

CS4225/CS5425 Big Data Systems for Data Science

Logistics + Introduction

Bryan Hooi
School of Computing
National University of Singapore
bhooi@comp.nus.edu.sg



Schedule

Week	Date	Topics	Tutorial	Due Dates
1	18 Aug	Overview and Introduction		
2	25 Aug	MapReduce - Introduction		
3	1 Sep	Polling Day Public Holiday		
4	8 Sep	MapReduce and Databases	Tutorial: MapReduce	Assignment 1 released
5	15 Sep	NoSQL Overview 1		
6	22 Sep	NoSQL Overview 2	Tutorial: NoSQL	
Recess				
7	6 Oct	Apache Spark 1		
8	13 Oct	Apache Spark 2	Tutorial: Spark	Assignment 1 due (15 Oct 11.59pm), Assignment 2 released
9	20 Oct	Stream Processing 1		
10	27 Oct	Stream Processing 2	Tutorial: Stream Processing	
11	3 Nov	Large Graph Processing 1		
12	10 Nov	NUS Well-Being Day		
13	17 Nov	Large Graph Processing 2	Tutorial: Graph Processing	Assignment 2 due (19 Nov 11.59pm)
	29 Nov	Final Exam		

Assessment

- 2 assignments on Hadoop and Spark (25% each)
- Final exam (50%) – held in-person.
- (Note: no weightage for in-lecture-quizzes)

Tutorials

- Start from Week 4
- Not counted for final grade
- Tutorial questions will be available on Canvas before the tutorial: we recommend attempting questions before tutorial
- Some questions are samples for tests

Final Exam

- Date: 29 Nov
 - Held in class; open book + notes, but no electronics usage
- Focus is on understanding and application, not facts / memorization
- Example questions
 - **Integrative:** Require you to combine knowledge from different chapters of the textbook
 - **Application:** Require you to apply your knowledge of fundamental concepts to reasonably practical scenarios.
 - **“Why not”:** Example: Tommy proposed a solution A to solve problem B. Explain the problem with solution A and how to overcome this problem

What will we learn?

- **We will learn to process different types of data:**
 - Large data volume
 - Graph data
 - Stream data (infinite/never-ending)
- **We will learn to use different models of computation:**
 - MapReduce/Spark
 - Large graph processing engines
 - Streams and online algorithms

What is Data Science?

- **Standard definition:**

“Data science is an **interdisciplinary** field about processes and systems to extract **knowledge or insights** from data in various forms.”

What is Data Science?

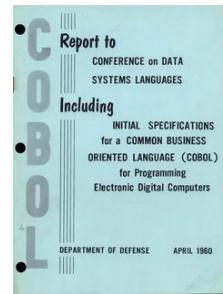
- **Standard definition:**

“Data science is an **interdisciplinary** field about processes and systems to extract **knowledge or insights** from data in various forms.”

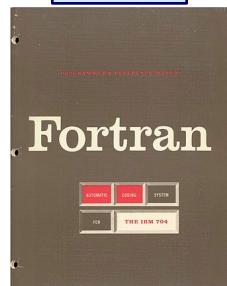
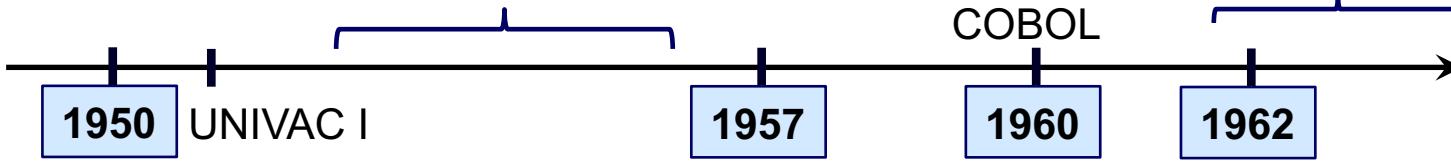
- **Historical view:** in 1962, statistician John Tukey described a field called “data analysis”, which emerged out of statistics, but focusing not just on testing hypotheses, but on **computing aspects of data analysis**: including collecting, storing, analyzing, presenting data, etc.



Early transistor computers (TX-0, IBM 608, 7000 series)



Early DBMS, statistical software



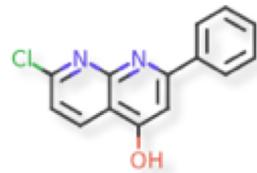
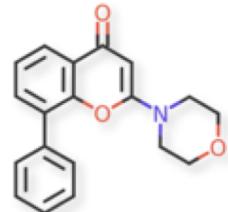
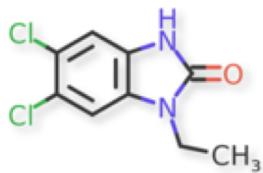
THE FUTURE OF DATA ANALYSIS¹
By JOHN W. TUKEY
Princeton University and Bell Telephone Laboratories



Q: What rule characterizes toxic molecules?

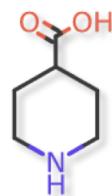
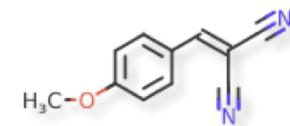
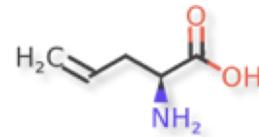


Toxic

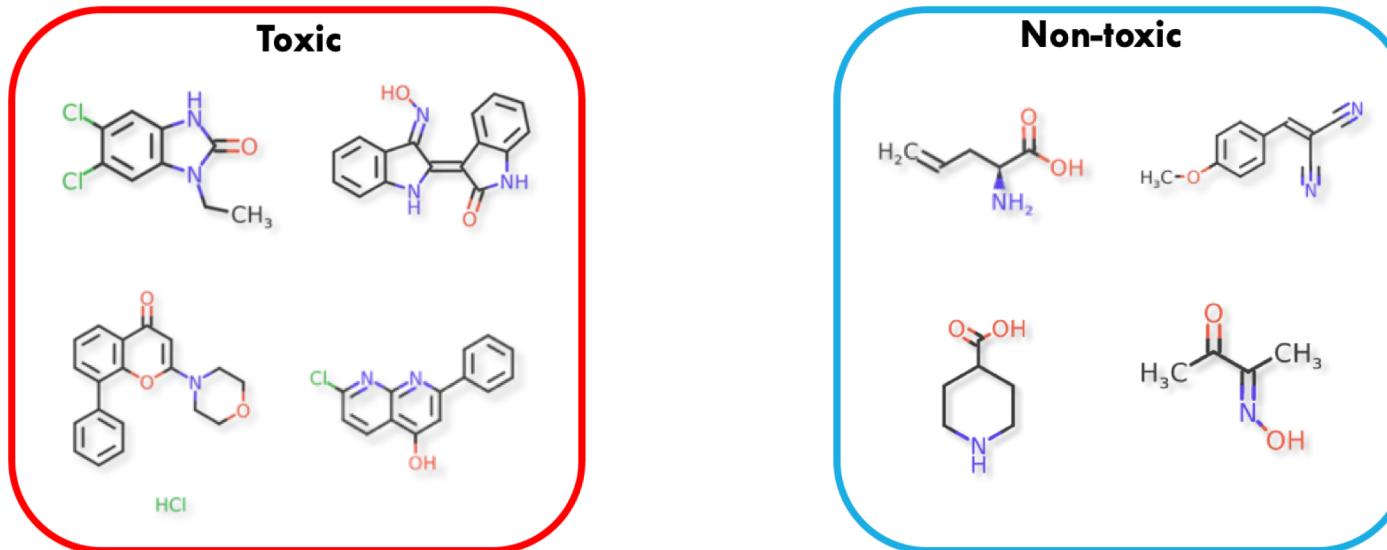


HCl

Non-toxic



Q: What rule characterizes toxic molecules?

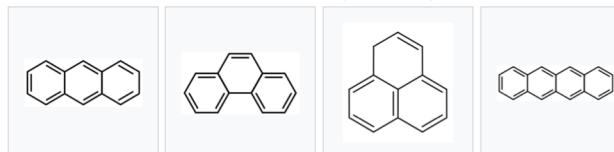


if > 1:
toxic
else
nontoxic

Polycyclic aromatic hydrocarbon

From Wikipedia, the free encyclopedia

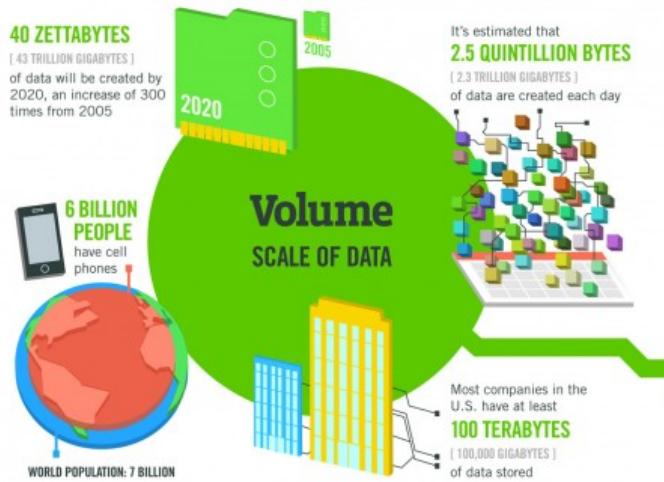
Principal PAH Compounds



Cancer [edit]

PAHs have been linked to [skin](#), [lung](#), [bladder](#), [liver](#), and [stomach](#) cancers in well-established animal model studies.^[72] Human carcinogens are identified in the section "Regulation and Oversight" below.

Challenges of Big Data: the 4 'V's



The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume**, **Velocity**, **Variety** and **Veracity**.

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015
4.4 MILLION IT JOBS will be created globally to support big data, with 1.9 million in the United States



The New York Stock Exchange captures
1 TB OF TRADE INFORMATION during each trading session



Modern cars have close to **100 SENSORS** that monitor items such as fuel level and tire pressure

Velocity ANALYSIS OF STREAMING DATA

By 2016, it is projected there will be

18.9 BILLION NETWORK CONNECTIONS

— almost 2.5 connections per person on earth



As of 2011, the global size of data in healthcare was estimated to be

150 EXABYTES [161 BILLION GIGABYTES]



30 BILLION PIECES OF CONTENT

are shared on Facebook every month



By 2014, it's anticipated there will be
420 MILLION WEARABLE, WIRELESS HEALTH MONITORS

4 BILLION+ HOURS OF VIDEO are watched on YouTube each month



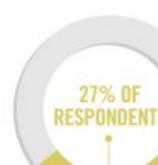
Variety DIFFERENT FORMS OF DATA



400 MILLION TWEETS are sent per day by about 200 million monthly active users

1 IN 3 BUSINESS LEADERS

don't trust the information they use to make decisions



27% OF RESPONDENTS in one survey were unsure of how much of their data was inaccurate

Poor data quality costs the US economy around

\$3.1 TRILLION A YEAR



Veracity UNCERTAINTY OF DATA

Veracity: ~3.3% of data in some of the most popular datasets are mislabelled

Pervasive Label Errors in Test Sets Destabilize Machine Learning Benchmarks

Curtis G. Northcutt*
ChipBrain, MIT

Anish Athalye
MIT

Jonas Mueller
Amazon



given: cat
corrected: frog



given: lobster
corrected: crab



given: ewer
corrected: teapot



given: white stork
corrected: black stork

[1] Northcutt, Curtis G., Anish Athalye, and Jonas Mueller. "Pervasive Label Errors in Test Sets Destabilize Machine Learning Benchmarks." NeurIPS 2021

Veracity: Importance of Data Quality

Forbes

ENTERPRISE & CLOUD

Andrew Ng Launches A Campaign For Data-Centric AI

Data is eating the world so Andrew Ng wants to make sure we radically improve its quality. “Data is food for AI,” says Ng, and he is launching a campaign to shift the focus of AI practitioners from model/algorithm development to the quality of the data they use to train the models.



The background of the image is a vast expanse of white and grey clouds against a clear blue sky. The clouds are dense and layered, creating a sense of depth. In the lower right quadrant, there are darker, more solid-looking clouds that appear to be overhanging land or water. The overall scene is one of a peaceful flight through the atmosphere.

Infrastructure: Cloud Computing

Utility Computing

- What?

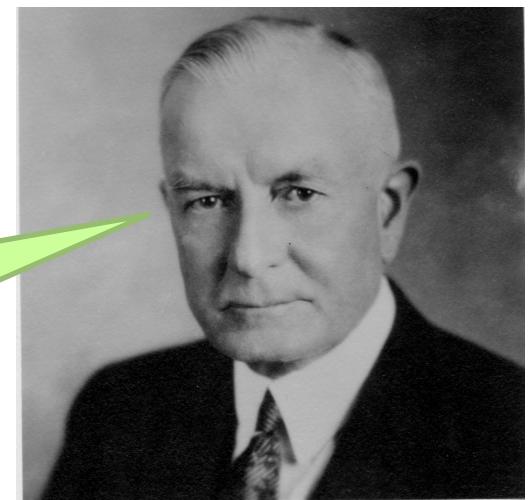
- Computing resources as a metered service (“pay as you go”)
- Ability to dynamically provision virtual machines

- Why?

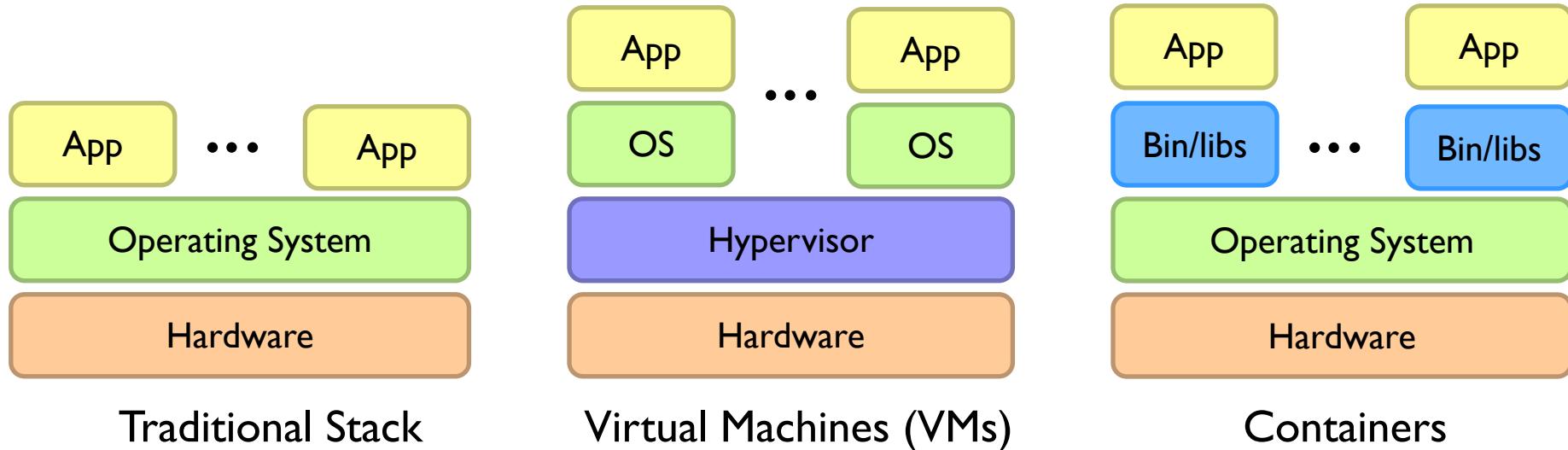
- Scalability: “infinite” capacity
- Elasticity: scale up or down on demand

I think there is a world market for about five computers.

Thomas J. Watson (attributed?)



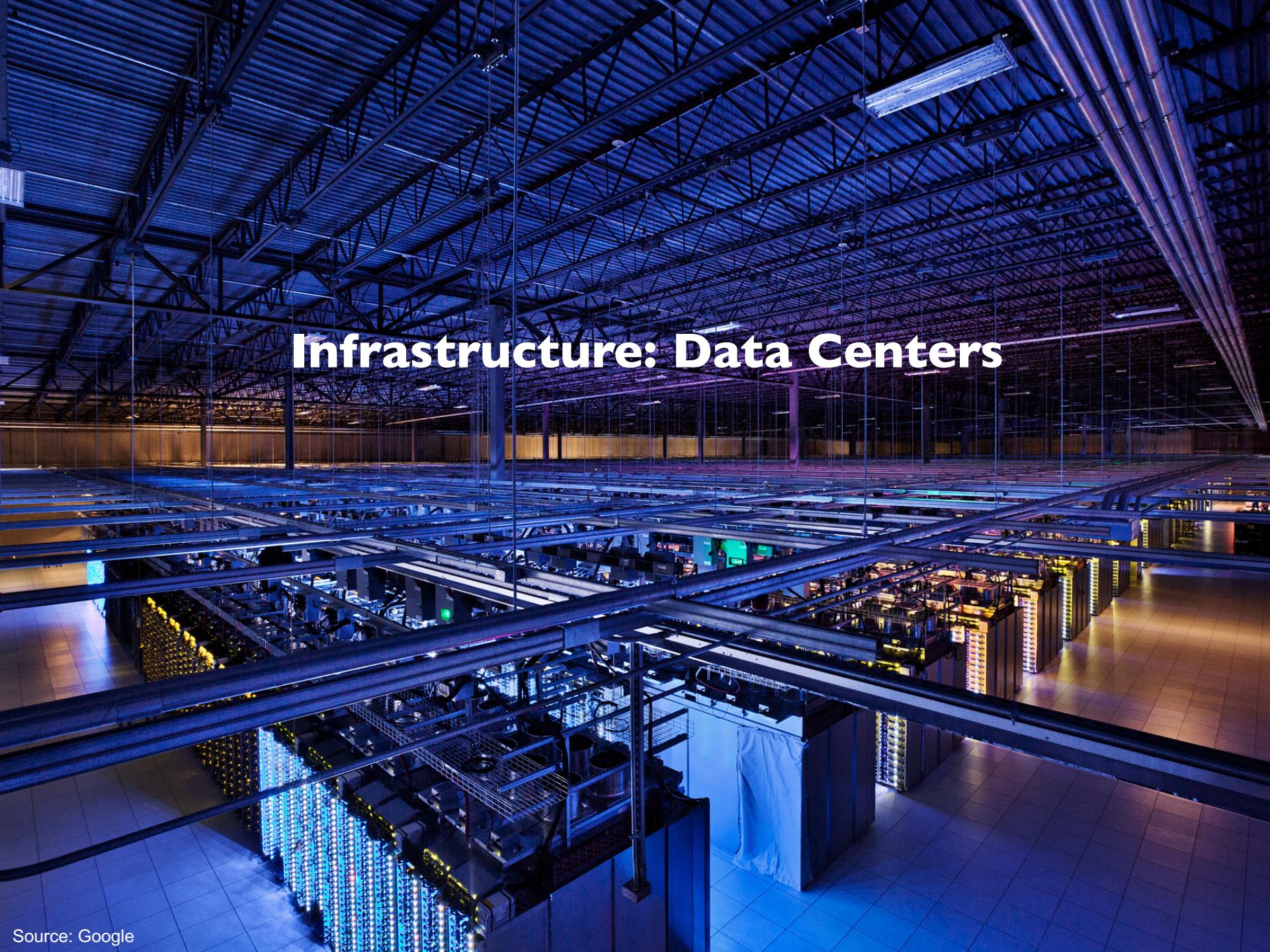
Enabling Technology: Virtualization and Containers



- **Virtual Machines:** enable sharing of hardware resources by running each application in an isolated virtual machine.
 - **High overhead** as each VM has its own OS.
- **Containers:** enable lightweight sharing of resources, as applications run in an isolated way, but still share the same OS.
 - A container is a **lightweight software package** that encapsulates an application and its environment.

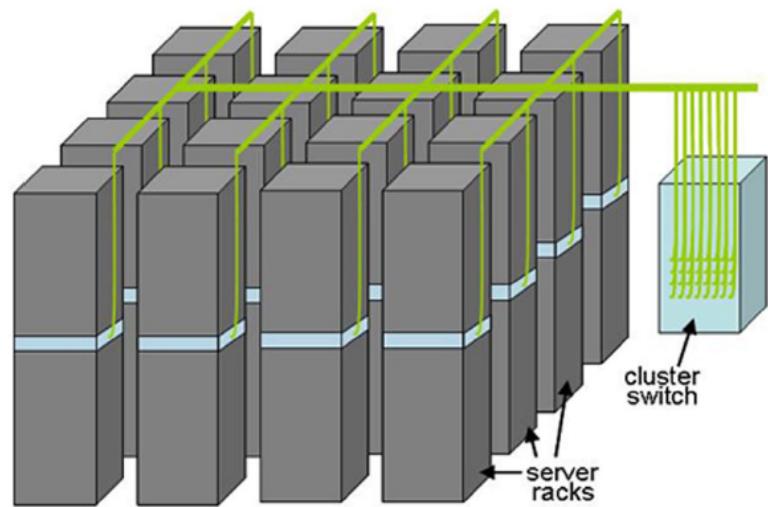
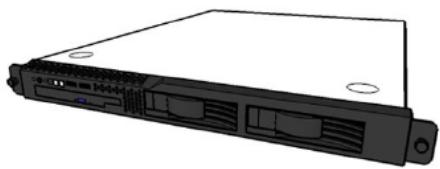
Everything as a Service

- Infrastructure as a Service (IaaS): Utility Computing
 - Why buy machines when you can rent cycles?
 - Examples: Amazon's EC2, Rackspace, Google Compute Engine
- Platform as a Service (PaaS)
 - Provides hosting for web applications and takes care of the hardware maintenance, upgrades, ...
 - Example: Google App Engine
- Software as a Service (SaaS)
 - Just run it for me!
 - Example: Gmail, Dropbox, Zoom



Infrastructure: Data Centers

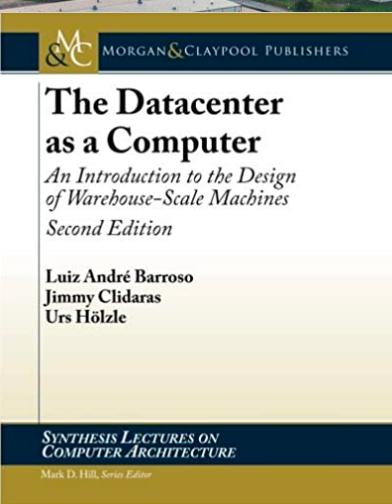
Building Blocks



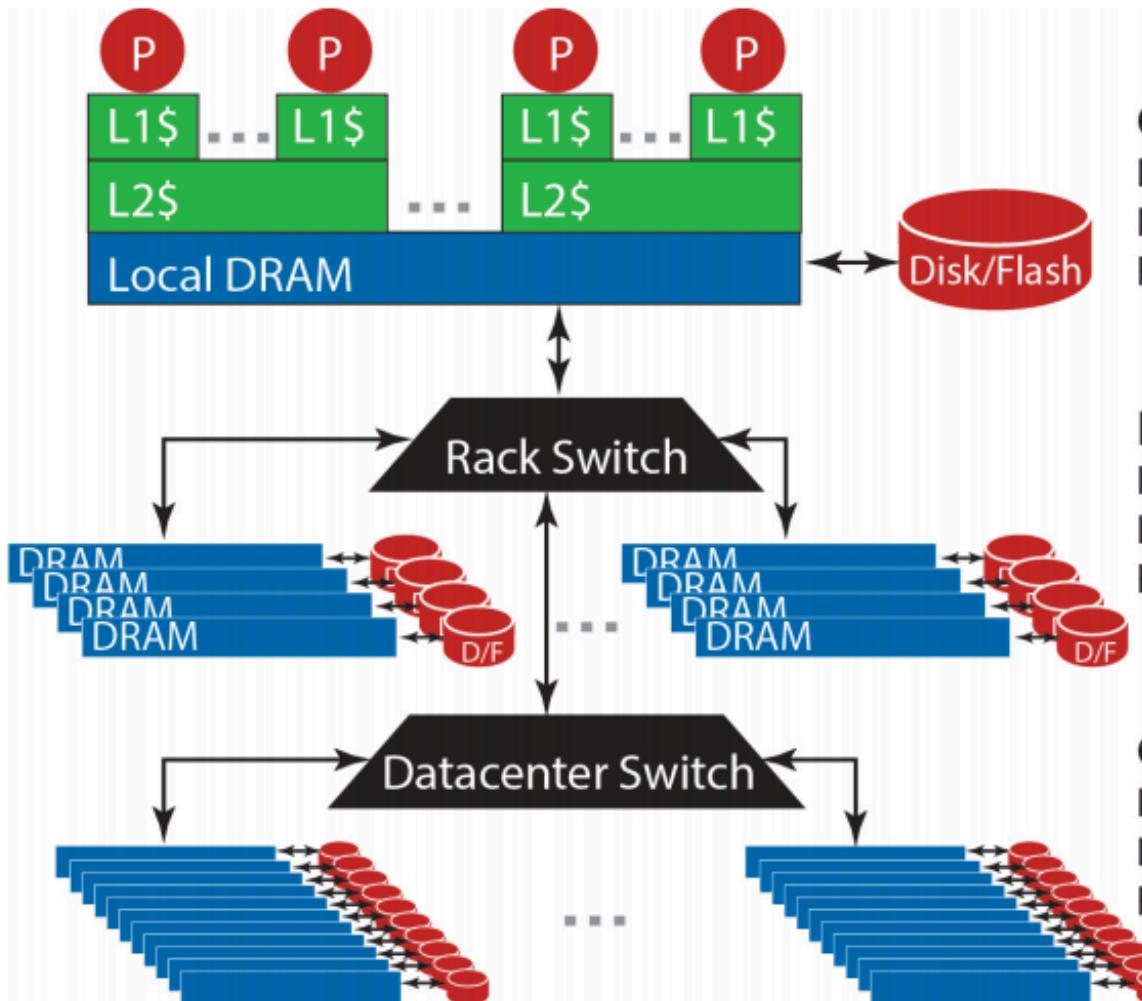


Key Idea: the datacenter is the computer

(Meaning: we can think of the machines in a data center as a big “processing unit” being used to solve a problem, rather than just as independent machines)



Storage Hierarchy



One Server

DRAM: 16 GB, 100 ns, 20 GB/s
Disk: 2TB, 10 ms, 200 MB/s
Flash: 128 GB, 100 us, 1 GB/s

Local Rack (80 servers)

DRAM: 1 TB, 300 us, 100 MB/s
Disk: 160 TB, 11 ms, 100 MB/s
Flash: 20 TB, 400 us, 100 MB/s

Cluster (30 racks)

DRAM: 30 TB, 500 us, 10 MB/s
Disk: 4.80 PB, 12 ms, 10 MB/s
Flash: 600 TB, 600 us, 10 MB/s

Bandwidth vs Latency

- **Bandwidth:** maximum amount of data that can be transmitted per unit time (e.g. in GB/s)
- **Latency:** time taken for 1 packet to go from source to destination (*one-way*) or from source to destination back to source (*round trip*), e.g. in ms
- When transmitting a *large* amount of data, bandwidth tells us roughly how long the transmission will take.
- When transmitting a *very small* amount of data, latency tells us how much delay there will be.
- **Throughput** is similar to bandwidth, but instead of referring to capacity, it refers to the rate at which some data was *actually transmitted* across the network during some period of time.



Low Bandwidth

High Bandwidth



Low Latency

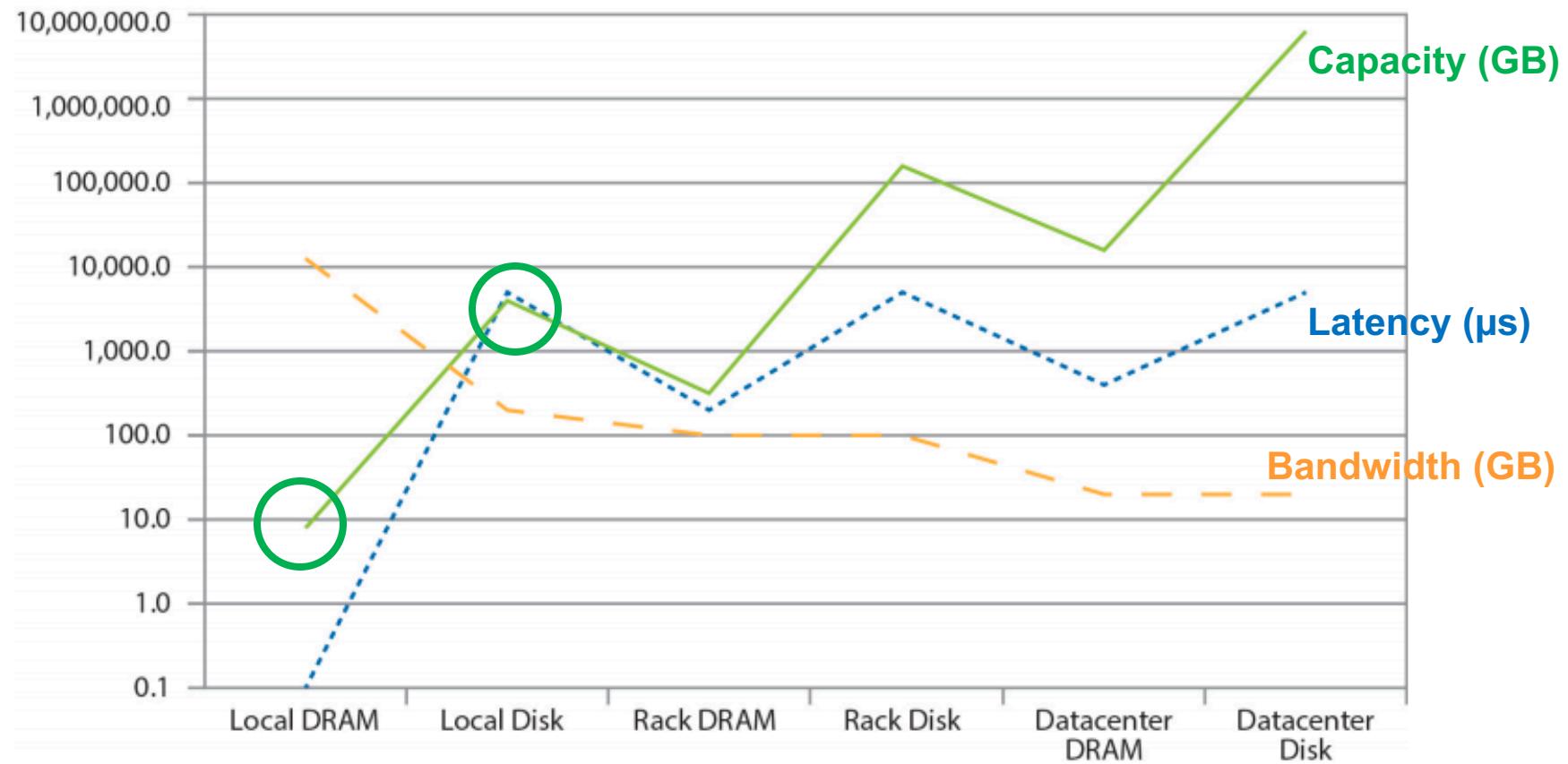


High Latency



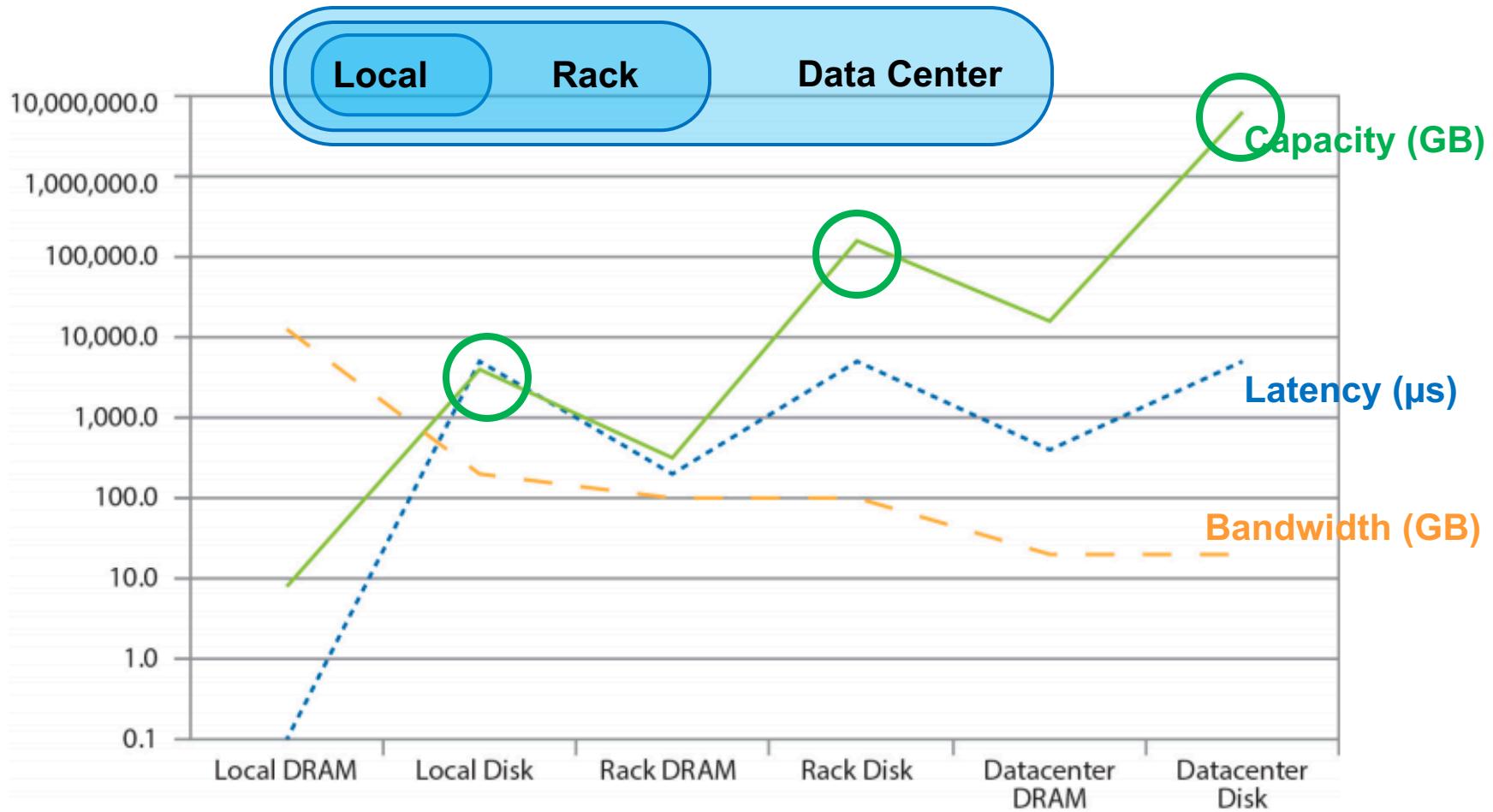
Adding many lanes to a highway: increases bandwidth, but does not decrease latency

Capacity of Storage Hierarchy



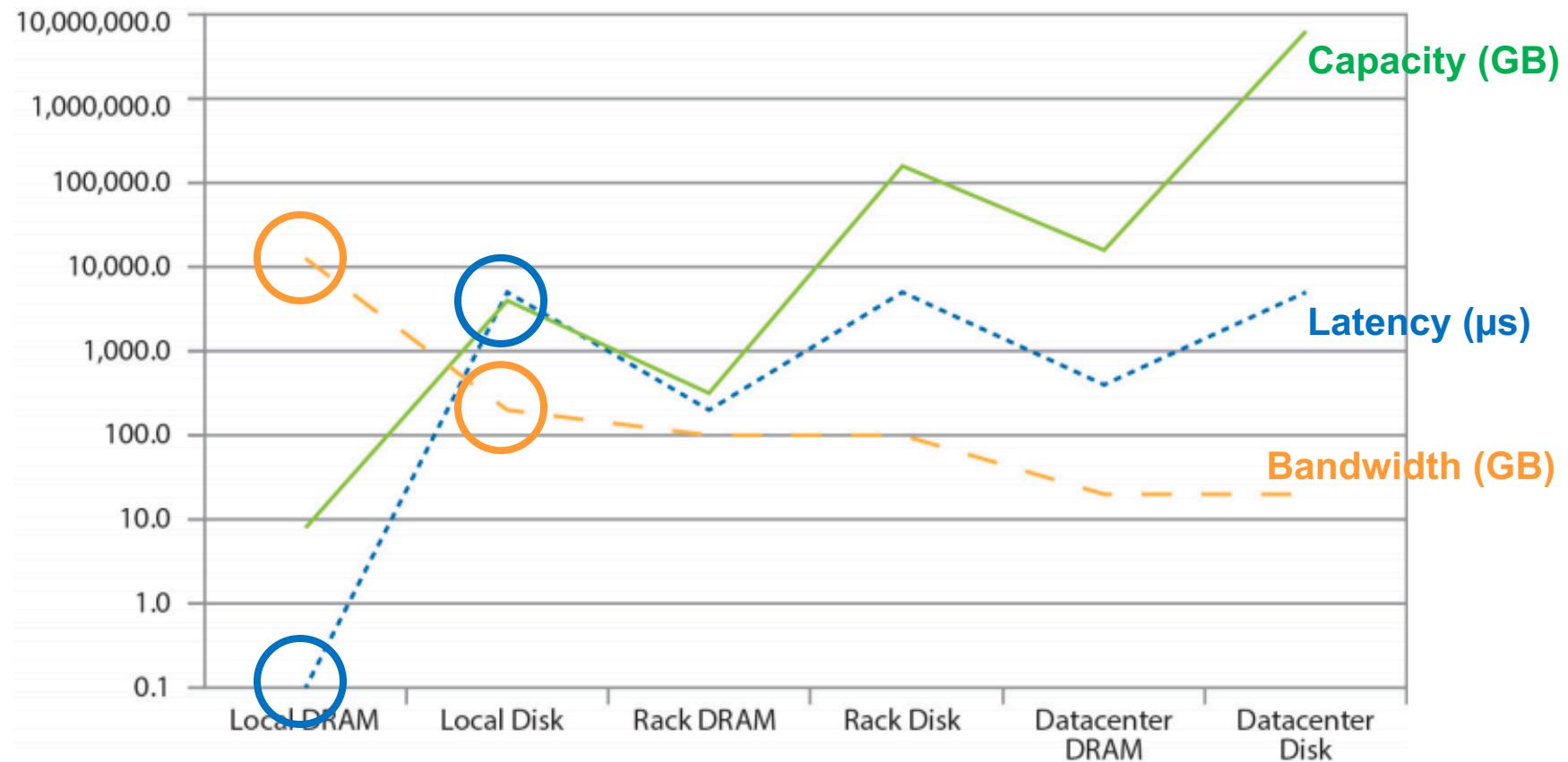
I. **Disk has much higher capacity than DRAM**

Capacity of Storage Hierarchy



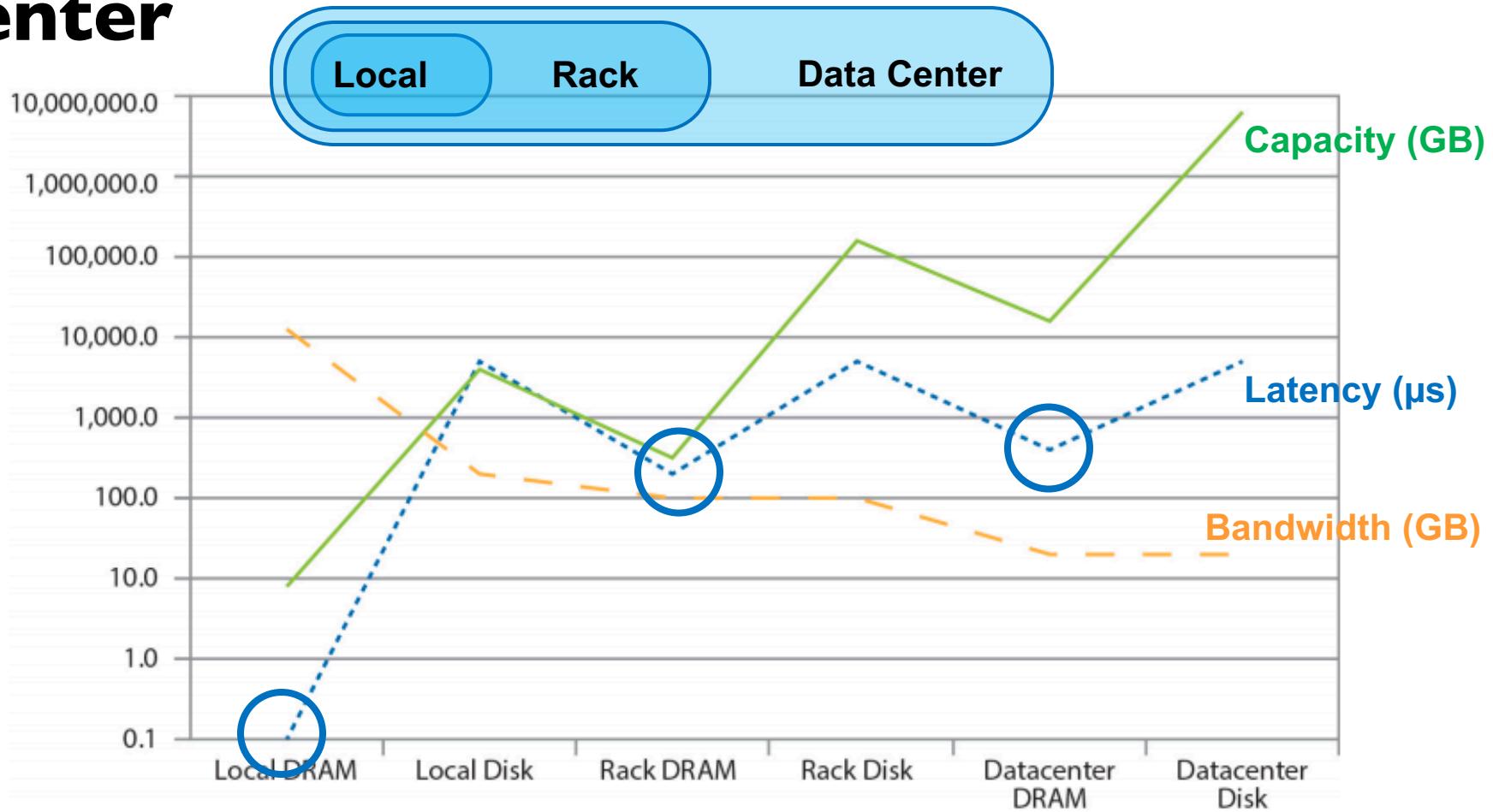
2. **Storage hierarchy:** capacity increases as we go from Local Server, to Rack, to Datacenter.

Speed of Moving Data Around Data Center



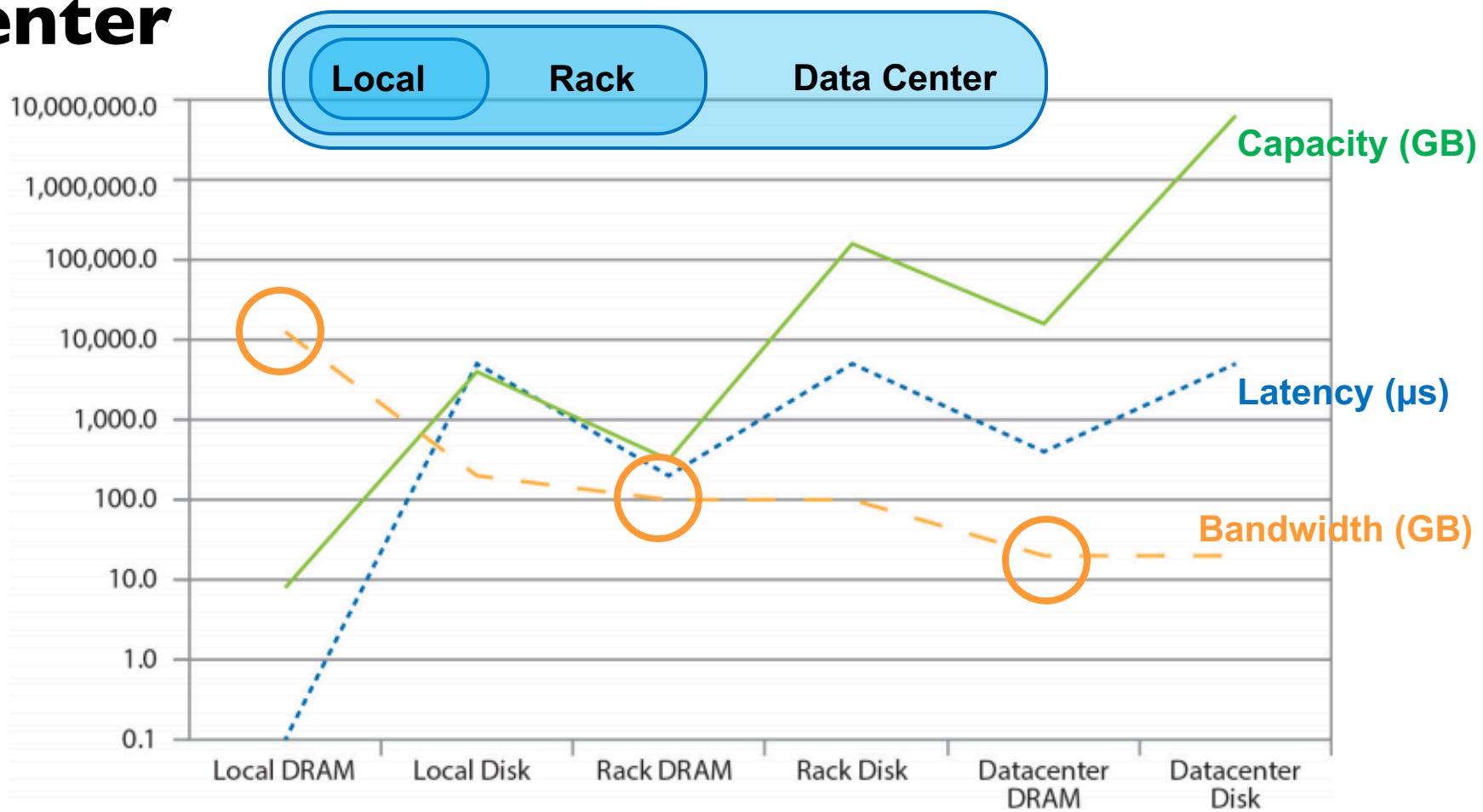
3. Disk reads are much more expensive than DRAM, both in terms of higher latency and lower bandwidth.

Speed of Moving Data Around Data Center



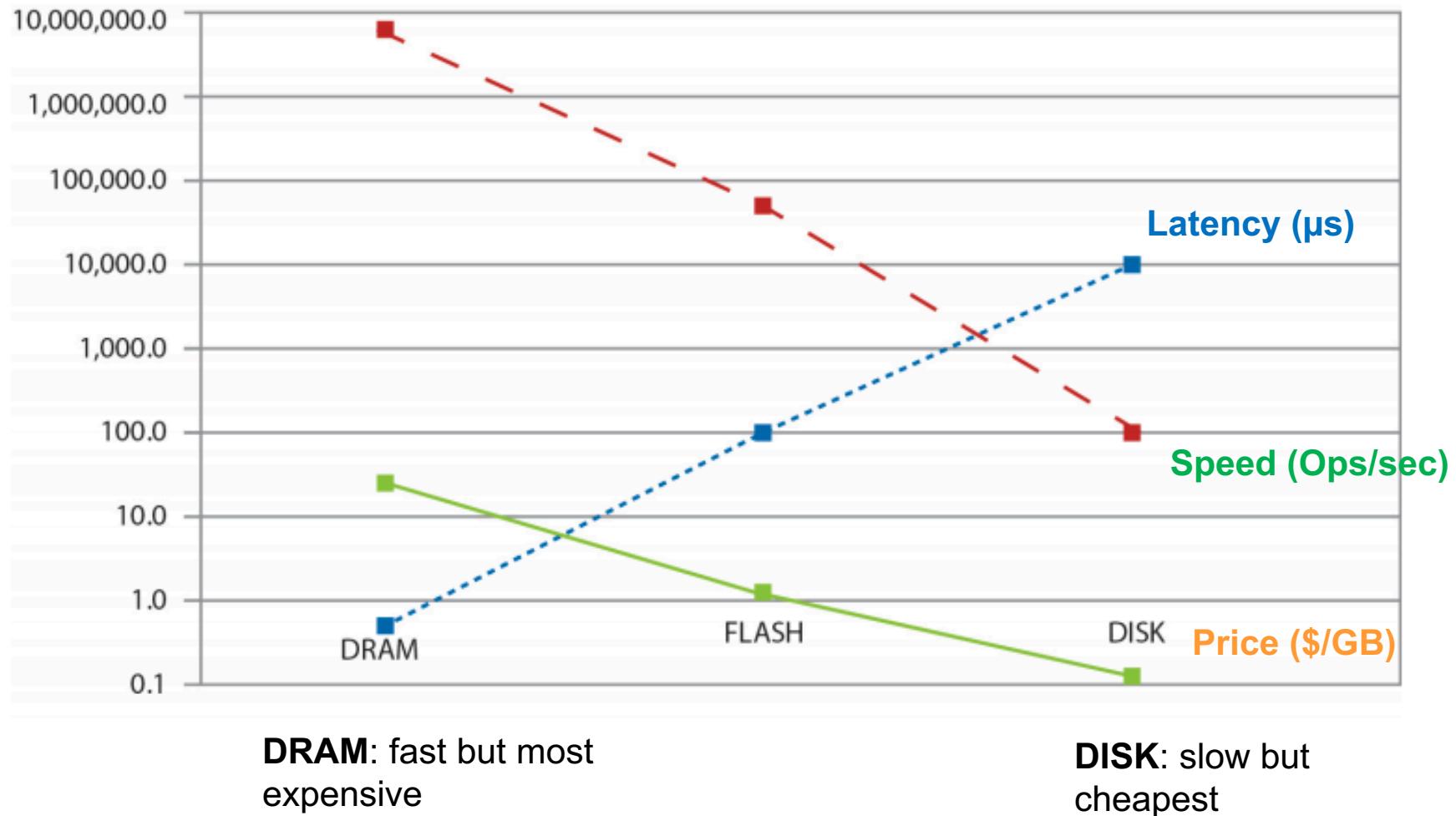
4. Costs increase over the storage hierarchy: **latency increases** as we go from Local to Rack to Datacenter.

Speed of Moving Data Around Data Center



5. Costs increase over the storage hierarchy: bandwidth decreases as we go from Local to Rack to Datacenter.

Price vs Speed Tradeoff



“Big Ideas” of Massive Data Processing in Data Centers

- Scale “out”, not “up”
 - scale ‘out’ = combining many cheaper machines; scale ‘up’: increasing the power of each individual machine
 - Also called ‘horizontal’ vs ‘vertical’ scaling
- Move processing to the data
 - Clusters have limited bandwidth: we should move the task to the machine where the data is stored
- Process data sequentially, avoid random access
 - Seek operations are expensive, disk throughput is reasonable
- Seamless scalability
 - E.g. if processing a certain dataset takes 100 machine hours, ideal scalability is to use a cluster of 10 machines to do it in about 10 hours.



What type of application domains are you most interested in?

Q: What type of application domains are you most interested in?



1. Business
2. Scientific
3. Biomedicine
4. Engineering
5. Finance / Economics
6. Web / E-commerce
7. Education
8. Other

MapReduce

a. Motivation

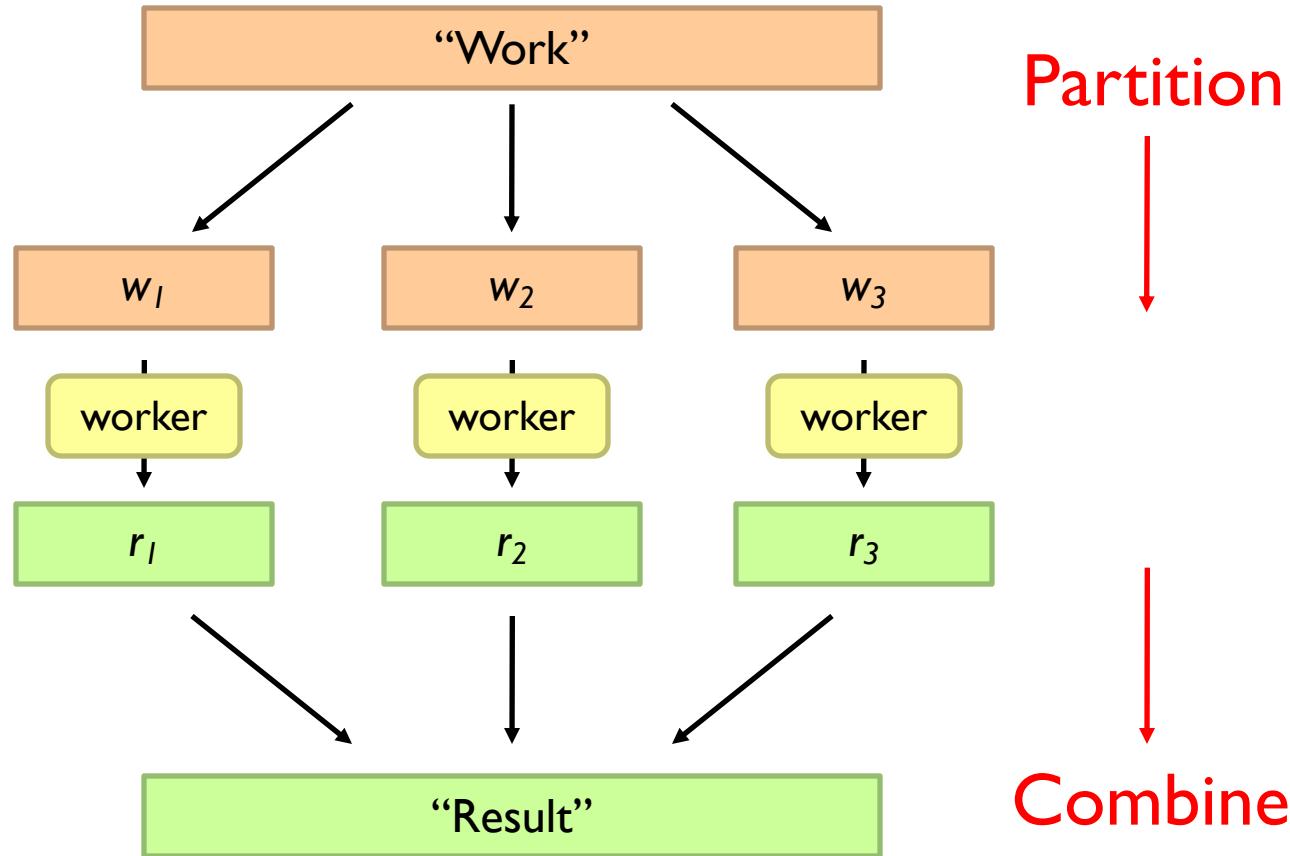
b. Basic MapReduce



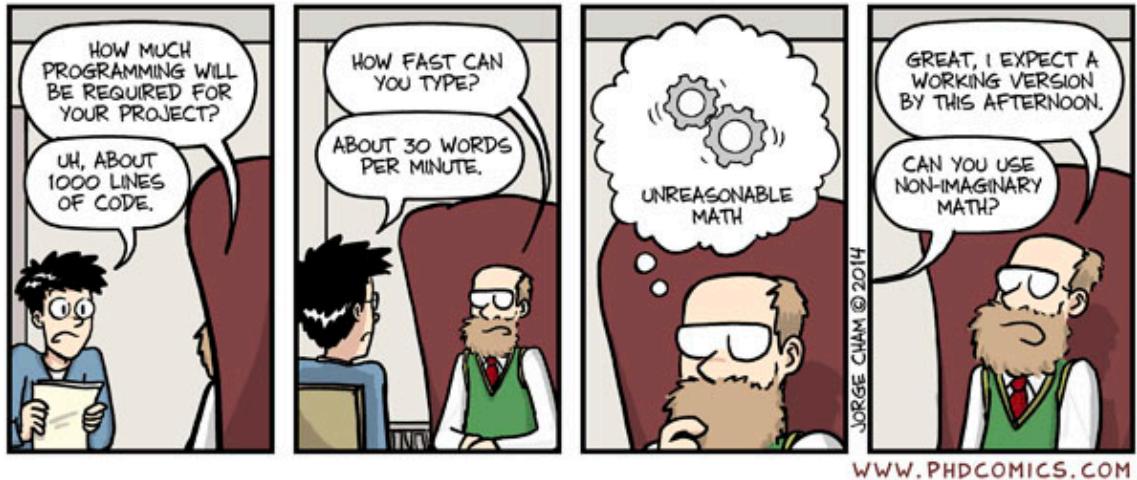
Motivation: Google Example

- 20+ billion web pages: ~400+ TB
- 1 computer reads 30-35 MB/sec from disk
 - ~4 months to read the web
- ~1,000 hard drives to store the web
- Takes even more to **do** something useful with the data!
- Infrastructure: cloud + data centers
 - Cluster of commodity nodes
 - Commodity network (ethernet) to connect them
- Solution:
 - Parallelization + divide and conquer

Divide and Conquer



Challenges



- How do we assign work units to workers?
- What if we have more work units than workers?
- What if workers need to share partial results?
- How do we aggregate partial results?
- How do we know all the workers have finished?
- What if workers die/fail?

Challenge I: Machine Failures

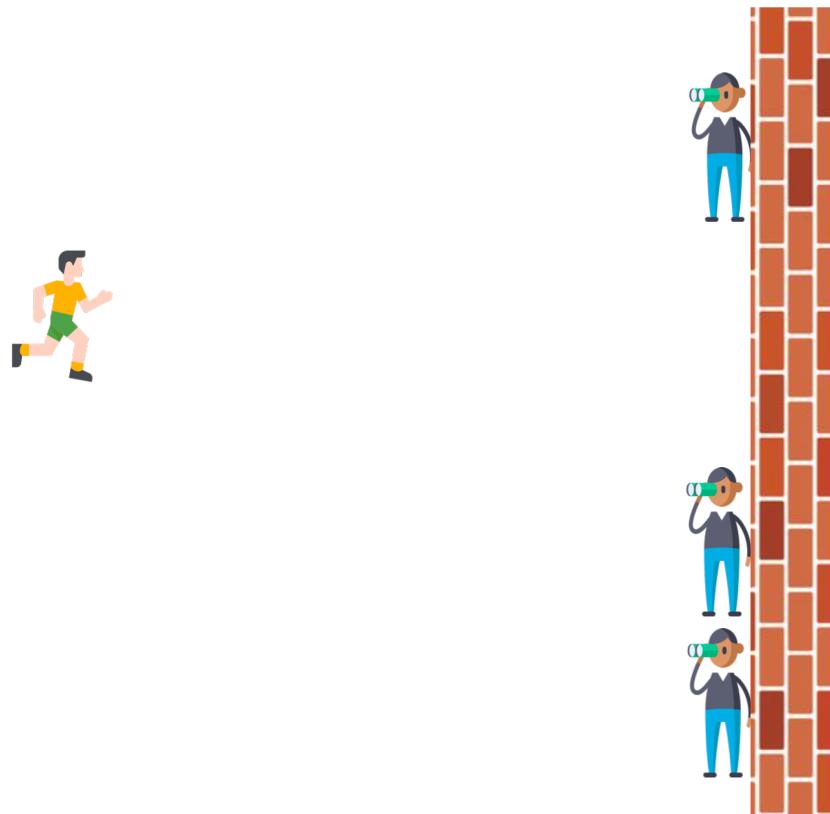
- One server may stay up 10 years (~ 3650 days)
- If you have 3650 servers, expect to loose 1/day
- People estimated Google had ~ 2.5 million machines in 2016
 - ~ 700 machines fail every day!

Challenge 2: Synchronization

- 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 data
- Thus, we need control mechanisms, such as **barriers**

Barrier

- Simple tool used when multiple processes are running at the same time
- It ensures that every process must stop at this point until all processes have reached the barrier



Challenge 3: Programming Difficulty

- Concurrency is difficult to reason about
 - At the scale of datacenters and across datacenters
 - In the presence of failures
 - In terms of multiple interacting services
- Not to mention debugging...
- The reality:
 - Lots of one-off solutions, custom code
 - Write your own dedicated library, then program with it
 - Burden on the programmer to explicitly manage everything

The datacenter *is* the computer

- It's all about the right level of abstraction
 - Moving beyond the single machine architecture
 - What's the “instruction set” (or “API”) 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

- a. Motivation
- b. Basic MapReduce



Typical Big Data Problem

- Iterate over a large number of records
- Extract something of interest from each
- Shuffle and sort intermediate results
- Aggregate intermediate results
- Generate final output

Map

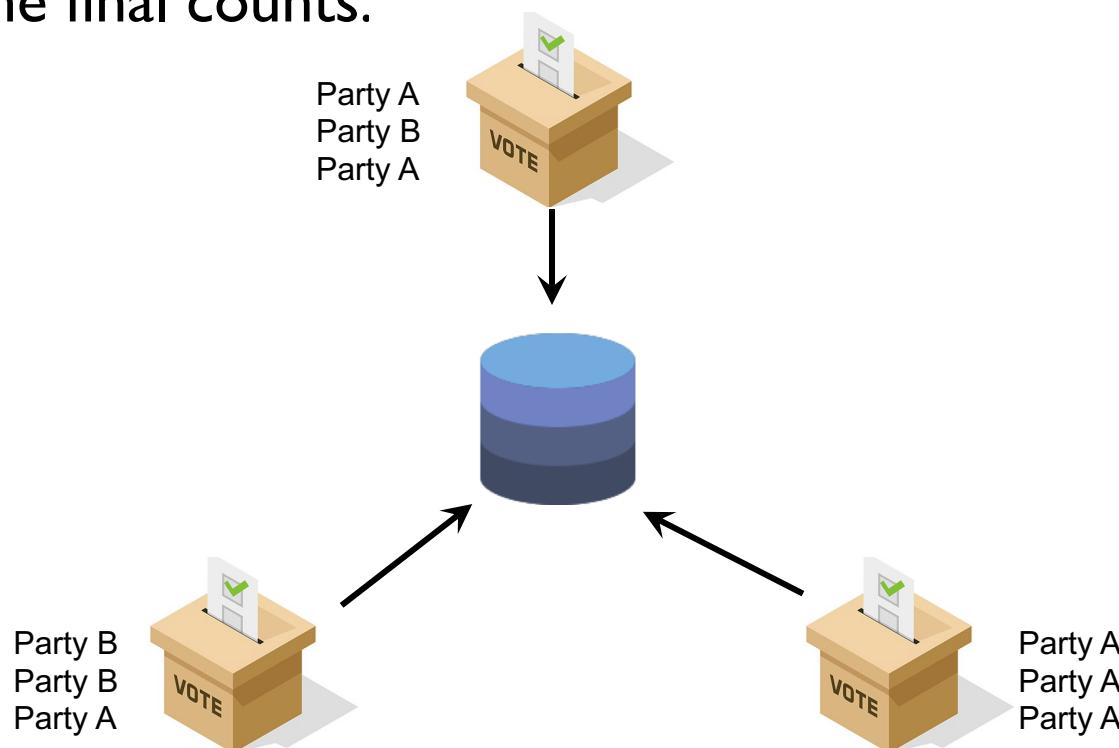
Shuffle

Reduce

Key idea: provide a functional abstraction for these two operations

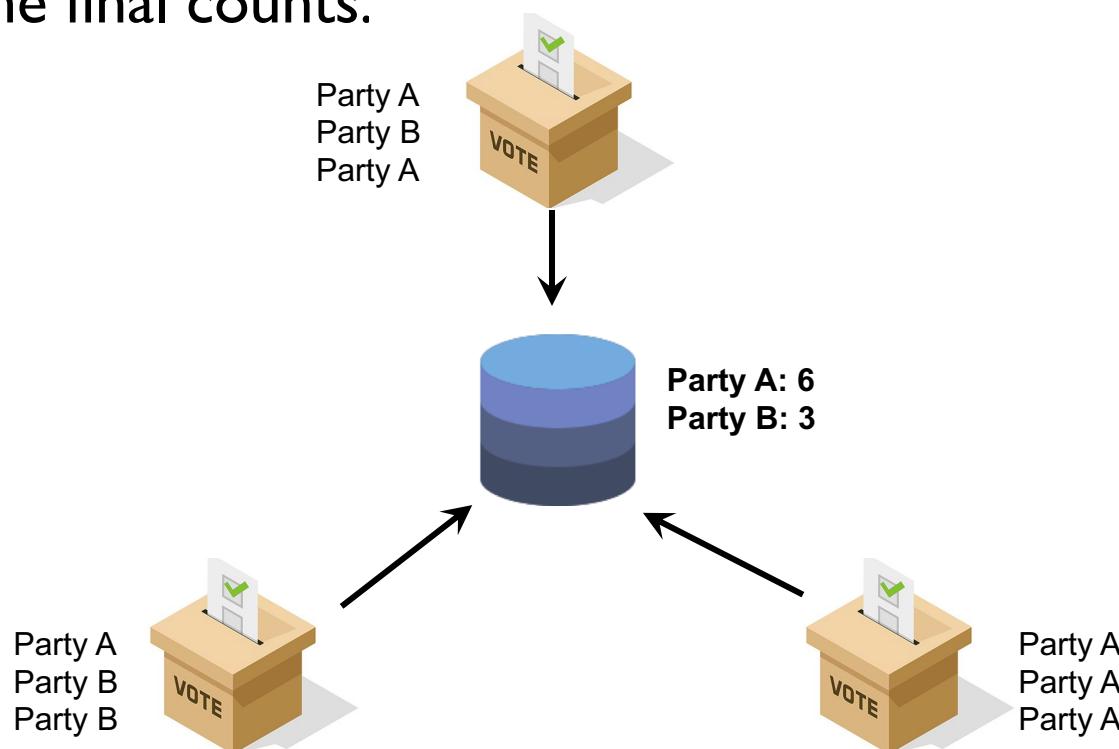
Basic Example: Tabulating Election Results from Multiple Polling Sites

- Imagine you are hired to develop the software for Singapore's election counting system.
- You need to **aggregate** vote counts from multiple stations into the final counts.



Basic Example: Tabulating Election Results from Multiple Polling Sites

- Imagine you are hired to develop the software for Singapore's election counting system.
- You need to **aggregate** vote counts from multiple stations into the final counts.



Tabulating Election Results: MapReduce

Map



Party A
Party B
Party A

Shuffle



Party A
Party B
Party B

Reduce



Party A: 6
Party B: 3



Party A
Party A
Party A

Tabulating Election Results: MapReduce

Map



A: 1
B: 1
A: 1

Shuffle



A: 1
B: 1
B: 1

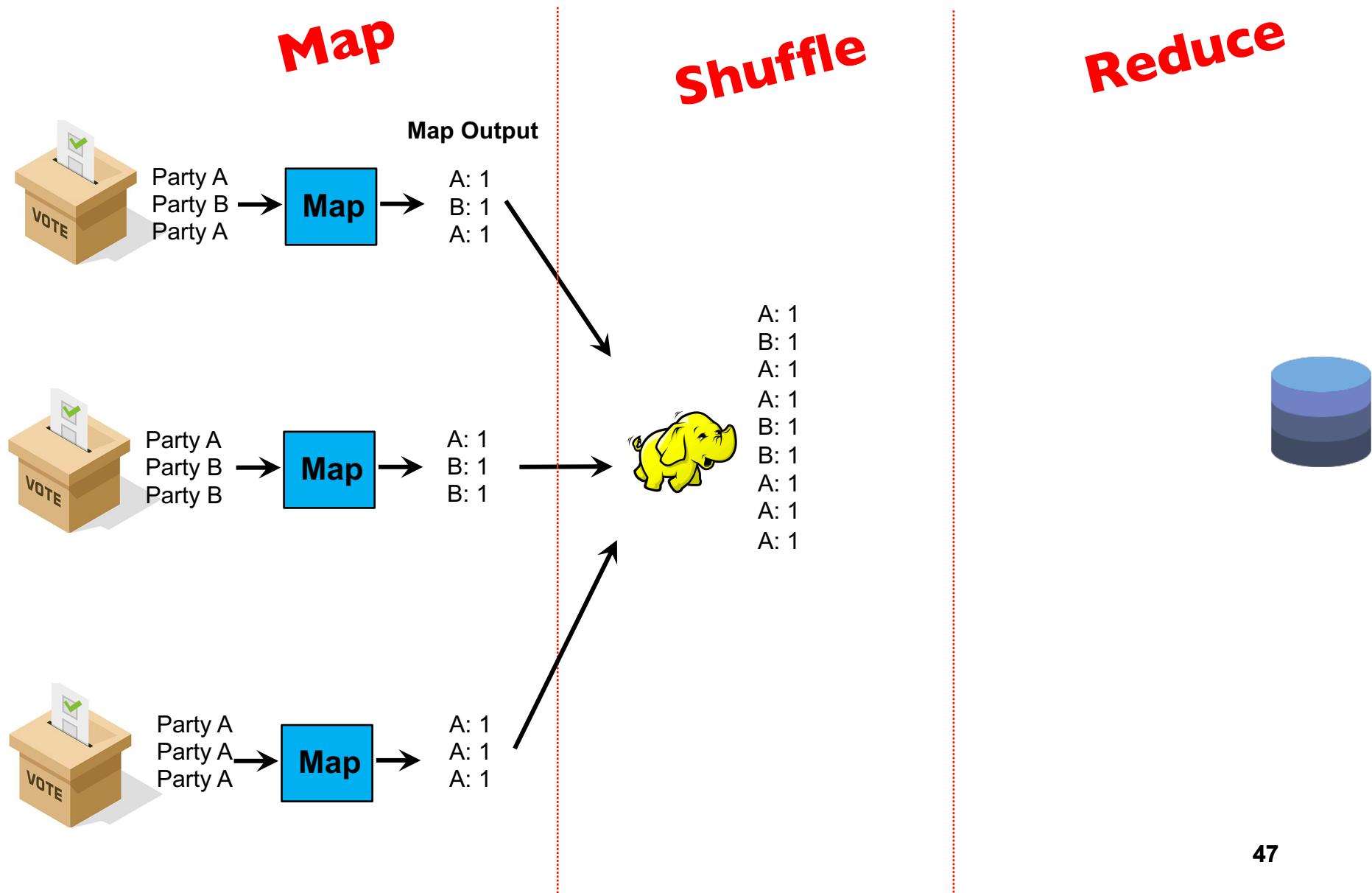
Reduce



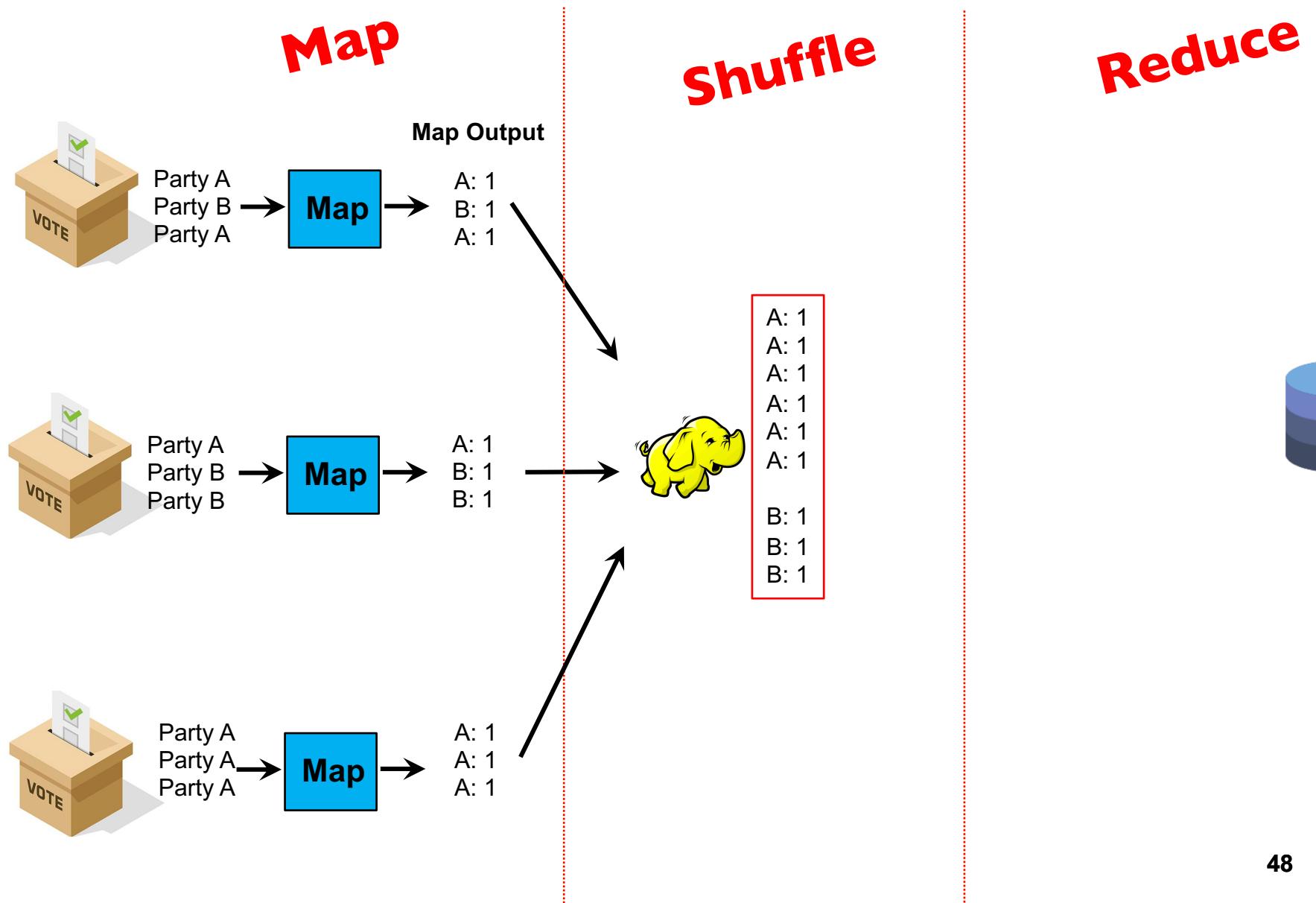
A: 1
A: 1
A: 1



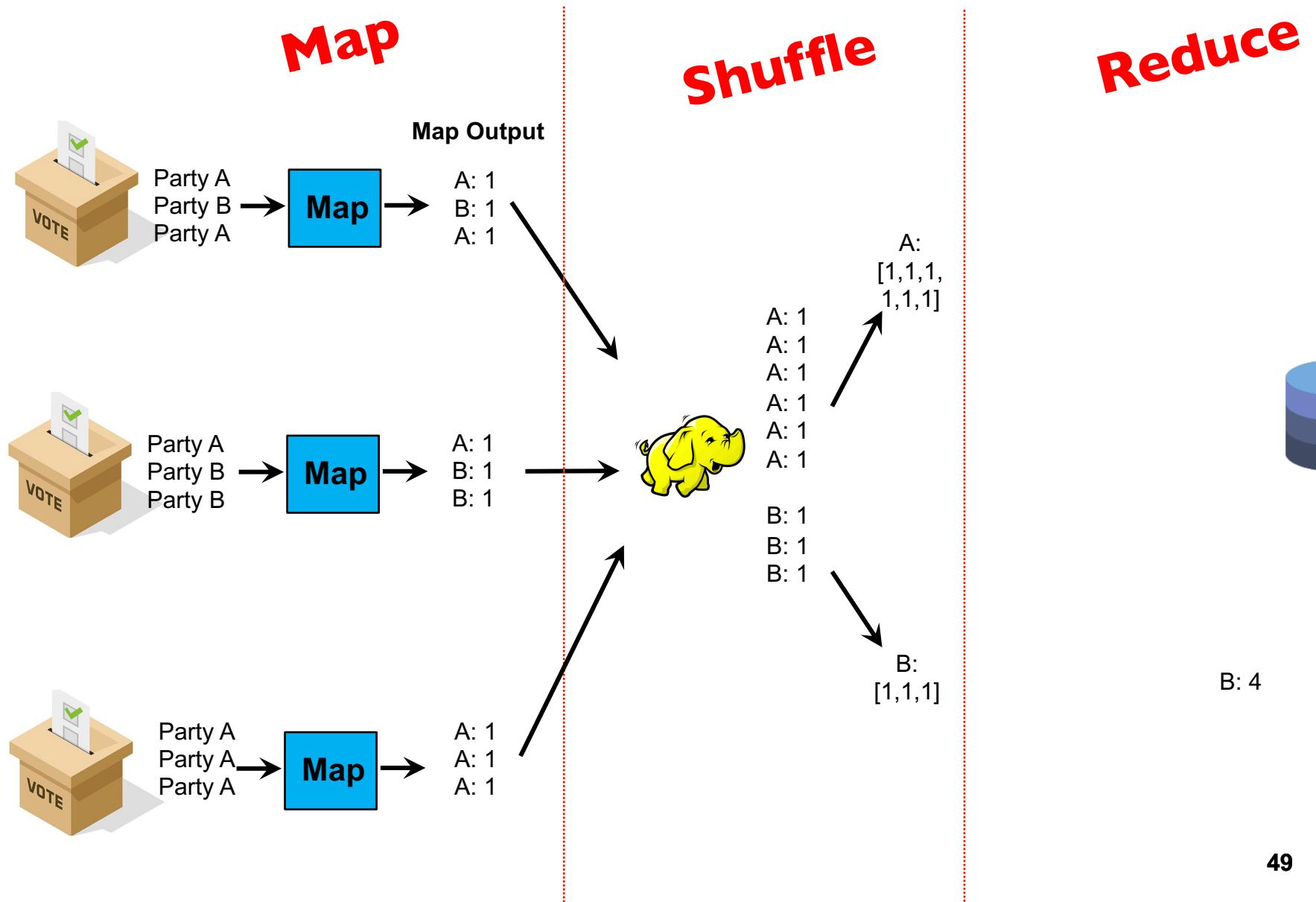
Tabulating Election Results: MapReduce



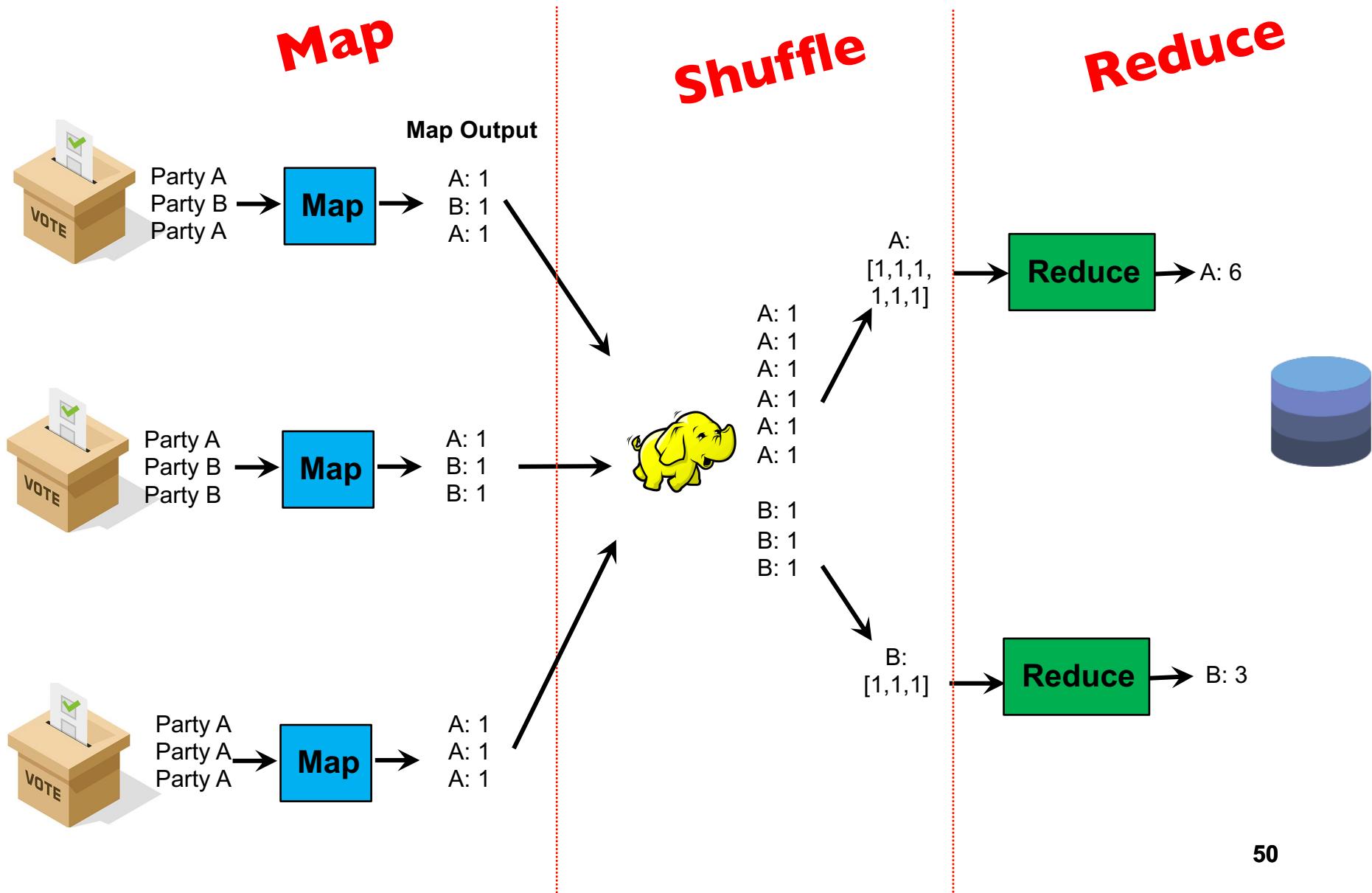
Tabulating Election Results: MapReduce



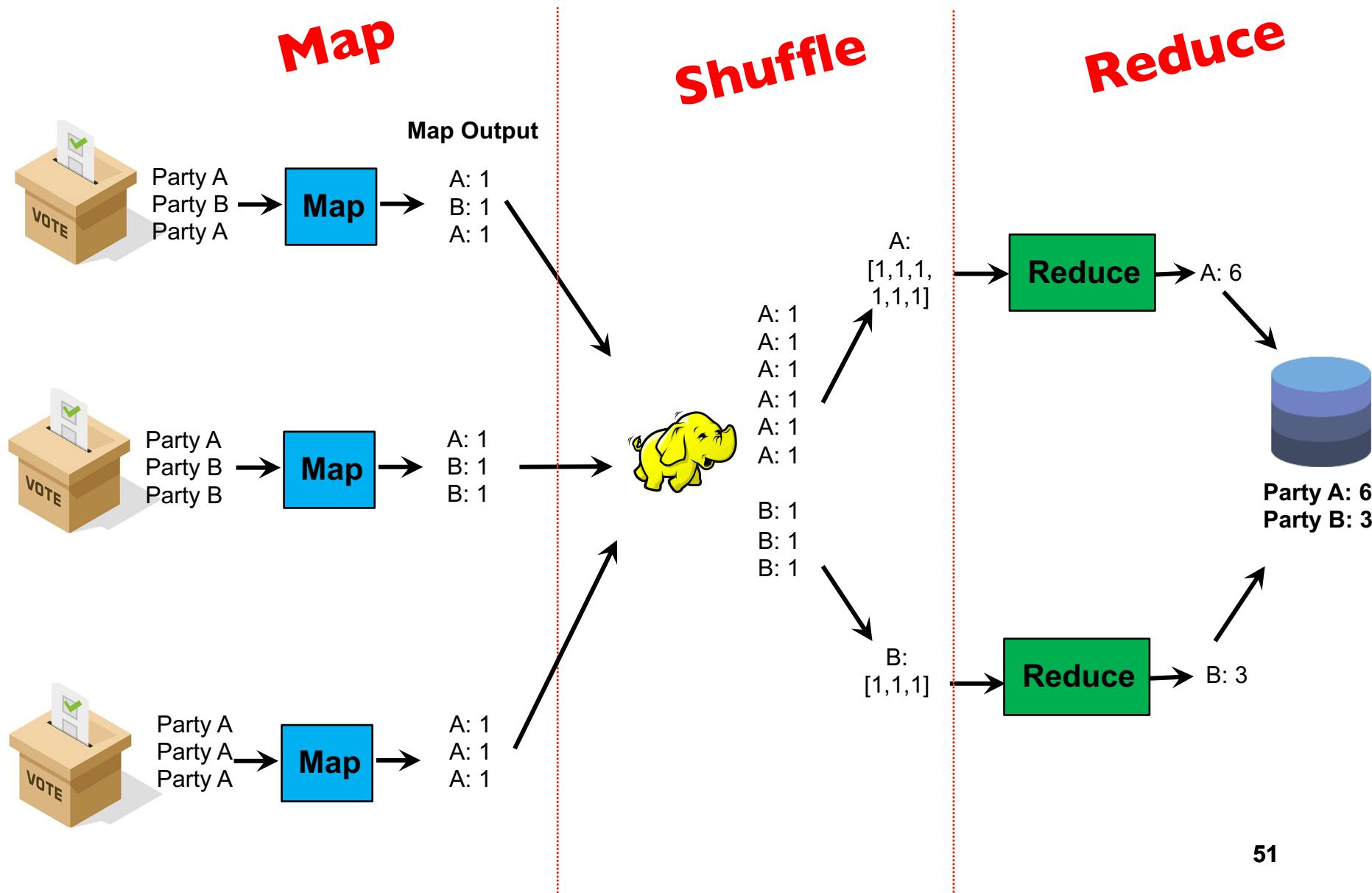
Tabulating Election Results: MapReduce



Tabulating Election Results: MapReduce

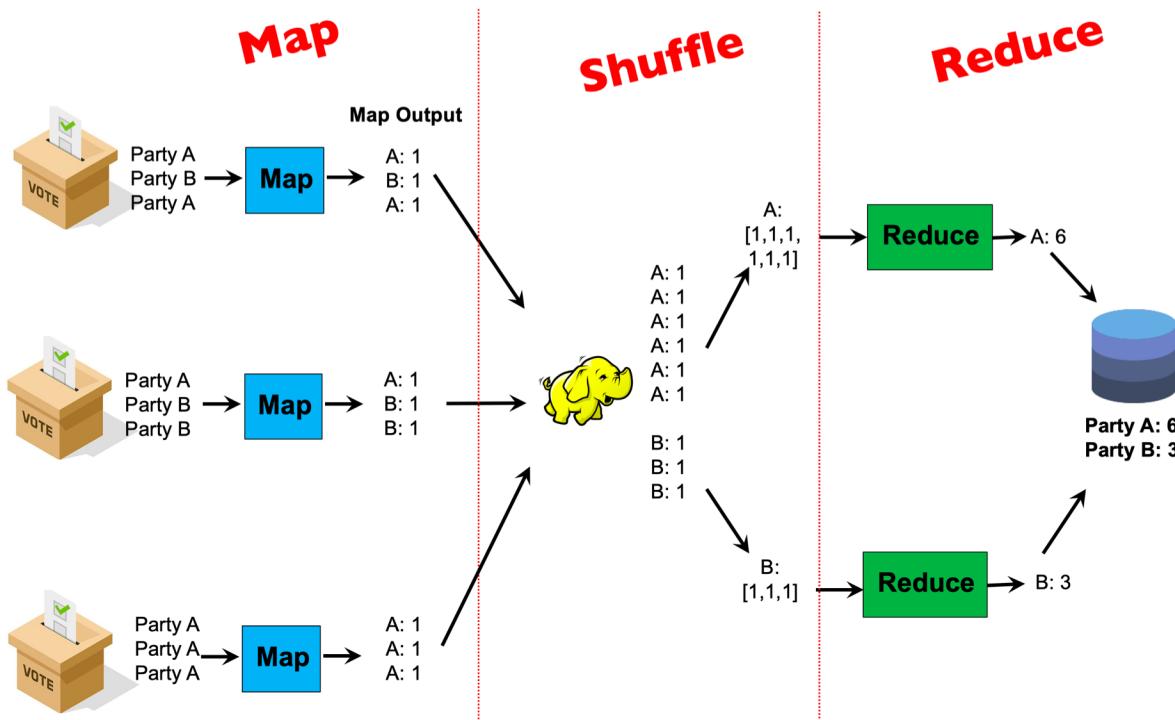


Tabulating Election Results: MapReduce



Writing MapReduce Programs

- Typical Interface: Programmers specify two functions:
 - map** (k_1, v_1) \rightarrow List(k_2, v_2)
 - reduce** ($k_2, \text{List}(v_2)$) \rightarrow List(k_3, v_3)
- All values with the same key are sent to the same reducer



Writing MapReduce Programs

- **Typical Interface:** Programmers specify two functions:
 - map** (k_1, v_1) \rightarrow List(k_2, v_2)
 - reduce** (k_2 , List(v_2)) \rightarrow List(k_3, v_3)
- All values with the same key are sent to the same reducer
- The execution framework handles many challenging issues...
 - How do we assign work units to workers?
 - What if we have more work units than workers?
 - What if workers need to share partial results?
 - How do we aggregate partial results?
 - How do we know all the workers have finished?
 - What if workers die/fail?

MapReduce Execution Framework

- Handles scheduling
 - Assigns workers to map and reduce tasks
- Handles “data distribution”
 - Moves processes to data
- Handles synchronization
 - Gathers, sorts, and shuffles intermediate data
- Handles errors and faults
 - Detects worker failures and restarts
- Everything happens on top of a distributed file system (later)