

Data-Intensive Distributed Computing

CS 431/631 451/651 (Fall 2019)

Part 1: MapReduce Algorithm Design (1/4)

Ali Abedi

These slides are available at <https://www.student.cs.uwaterloo.ca/~cs451/>



This work is licensed under a Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States
See <http://creativecommons.org/licenses/by-nc-sa/3.0/us/> for details

Agenda for Today

Who am I?

What is big data?

Why big data?

What is this course about?

Administrivia

Who am I?

PhD from Waterloo (2017)
Systems and Networking Research Group



Big Data



1950s

**\$3398
10MB**

THE HARD DISK YOU'VE BEEN WAITING FOR

XP introduces a complete micro-size disk system with more...

- **MORE STORAGE**
- **MORE SPEED**
- **MORE VALUE**
- **MORE SUPPORT**

DO users... The XCOMP subsystem is now available with 10 megabytes of storage, 5 megabytes available at \$2,999.00. Compare the price and features of any other 5 1/4-inch or even 8-inch system, and you'll agree that XCOMP's value is unbeatable.

OUTPERFORMS OTHER HARD DISKS

Floppy disk and larger, more expensive hard disks are no match for this powerful little system. More data is available on every seek: 64K on 10MB and 32K on 5MB. Faster seek time too — an average of 70MS, II provides solid performance anywhere with only 20 watts of power. Data is protected in the sealed enclosure, and the landing zone for heads provides another margin of safety. The optional power board plugs directly into the S100 bus and provides power

MORE SOFTWARE

Included with the system is software for testing, formatting, I/O drivers for CP/M[®], plus an automatic CP/M driver attach program. Support software and drivers for MP/M[®] and Oasis[®] are also available. The sophisticated formatting program assigns alternate sectors for any weak sectors detected during formatting, assuring the lowest possible error rate — at least ten times better than floppies.

WARRANTY

The system has a full one-year warranty on parts and workmanship.

ALSO AVAILABLE FROM XCOMP

- General Purpose controllers (8 bit interface), with easy interface to microprocessor-based systems.
- GP controller adapter that plugs directly into most Z80 computers.
- ST/R GP controller for the 5MB and 10MB drive above, with ST506 type interface.
- SM/R GP controller for SA1000 interface.
- ST/S, SQ/S, and SM/S, same as above, for the S100 bus.

Quantity discounts available. Distributor, Dealer, and OEM inquiries invited.

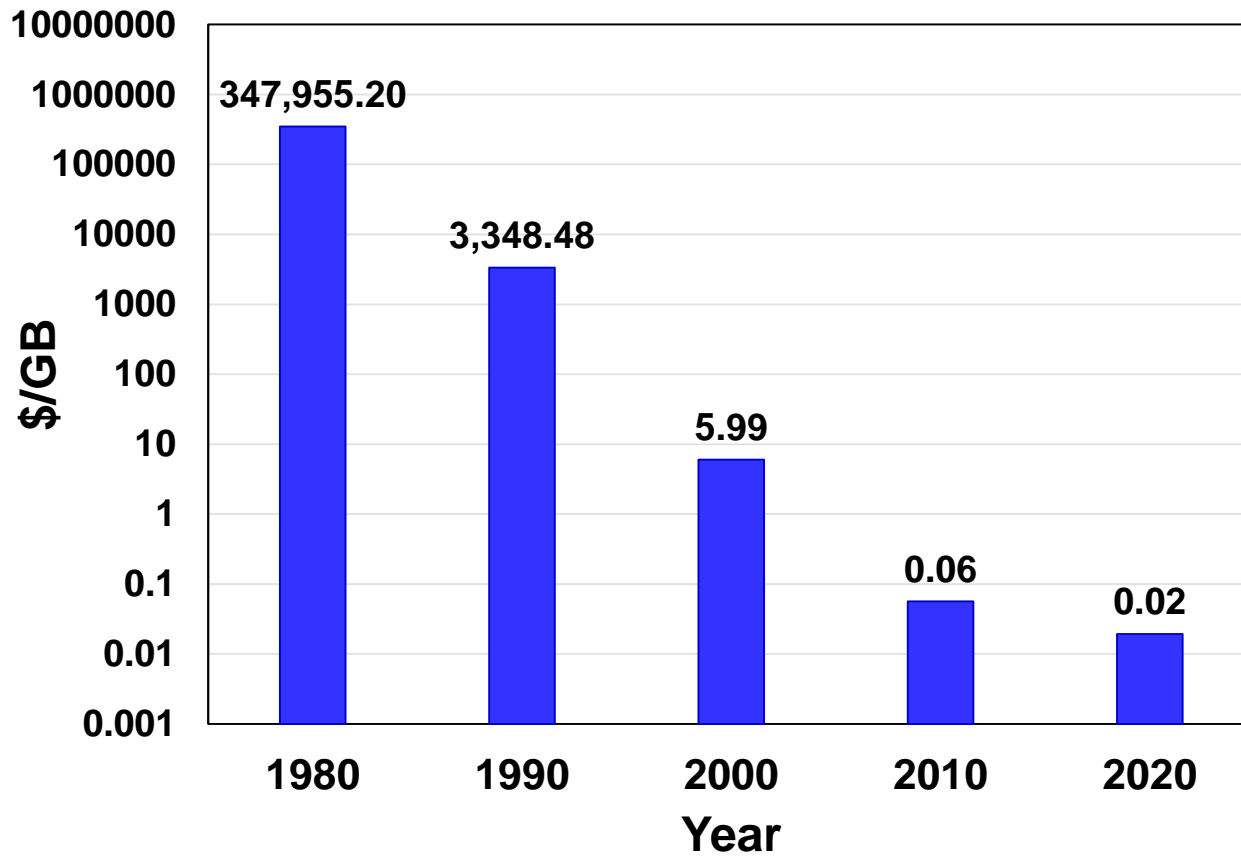
1980s



Today

Storage evolution over time

Storage cost over time





Processes 20 PB a day (2008)
 Crawls 20B web pages a day (2012)
 Search index is 100+ PB (5/2014)
 Bigtable serves 2+ EB, 600M QPS (5/2014)



400B pages,
 10+ PB (2/2014)



19 Hadoop clusters: 600
 PB, 40k servers (9/2015)



150 PB on 50k+ servers
 running 15k apps (6/2011)



Hadoop: 10K nodes, 150K
 cores, 150 PB (4/2014)

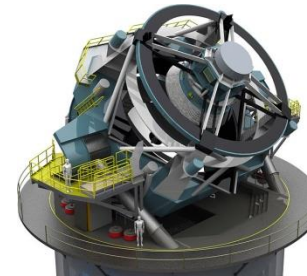
300 PB data in Hive +
 600 TB/day (4/2014)



S3: 2T objects, 1.1M
 request/second (4/2013)



LHC: ~15 PB a year



LSST: 6-10 PB a year
 (~2020)

640K ought to be
 enough for
 anybody.



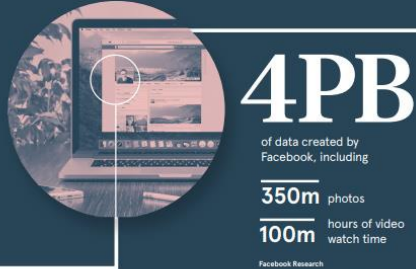
SKA: 0.3 – 1.5 EB
 per year (~2020)



How much data?

A DAY IN DATA

The exponential growth of data is undisputed, but the numbers behind this explosion – fuelled by internet of things and the use of connected devices – are hard to comprehend, particularly when looked at in the context of one day



DEMYSIFYING DATA UNITS

From the more familiar 'bit' or 'megabyte', larger units of measurement are more frequently being used to explain the masses of data

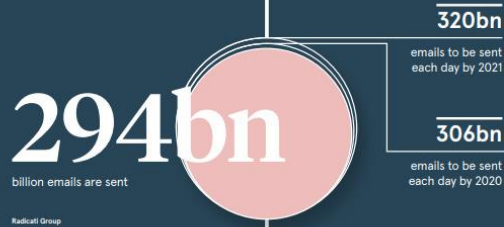
| Unit | Value | Size |
|---------------------|--------------------------|---|
| b bit | 0 or 1 | 1/8 of a byte |
| B byte | 8 bits | 1 byte |
| KB kilobyte | 1,000 bytes | 1,000 bytes |
| MB megabyte | 1,000 ² bytes | 1,000,000 bytes |
| GB gigabyte | 1,000 ³ bytes | 1,000,000,000 bytes |
| TB terabyte | 1,000 ⁴ bytes | 1,000,000,000,000 bytes |
| PB petabyte | 1,000 ⁵ bytes | 1,000,000,000,000,000 bytes |
| EB exabyte | 1,000 ⁶ bytes | 1,000,000,000,000,000,000 bytes |
| ZB zettabyte | 1,000 ⁷ bytes | 1,000,000,000,000,000,000,000 bytes |
| YB yottabyte | 1,000 ⁸ bytes | 1,000,000,000,000,000,000,000,000 bytes |

*A lowercase "b" is used as an abbreviation for bits, while an uppercase "B" represents bytes.

463EB

of data will be created every day by 2025

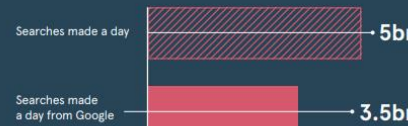
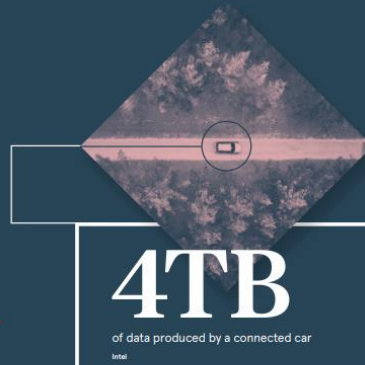
IOC



320bn
emails to be sent each day by 2021

306bn
emails to be sent each day by 2020

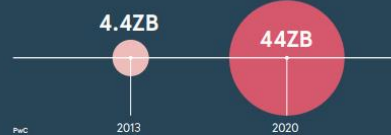
3.9bn
people use emails



28PB
to be generated from wearable devices by 2020

Statista

ACCUMULATED DIGITAL UNIVERSE OF DATA



**Data
generation**



**Storage
cost**



Big Data

Why big data? Science
Business
Society



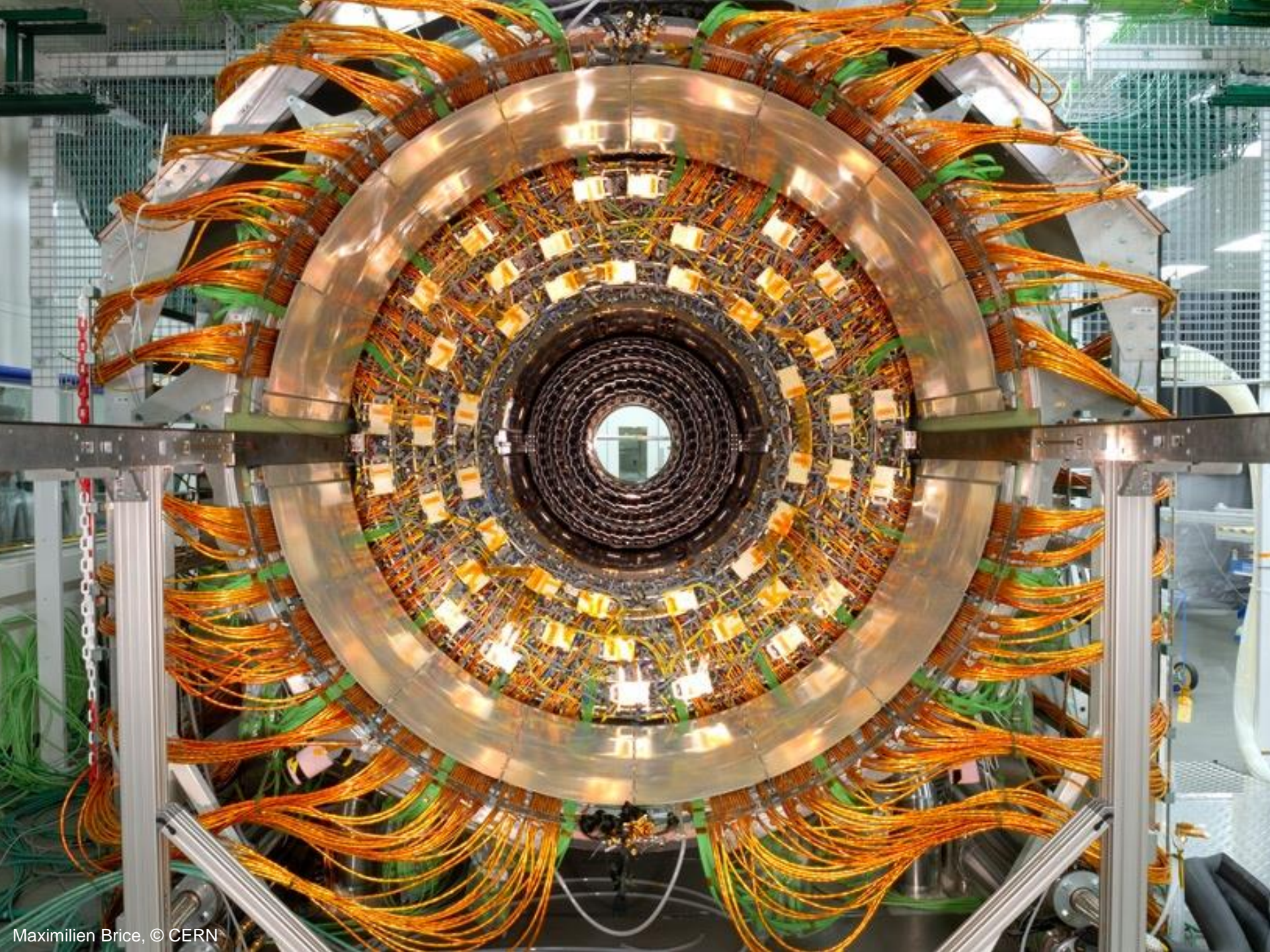


Science

Emergence of the 4th Paradigm

Data-intensive e-Science





Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC

The [ATLAS Collaboration](#)

(Submitted on 31 Jul 2012 (v1), last revised 31 Aug 2012 (this version, v2))

A search for the Standard Model Higgs boson in proton–proton collisions with the ATLAS detector at the LHC is presented. The datasets used correspond to integrated luminosities of approximately 4.8 fb^{-1} collected at $\sqrt{s} = 7 \text{ TeV}$ in 2011 and 5.8 fb^{-1} at $\sqrt{s} = 8 \text{ TeV}$ in 2012. Individual searches in the channels $H \rightarrow ZZ^{(*)} \rightarrow \text{llll}$, $H \rightarrow \gamma \gamma$ and $H \rightarrow WW \rightarrow e \nu \mu \nu$ in the 8 TeV data are combined with previously published results of searches for $H \rightarrow ZZ^{(*)} \rightarrow \text{llll}$ and $H \rightarrow \gamma \gamma$ channels in the 7 TeV data. Clear evidence for the production of a neutral boson with a measured mass of $126.0 \pm 0.4(\text{stat}) \pm 0.4(\text{sys}) \text{ GeV}$ is presented. This observation, which has a significance of 5.9 standard deviations, corresponding to a background fluctuation probability of 1.7×10^{-9} , is compatible with the production and decay of the Standard Model Higgs boson.

Comments: 24 pages plus author list (38 pages total), 12 figures, 7 tables, revised author list, matches version to appear in Physics Letters B

Subjects: **High Energy Physics – Experiment (hep-ex)**

Journal reference: Phys.Lett. B716 (2012) 1–29

DOI: [10.1016/j.physletb.2012.08.020](https://doi.org/10.1016/j.physletb.2012.08.020)

Report number: CERN-PH-EP-2012-218

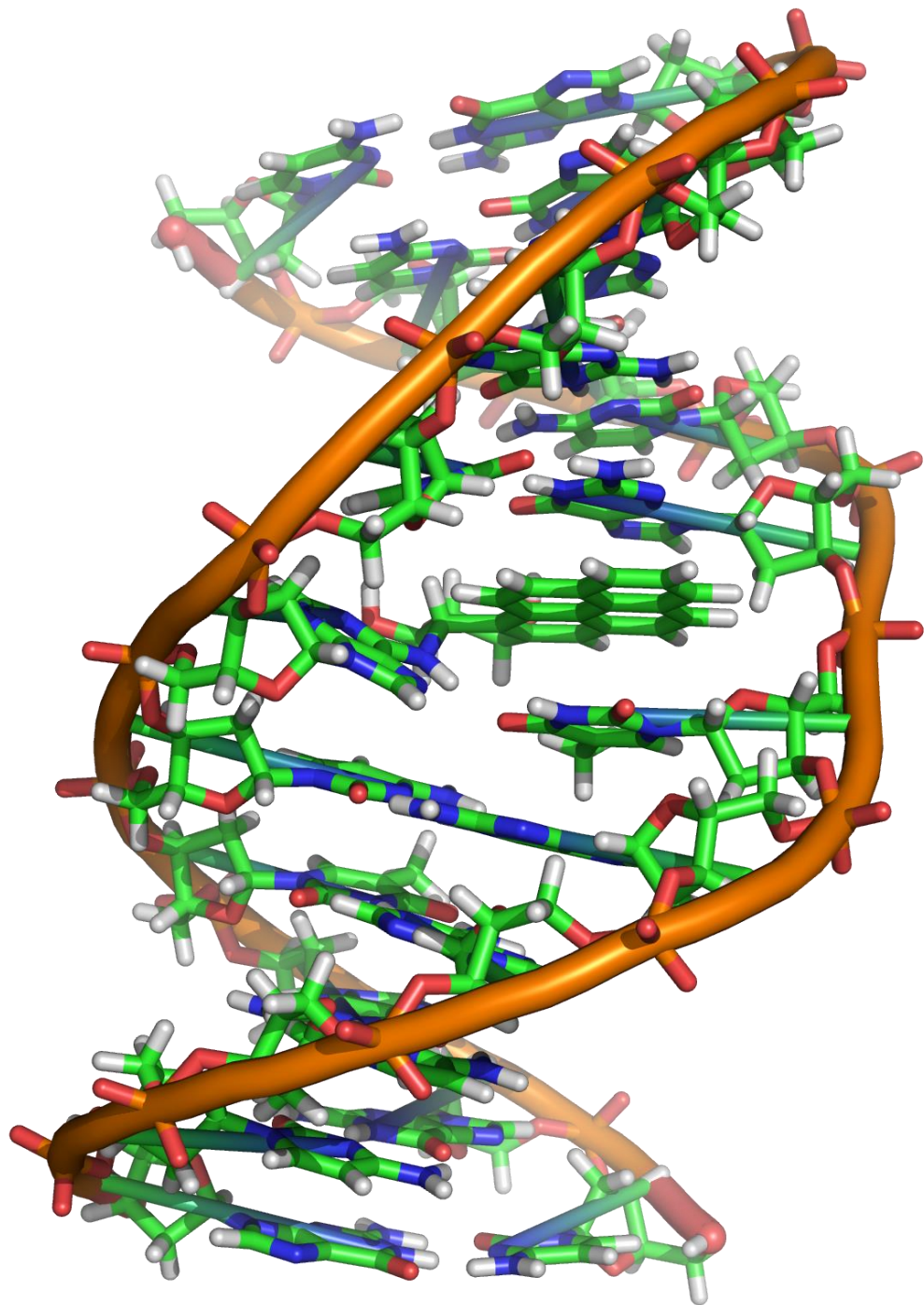
Cite as: [arXiv:1207.7214](https://arxiv.org/abs/1207.7214) [hep-ex]
(or [arXiv:1207.7214v2](https://arxiv.org/abs/1207.7214v2) [hep-ex] for this version)

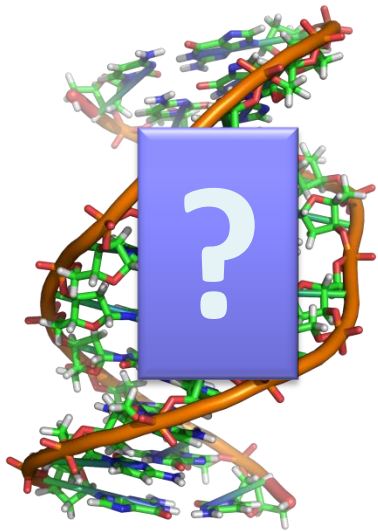
Submission history

From: Atlas Publications [[view email](#)]

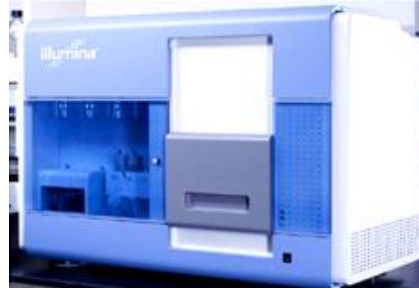
[v1] Tue, 31 Jul 2012 11:59:59 GMT (334kb)

[v2] Fri, 31 Aug 2012 19:29:54 GMT (334kb)





Subject genome



Sequencer

```
GATGCTTACTATGCGGGCCCC
CGGTCTAATGCTTACTATGC
GCTTACTATGCGGGCCCCCTT
AATGCTTACTATGCGGGCCCCCTT
TAATGCTTACTATGC
AATGCTTAGCTATGCGGGC
AATGCTTACTATGCGGGCCCCCTT
AATGCTTACTATGCGGGCCCCCTT
CGGTCTAGATGCTTACTATGC
AATGCTTACTATGCGGGCCCCCTT
CGGTCTAATGCTTAGCTATGC
ATGCTTACTATGCGGGCCCCCTT
```

Reads

Human genome: 3 gbp
A few billion short reads
(~100 GB compressed data)

Business

Data-driven decisions

Data-driven products



Business Intelligence

An organization should retain data that result from carrying out its mission and exploit those data to generate insights that benefit the organization, for example, market analysis, strategic planning, decision making, etc.

Duh!?

This is not a new idea!

In the 1990s, Wal-Mart found that customers tended to buy diapers and beer together. So they put them next to each other and increased sales of both.*

So what's changed?

More compute and storage
Ability to gather behavioral data

* BTW, this is completely apocryphal. (But it makes a nice story.)

Virtuous Product Cycle

a useful service

\$

(hopefully)

transform insights
into action

analyze user behavior
to extract insights

Google. Facebook. Twitter. Amazon. Uber.

data products

data science

net·flix·ing

/ˈnetfliks-ɪŋ/ v

1. The act of watching an entire season of a show in one sitting.
2. A totally valid excuse for avoiding social obligations.

“Sorry, I can’t make it to the party tonight. I am *netflixing*.”





Chinese English Spanish Detect language

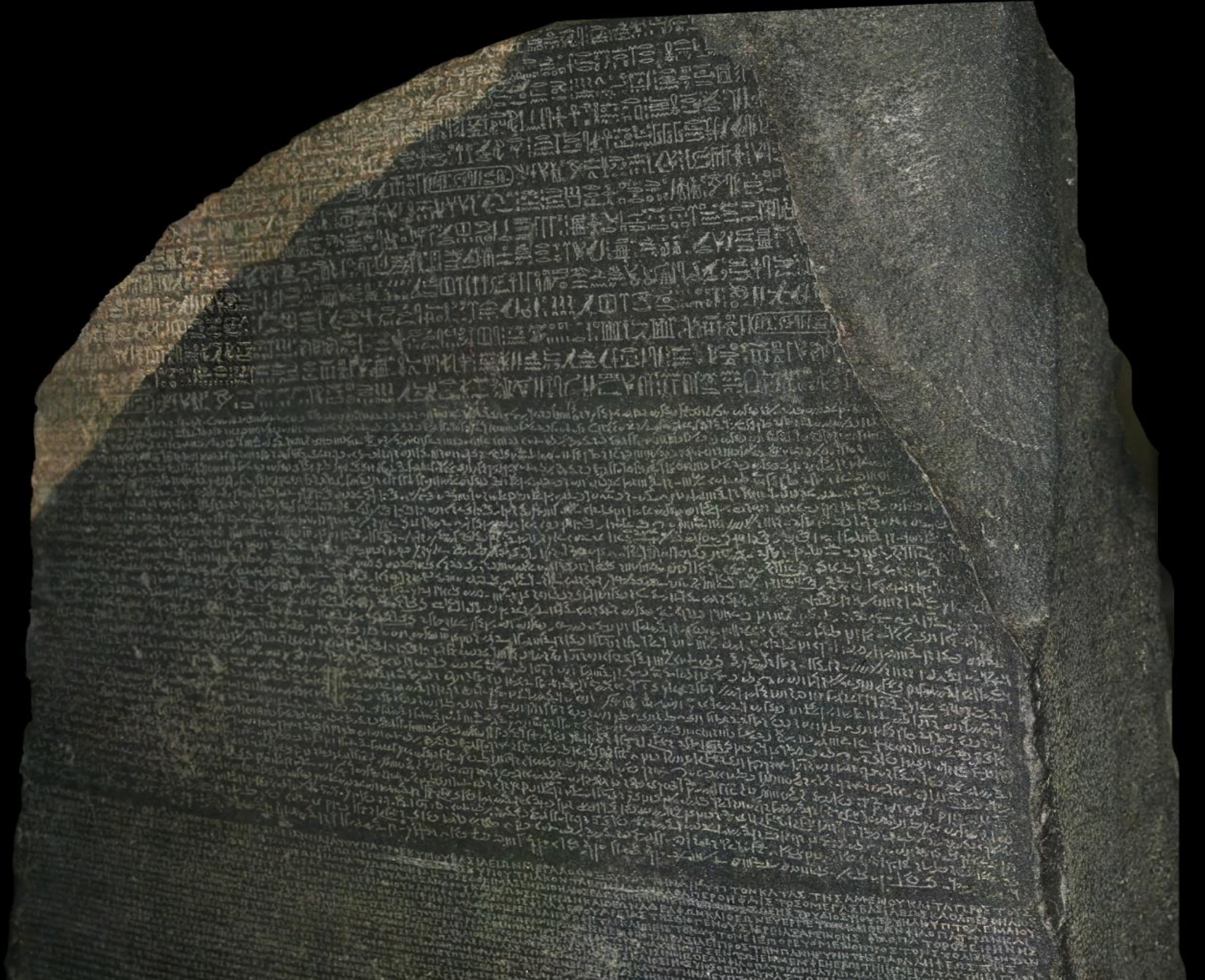
English Chinese (Traditional) Chinese (Simplified)

Translate

How does Google translate English into Chinese?
47/5000

Google如何将英语翻译成中文?
Suggest an edit

Google rúhé jiāng yīngyǔ fānyì chéng zhōngwén?



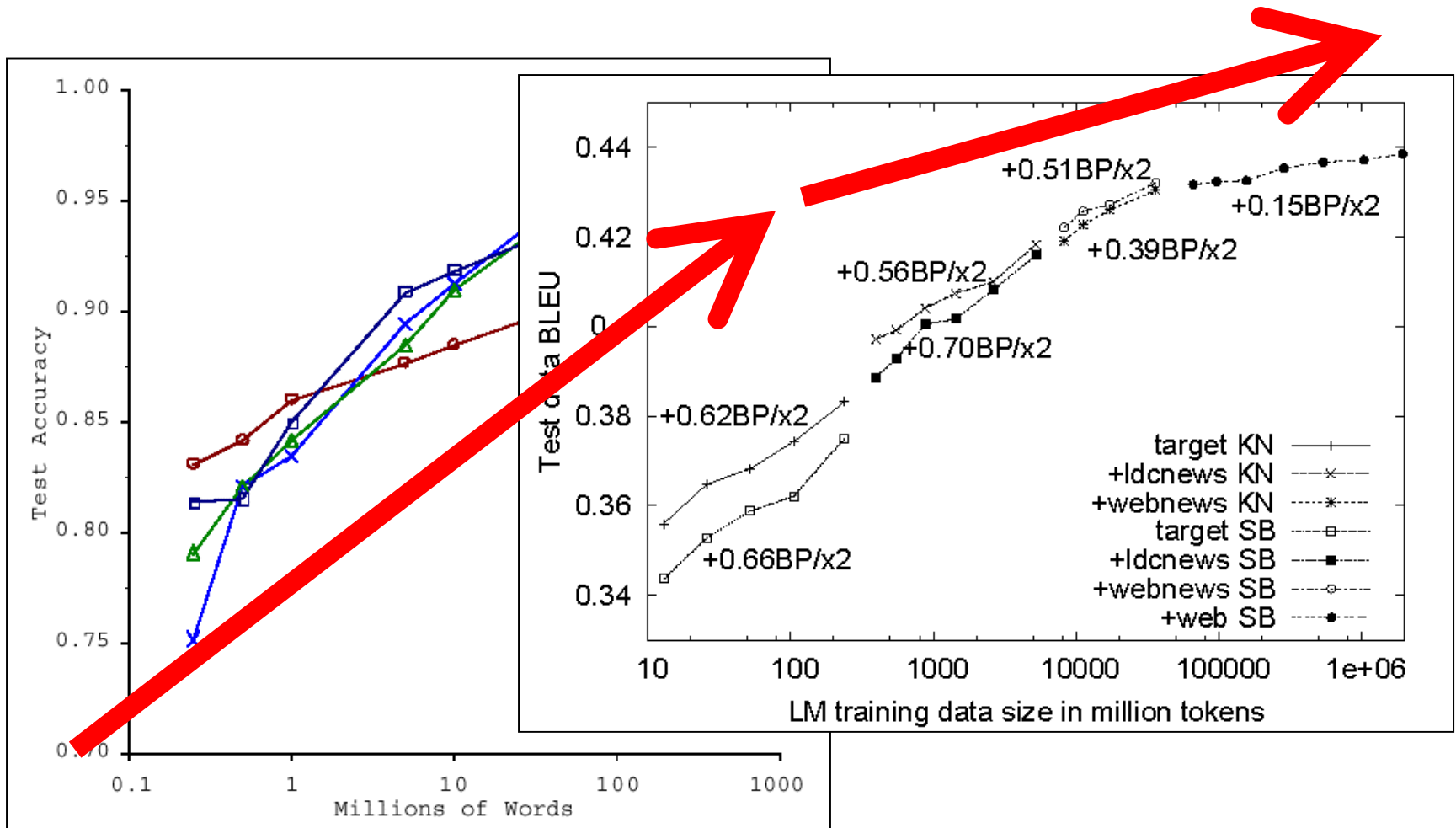
Hieroglyphic text in columns, including the famous cartouches of Ptolemy and Cleopatra. The text is arranged in approximately 15 vertical columns, with some lines containing decorative elements like lotus flowers.

Large block of Greek text in a cursive script, likely a translation of the hieroglyphs above. It consists of many lines of text, with some lines starting with large, decorative initial letters.

Small block of Greek text at the bottom of the fragment, possibly a fragment of a larger inscription or a specific section of the text.

Source: Wikipedia (Rosetta Stone)

No data like more data!



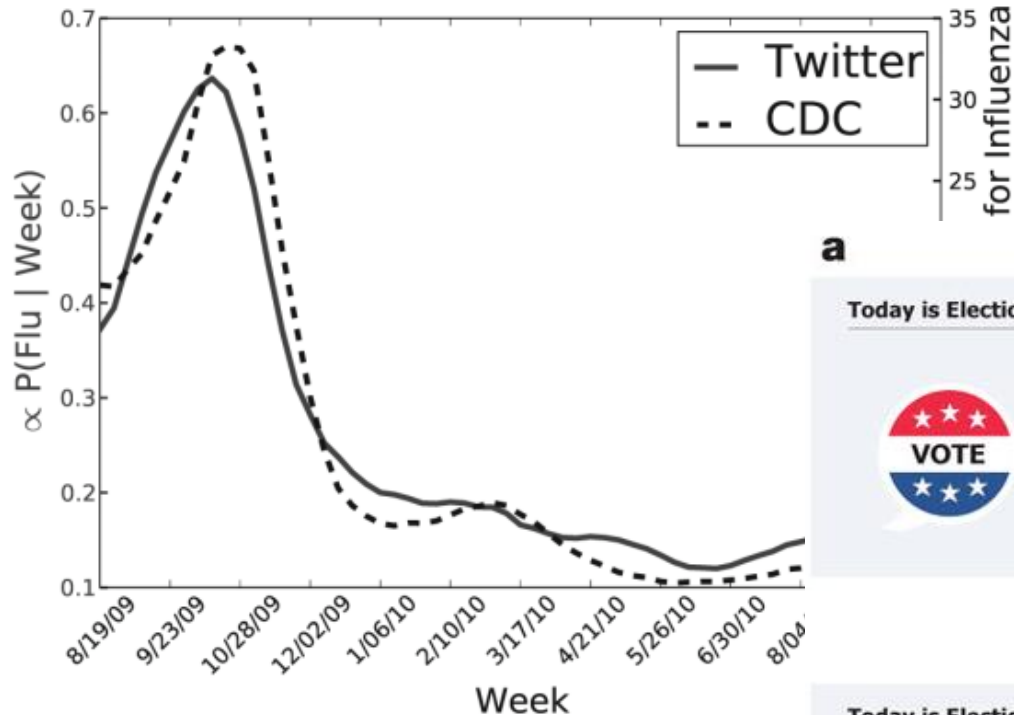
Society

Humans as social sensors

Computational social science



Predicting X with Twitter



a

Informational message



Social message



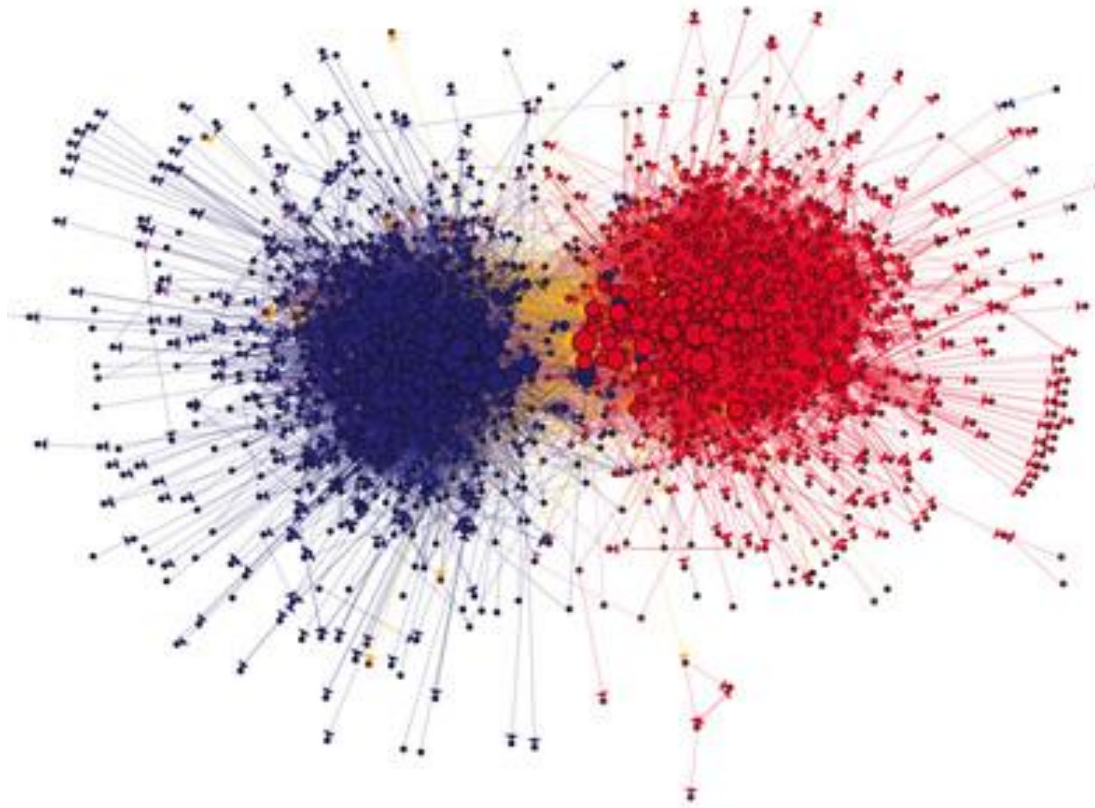
2010 US Midterm Elections:
60m users shown "I Voted" Messages

Summary: increased turnout by
60k directly and 280k indirectly

*Woah! You should feel
unsettled about this!*

Political Mobilization on Facebook

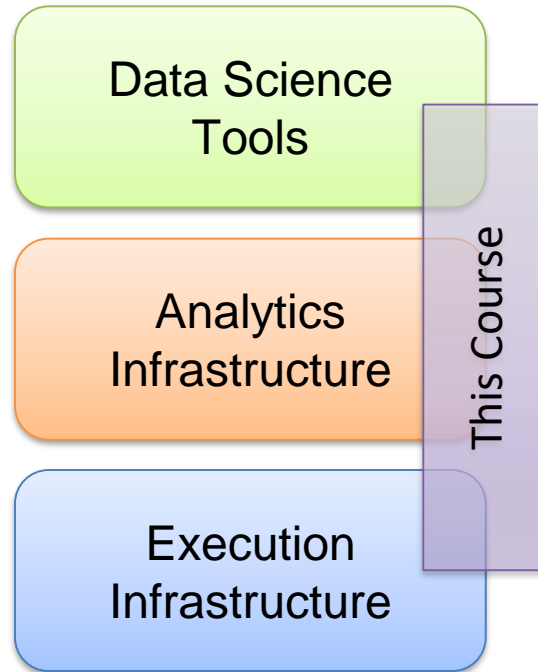
The Political Blogosphere and the 2004 U.S. Election





Source: Popular Internet Meme

What is this course about?

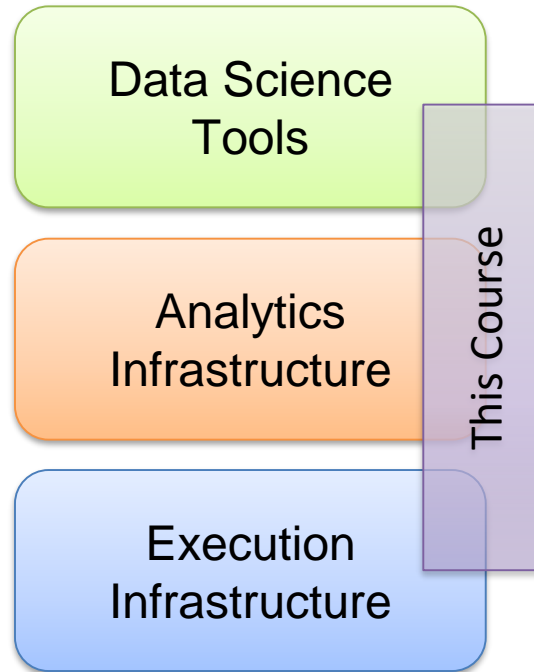


“big data stack”

Buzzwords

data science, data analytics,
business intelligence, data
warehouses and data lakes

MapReduce, Spark, Flink,
Pig, Dryad, Hive, Dryad,
noSQL, Pregel, Giraph,
Storm/Heron



“big data stack”

Text: frequency estimation,
language models, inverted
indexes

Graphs: graph traversals,
random walks (PageRank)

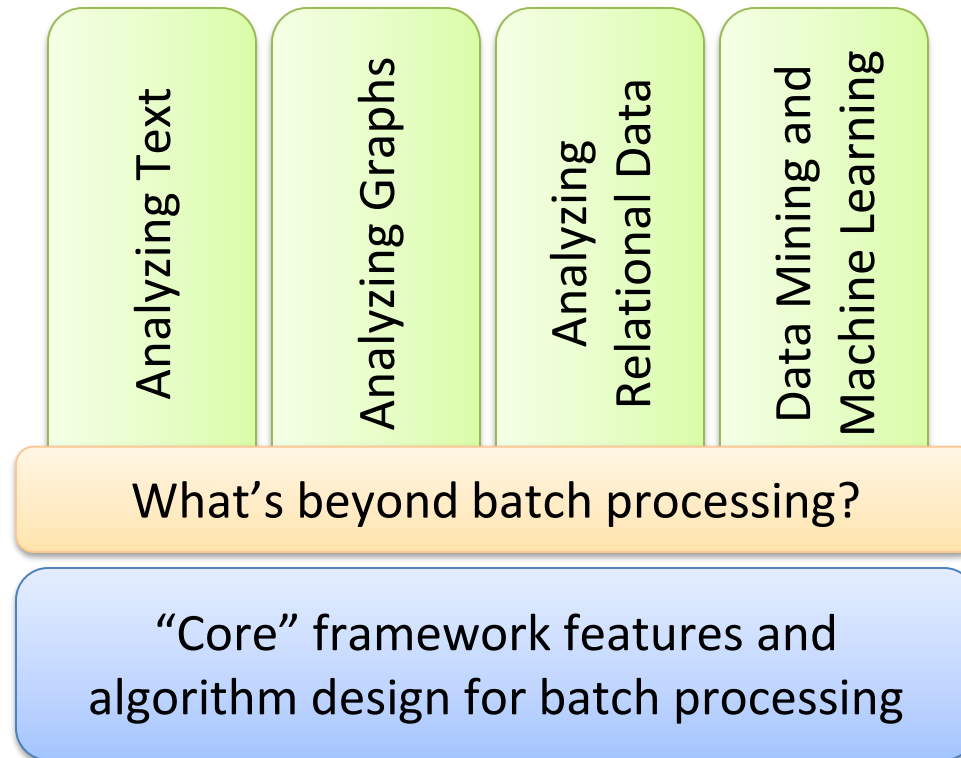
Relational data: SQL, joins,
column stores

Data mining: hashing,
clustering (k -means),
classification,
recommendations

Streams: probabilistic data
structures (Bloom filters,
CMS, HLL counters)

This course focuses on algorithm design and “thinking at scale”

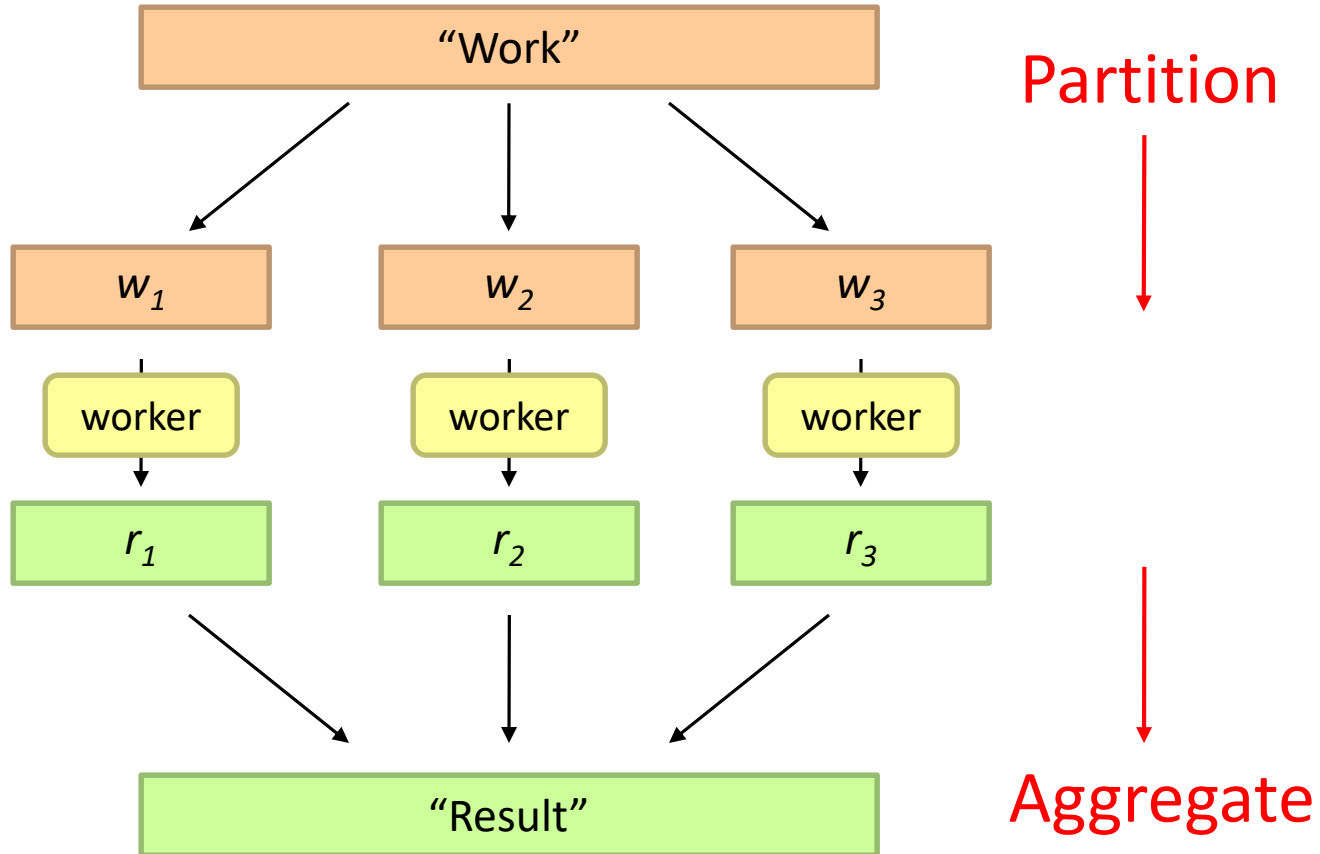
Structure of the Course



Tackling Big Data

A wide-angle, high-angle photograph of a massive server room. The room is filled with rows of server racks, each with numerous glowing lights. The ceiling is a complex network of metal beams and pipes, with several long, rectangular light fixtures hanging from it. The floor is a light-colored, tiled surface. The overall atmosphere is one of a high-tech, industrial environment. The text "Tackling Big Data" is overlaid in the center of the image.

Divide and Conquer



Parallelization Challenges

How do we assign work units to workers?

What if we have more work units than workers?

What if workers need to communicate partial results?

What if workers need to access shared resources?

How do we know when a worker has finished? (Or is simply waiting?)

What if workers die?

Difficult because:

We don't know the order in which workers run...

We don't know when workers interrupt each other...

We don't know when workers need to communicate partial results...

We don't know the order in which workers access shared resources...

What's the common theme of all of these challenges?

Common Theme?

Parallelization challenges arise from:

- Need to communicate partial results

- Need to access shared resources

(In other words, sharing state)

How do we tackle these challenges?

“Current” Tools

Basic primitives

Semaphores (lock, unlock)

Conditional variables (wait, notify, broadcast)

Barriers

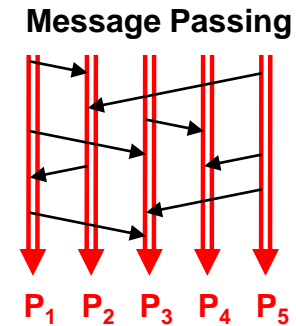
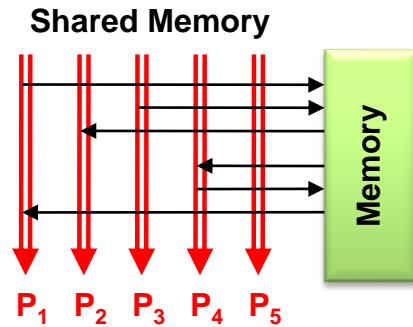
Awareness of Common Problems

Deadlock, livelock, race conditions...

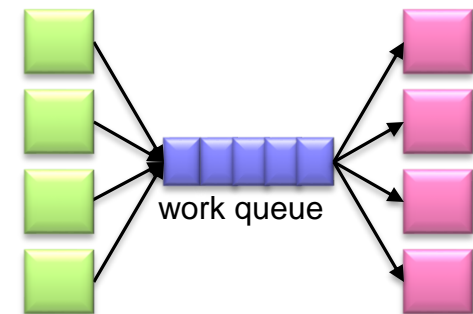
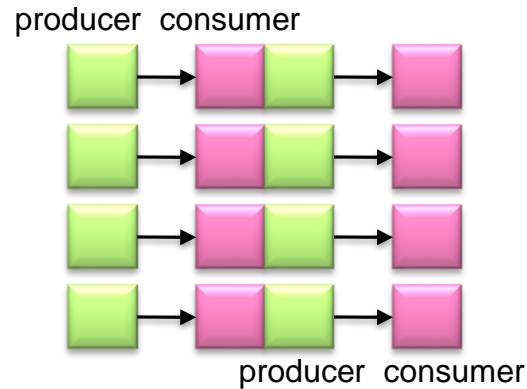
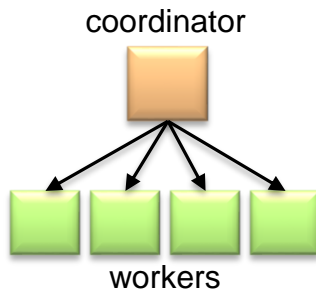
Dining philosophers, sleeping barbers, cigarette smokers...

“Current” Tools

Programming Models



Design Patterns



When Theory Meets Practices

Concurrency is already difficult to reason about...

Now throw in:

The scale of clusters and (multiple) datacenters

The presence of hardware failures and software bugs

The presence of multiple interacting services

The reality:

Lots of one-off solutions, custom code

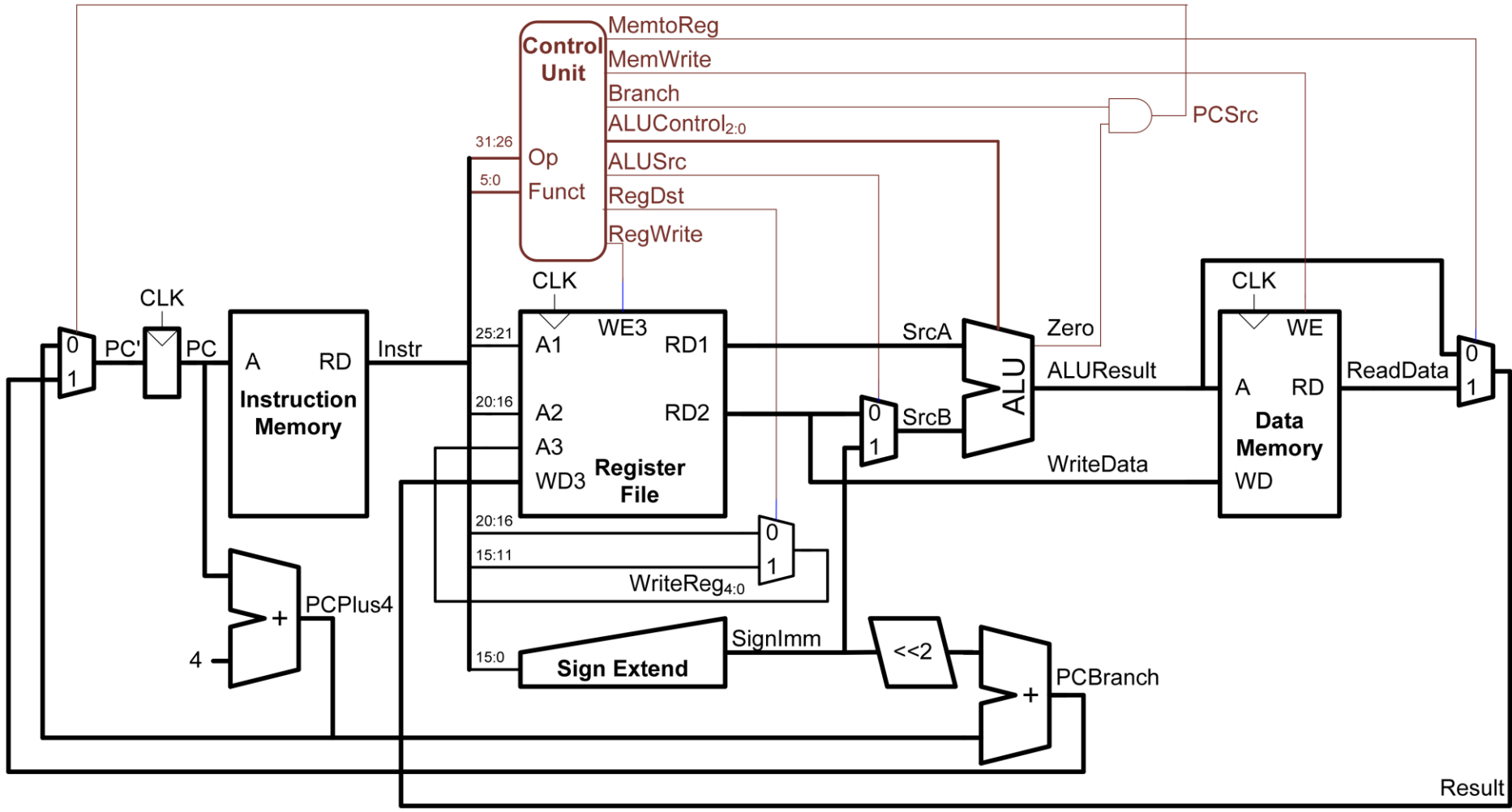
Write you own dedicated library, then program with it

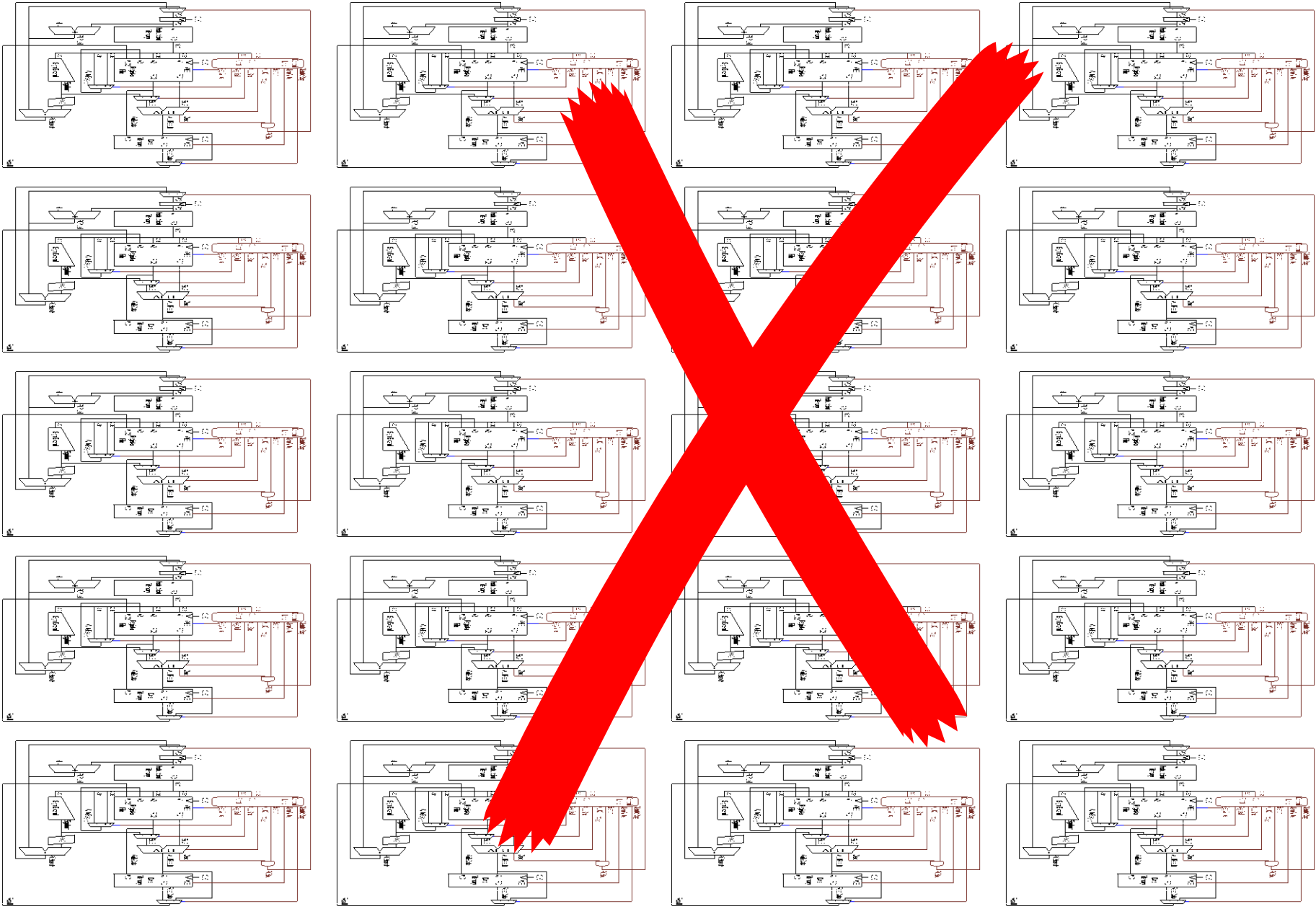
Burden on the programmer to explicitly manage everything

Bottom line: it's hard!



Source: Ricardo Guimarães Herrmann





An aerial photograph of a large industrial datacenter facility during sunset. The sun is a bright orange orb in the upper left, casting a warm glow over the scene. The facility consists of several large, white, rectangular buildings with flat roofs, arranged in a grid-like pattern. A prominent building in the lower right is filled with rows of server racks. A parking lot in the foreground is filled with numerous white semi-trailers. The surrounding landscape is a mix of green fields and brown agricultural land, with some small ponds visible. The sky is a gradient of orange and yellow, transitioning to a pale blue at the horizon.

The datacenter *is* the computer!

The datacenter *is* the computer!

It's all about the right level of abstraction

Moving beyond the von Neumann architecture

What's the "instruction set" of the datacenter computer?

Hide system-level details from the developers

No more race conditions, lock contention, etc.

No need to explicitly worry about reliability, fault tolerance, etc.

Separating the *what* from the *how*

Developer specifies the computation that needs to be performed

Execution framework ("runtime") handles actual execution

MapReduce is the first instantiation of this idea... but not the last!



Questions?



Course Administrivia

Four in One!

CS 451/651 431 631 all meet together

CS 451: version for undergrads (most students)

CS 651: version for grad students

CS 431: version for undergrads

CS 631: version for grad students

Two in One!

CS 451/651

CS 451: version for CS undergrads (most students)

CS 651: version for CS grads

CS 431/631

CS 431: version for non-CS undergrads

CS 631: version for non-CS grads

Course instructors

Ali Abedi: The guy talking right now

TAs: Ryan Clancy, Zheng Ma, Yuqing Xie, Wei Tu, Abdul Naik

Important Coordinates

Course website:

<https://www.student.cs.uwaterloo.ca/~cs451/>

Lots of info there, read it!

("I didn't see it" will not be accepted as an excuse)

Communicating with us:

[Piazza for general/private questions \(link on course homepage\)](#)

Bespin

<http://bespin.io/>

Course Design

This course focuses on algorithm design and “thinking at scale”

Not the “mechanics” (API, command-line invocations, et.)
You’re expected to pick up MapReduce/Spark with minimal help

Components of the final grade:

6 (CS 431/631) or 8 (CS 451/651) individual assignments

Final exam

Additional group final project (CS 631/651)

Expectations (CS 451)

Your background:

Pre-reqs: CS 341, CS 348, CS 350

Comfortable in Java and Scala (or be ready to pick it up quickly)

Know how to use Git

Reasonable “command-line”-fu skills

Experience in compiling, patching, and installing open source software

Good debugging skills

You are:

Genuinely interested in the topic

Be prepared to put in the time

Comfortable with rapidly-evolving software

MapReduce/Spark Environments (CS 451)

See “Software” page in course homepage for instructions

Single-Node Hadoop: Linux Student CS Environment

Everything is set up for you, just follow instructions

We'll make sure everything works

Single-Node Hadoop: Local installations

Install all software components on your own machine

Requires at least 4GB RAM and plenty of disk space

Works fine on Mac and Linux, YMMV on Windows

Important: For your convenience only!

We'll provide basic instructions, but not technical support

Distributed Hadoop: Datasci Cluster

Assignment Mechanics (CS 451)

We'll be using private GitHub repos for assignments

Complete your assignments, push to GitHub
We'll pull your repos at the deadline and grade

Note late policy (details on course homepage)

Late by up to 24 hours: 25% reduction in grade
Late 24-48 hours: 50% reduction in grade
Late by more the 48 hours: not accepted

By assumption, we'll pull and mark at deadline:
If you want us to hold off, you must let us know!

Important: Register for (free) GitHub educational account!

https://education.github.com/discount_requests/new

Assignment Mechanics (CS 431)

Assignments will use Python and Jupyter

Everything you need to know is in the assignment itself

Assignments will generally be submitted using Git

Details are on the course website for the appropriate assignment

Assignment Mechanics (CS 431)

Note late policy (details on course homepage)

Late by up to 24 hours: 25% reduction in grade

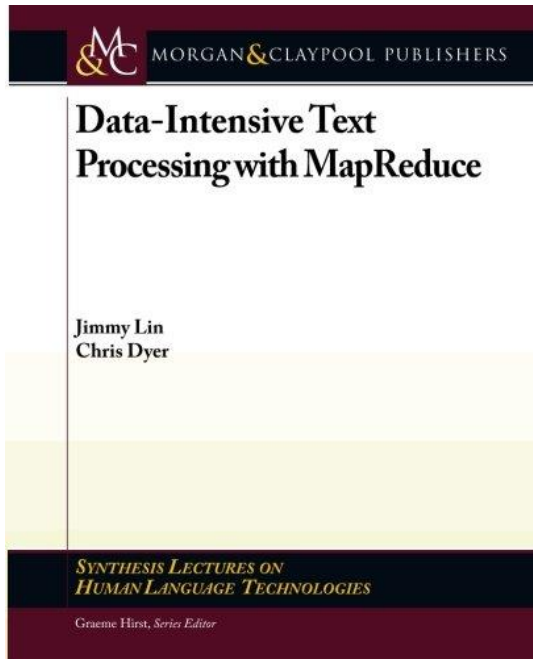
Late 24-48 hours: 50% reduction in grade

Late by more the 48 hours: not accepted

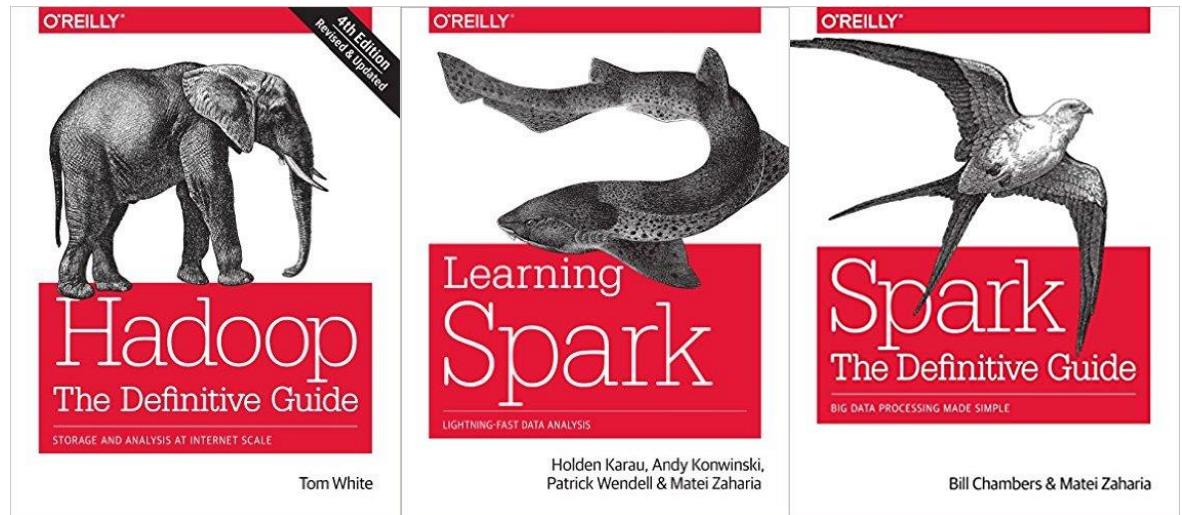
By assumption, we'll pull and mark at deadline:
If you want us to hold off, you must let us know!

Course Materials

One (required) textbook +
Three (optional but recommended) books +
Additional readings from other sources as appropriate



(optional but recommended)



Note: 4th Edition

If you're not (yet) registered:

Register for the wait list at:

By sending Ali an email at ali.abedi@uwaterloo.ca

Priority for unregistered students

CS students

Have all the pre-reqs

Final opportunity to take the course (e.g., 4B students)

Continue to attend class until final decision

Once the course is full, it is *full*

Note: late registration is not an excuse for late assignments