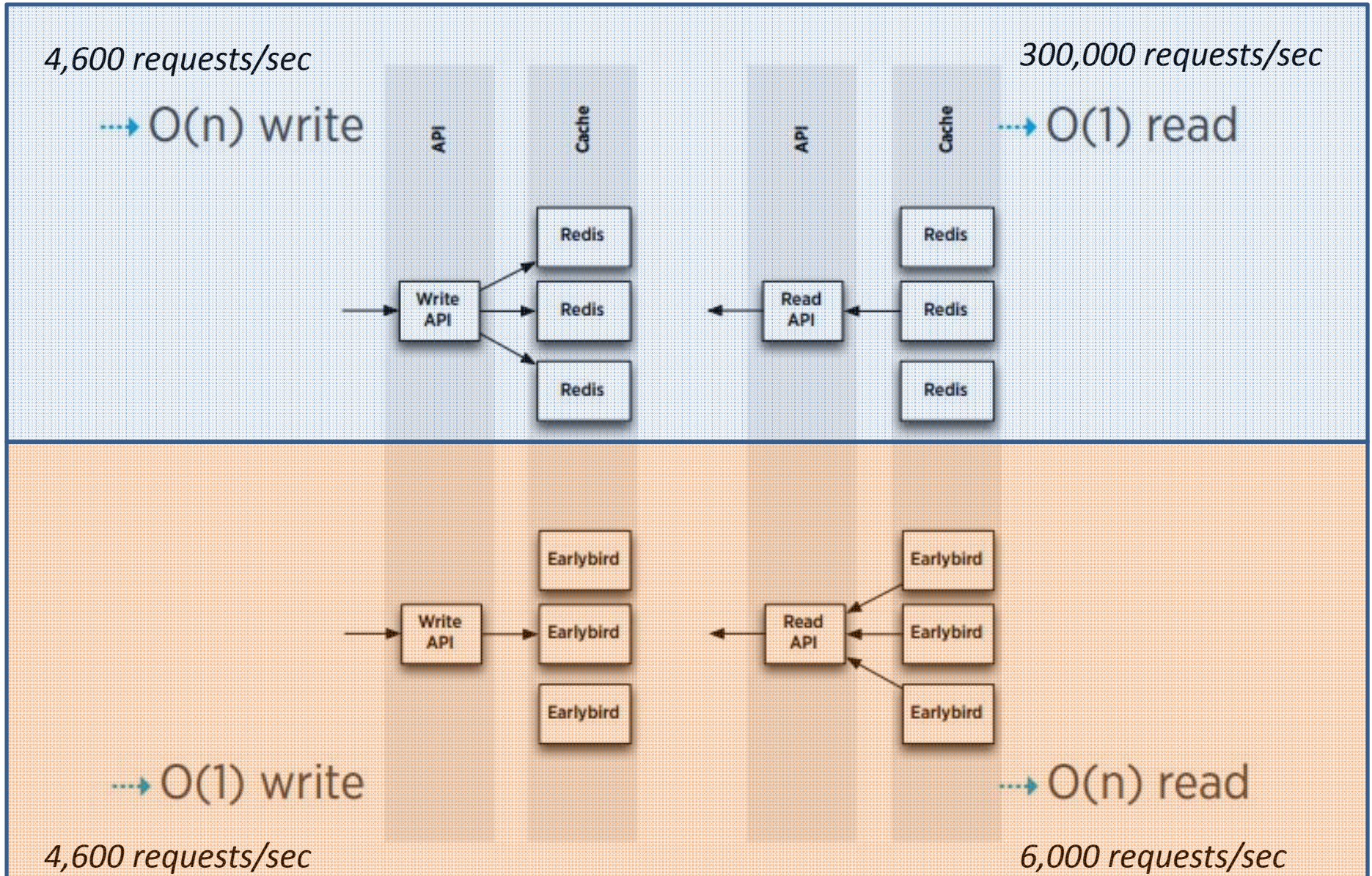- 300 thousand queries per second



**1ms @p50**
**4ms @p99**

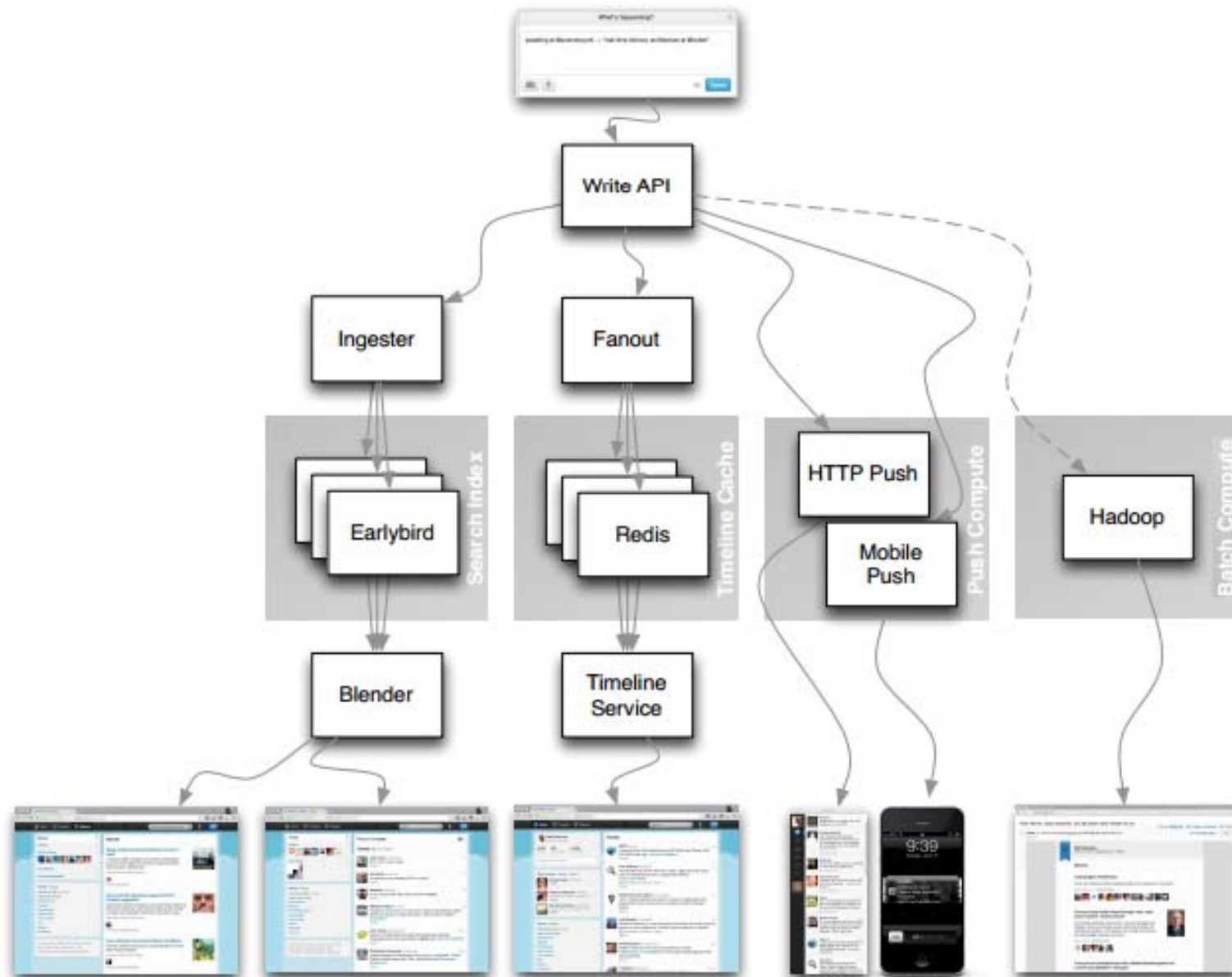# Supporting text search

- Information retrieval
  - Earlybird: Lucene clone
  - Write once
  - Query many

# Timeline vs. Search

# Twitter: Full Architecture

# Big Data at Twitter

- 150 million active worldwide users
- 400 million tweets per day
  - 4,600 tweets per second
  - max: 143,199 tweets per second

- 300 thousand queries/sec for user timelines

- 6 thousand queries/sec for custom search

# "Procesamiento Masivo de Datos"
## About the Course …

# What the Course Is/Is Not

- Data-intensive not Compute-intensive

- Distributed tasks not networking

- Commodity hardware not big supercomputers

- General methods not specific algorithms

- Practical methods with a little theory

# What the Course *Is*!

- Principles of Distributed Computing [2 weeks]
- Distributed Processing Frameworks [3 weeks]
- Information Retrieval [3 weeks]
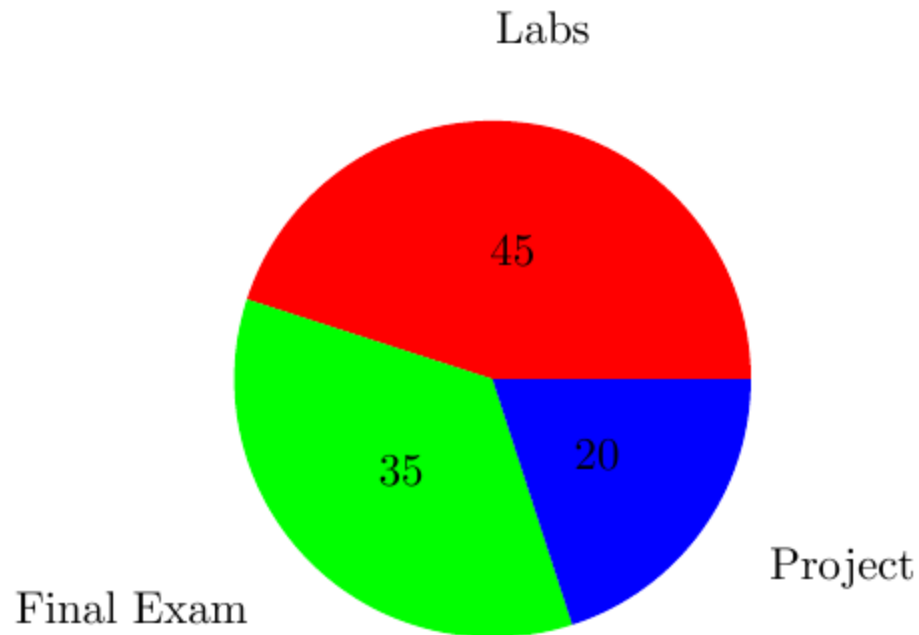- Principles of Distributed Databases [3 weeks]
- Projects [2 weeks]

# Course Structure

- ~1.5 hours of lectures per week [Monday]
- 1.5 hours of labs per week [Wednesday]
    - To be turned in by Friday evening
    - Mostly Java
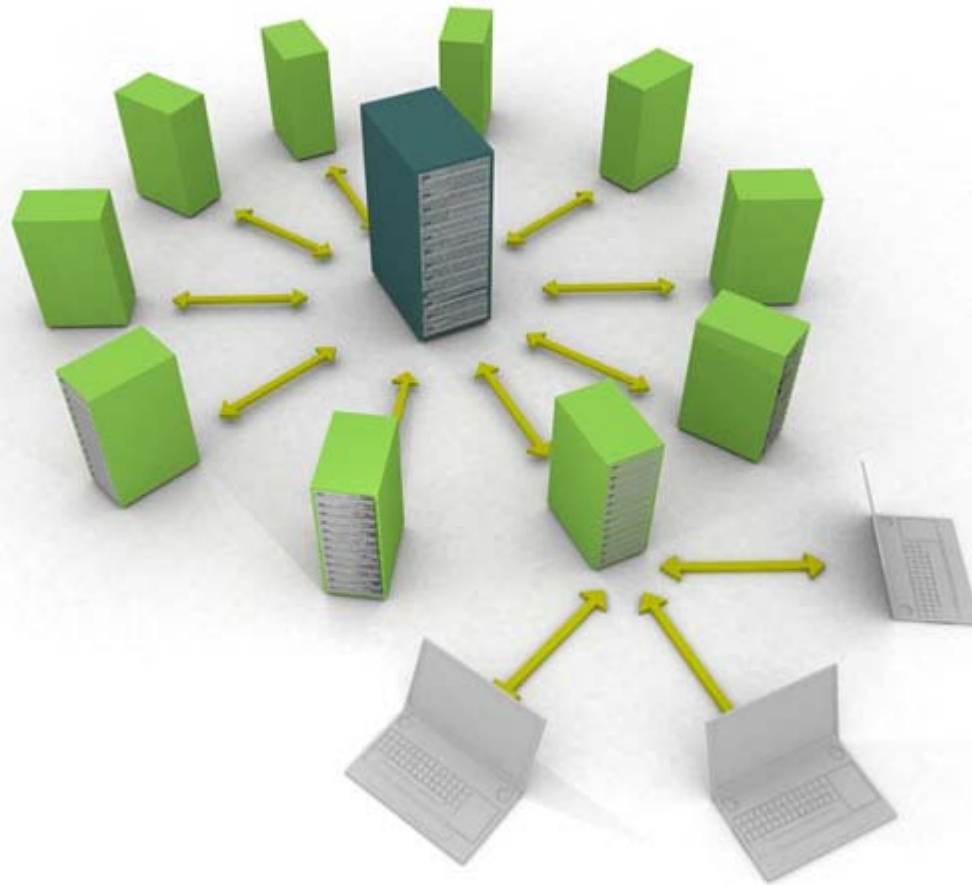    - In Lab on 3$^{rd}$ floor, edificio norte

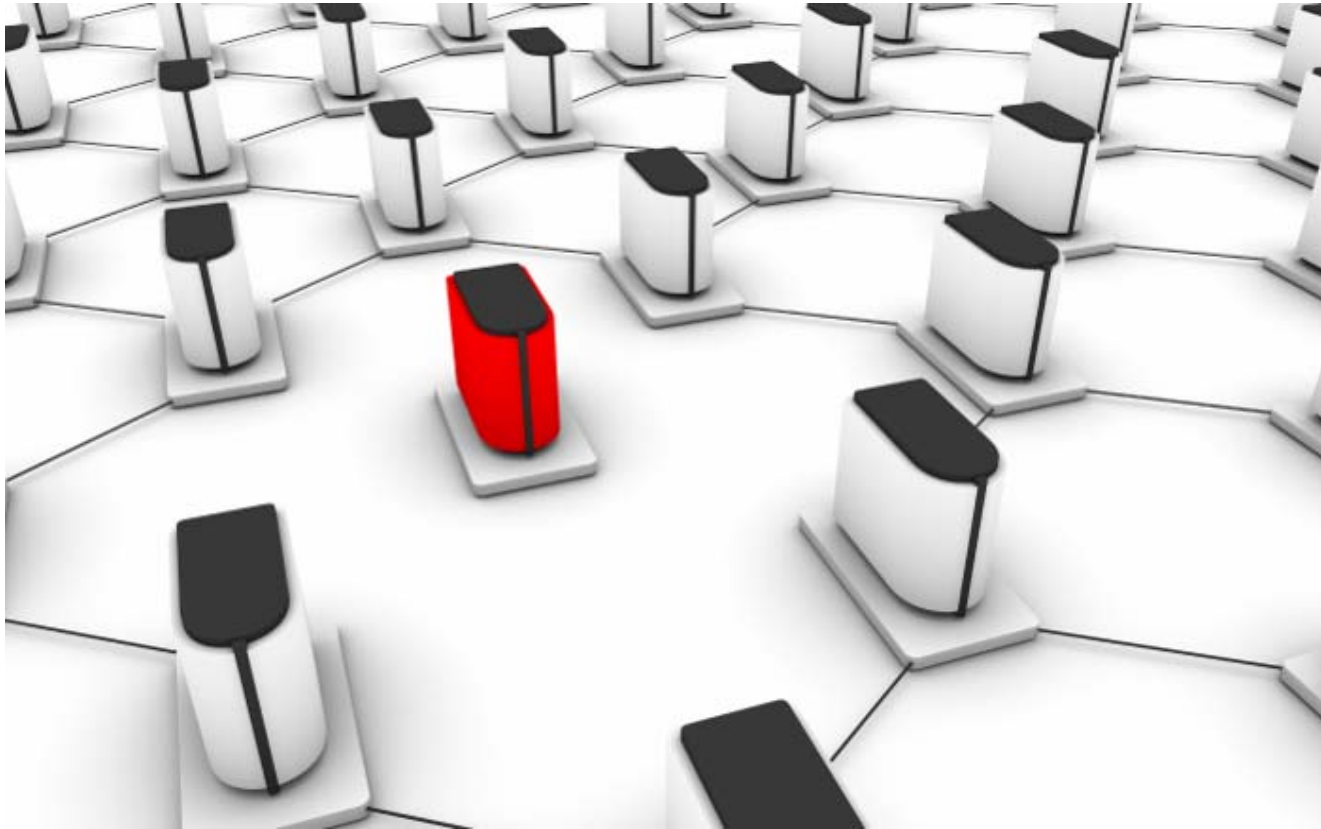http://aidanhogan.com/teaching/cc5212-1-2016/

# Course Marking

- 45% for Weekly Labs (~3% a lab!)
- 20% for Small Class Project
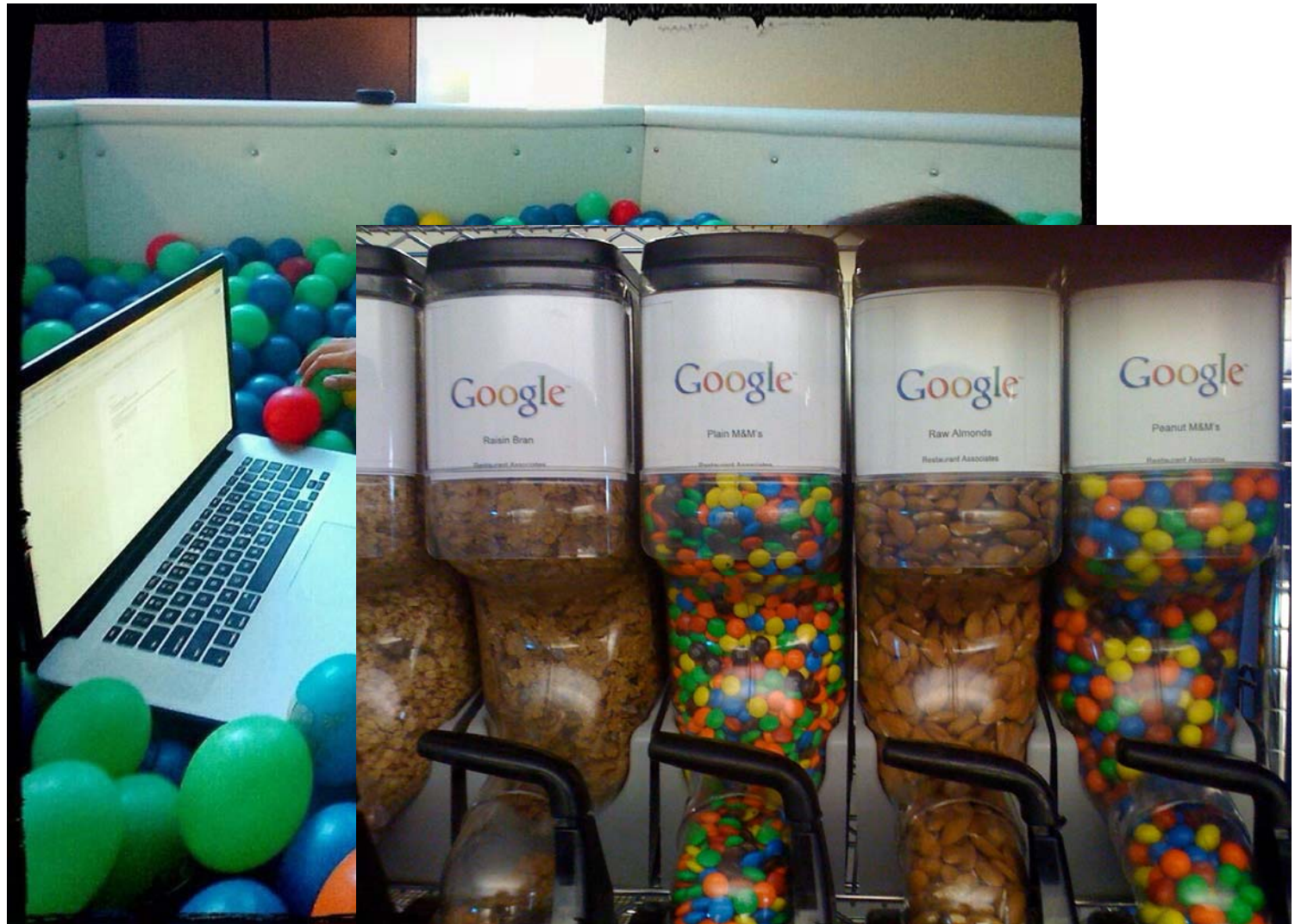- 35% for Final Exam [more challenging]

# Outcomes!

# Outcomes!

# Outcomes!

# Outcomes!

# Outcomes!

# Outcomes!

# Questions?