# Big Data Infrastructure

CS 489/698 Big Data Infrastructure (Winter 2016)

## Week 1: Introduction (2/2)

January 7, 2016

Jimmy Lin

David R. Cheriton School of Computer Science

University of Waterloo

These slides are available at http://lintool.github.io/bigdata-2016w/

# Why big data?

Science
Engineering
Commerce
Society
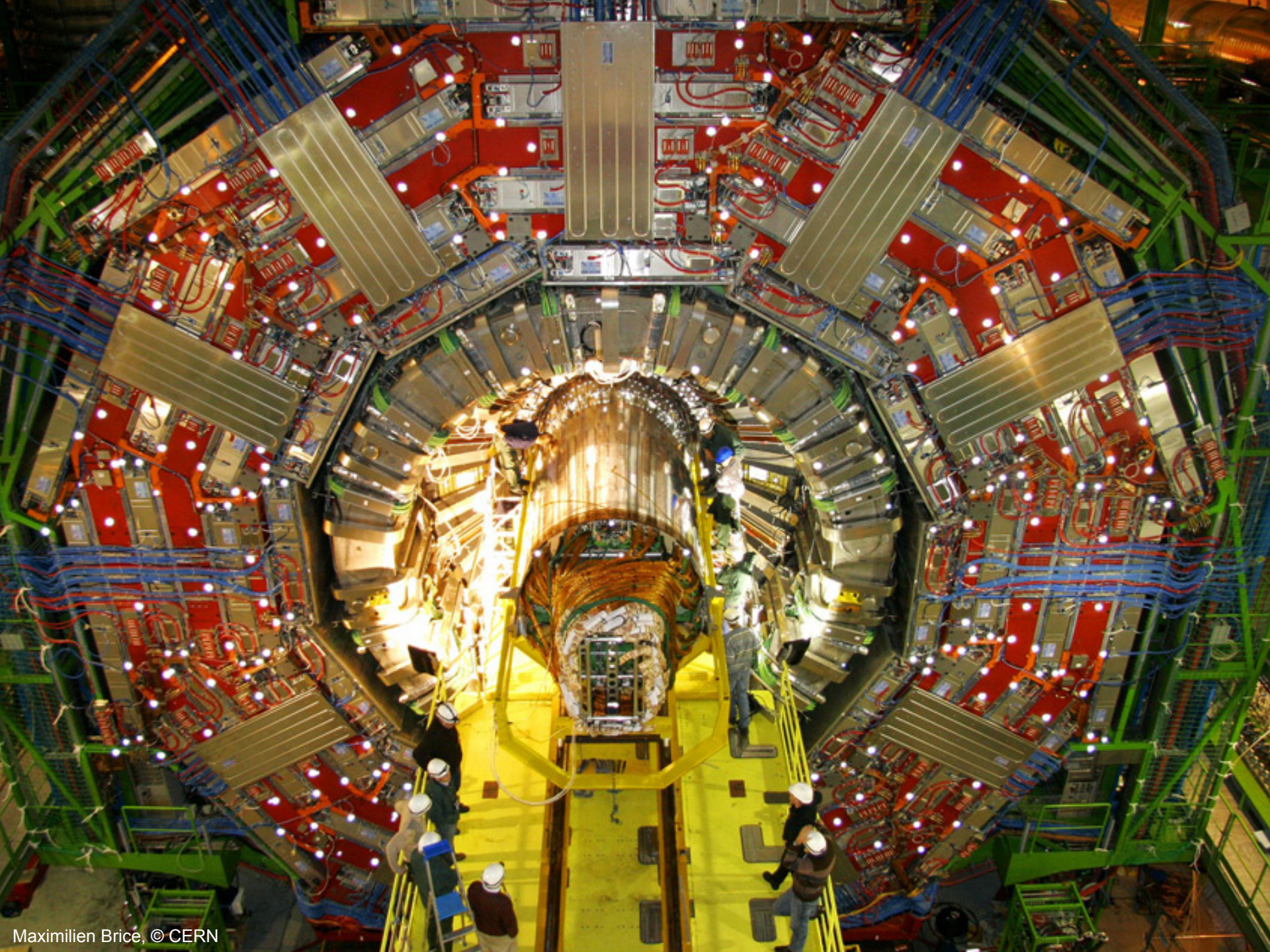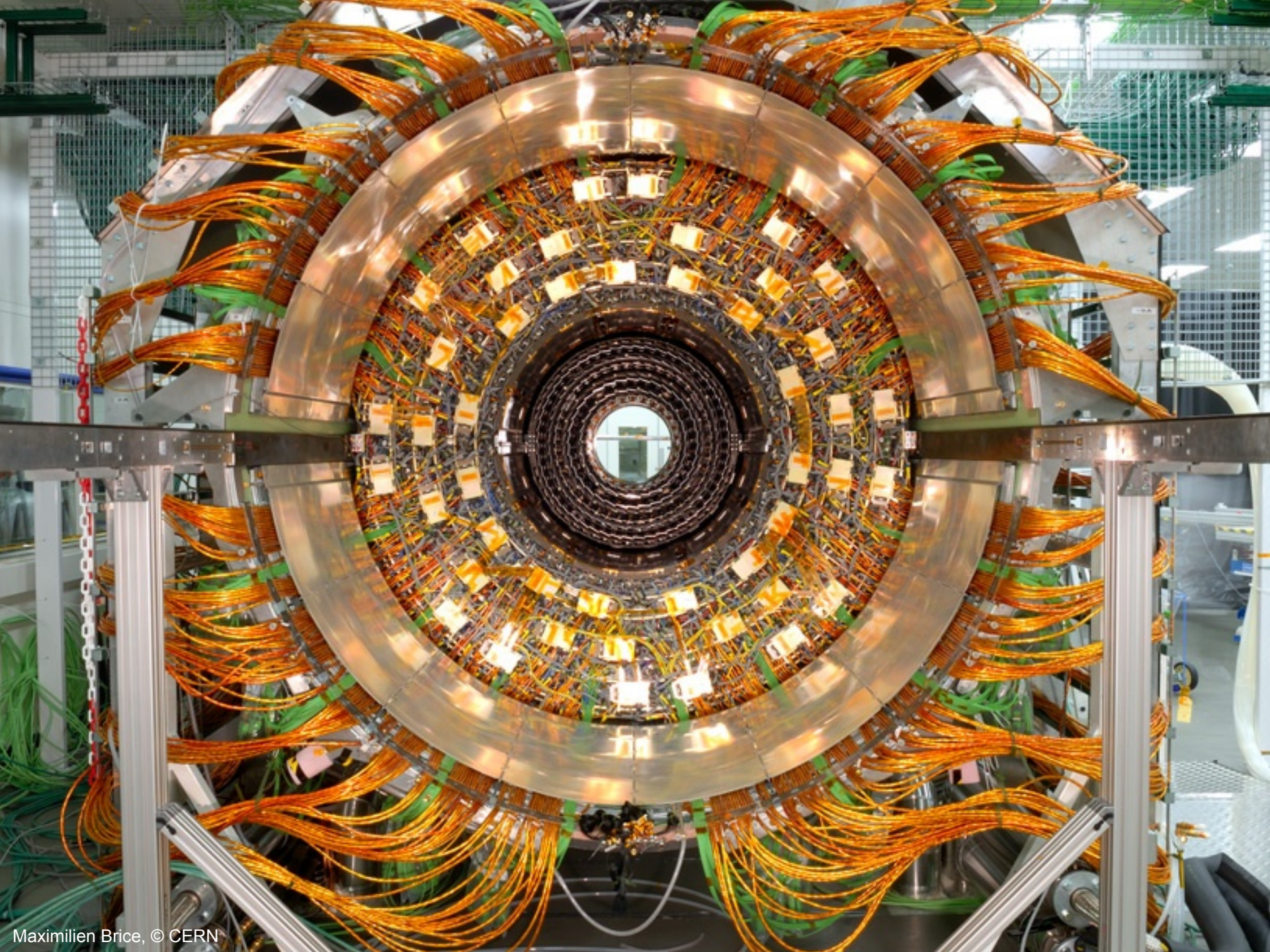
Science

Emergence of the 4th Paradigm

Data-intensive e-Science

Maximilien Brice, © CERN

**Subject genome**

**Sequencer**

**Reads**

GATGCTTACTATGCGGGCCCC
CGGTCTAATGCTTACTATGC
GCTTACTATGCGGGCCCCTT
AATGCTTACTATGCGGGCCCCTT
TAATGCTTACTATGC
AATGCTTAGCTATGCGGGC
AATGCTTACTATGCGGGCCCCTT
AATGCTTACTATGCGGGCCCCTT
CGGTCTAGATGCTTACTATGC
AATGCTTACTATGCGGGCCCCTT
CGGTCTAATGCTTAGCTATGC
ATGCTTACTATGCGGGCCCCTT

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

# Engineering

The unreasonable effectiveness of data

Search, recommendation, prediction, …

# Translate

How does Google's translation system work?   ✕

如何谷歌的翻译系统的工作?

Wrong?

Rúhé gǔgē de fānyì xìtǒng de gōngzuò?

# No data like more data!

s/knowledge/data/g;



Test data BLEU vs LM training data size in million tokens:
- +0.62BP/x2
- +0.66BP/x2
- +0.70BP/x2
- +0.56BP/x2
- +0.39BP/x2
- +0.51BP/x2
- +0.15BP/x2

Legend:
- target KN
- +ldcnews KN
- +webnews KN
- target SB
- +ldcnews SB
- +webnews SB
- +web SB

(Banko and Brill, ACL 2001)
(Brants et al., EMNLP 2007)

Know thy customers

Data → Insights → Competitive advantages

# Commerce

# 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!?

# 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**

# Society

Humans as social sensors

Computational social science

# Predicting *X* with Twitter



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

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



# Political Mobilization on Facebook

(Paul and Dredze, ICWSM 2011; Bond et al., Nature 2011)

# Tackling Big Data

$k_1$ $v_1$  $k_2$ $v_2$  $k_3$ $v_3$  $k_4$ $v_4$  $k_5$ $v_5$  $k_6$ $v_6$

map    map    map    map

a 1  b 2    c 3  c 6    a 5  c 2    b 7  c 8

combine    combine    combine    combine

a 1  b 2    c 9    a 5  c 2    b 7  c 8

partition    partition    partition    partition

**Shuffle and Sort:** aggregate values by keys

a  1 5    b  2 7    c  2 9 6 8

reduce    reduce    reduce

$r_1$ $s_1$    $r_2$ $s_2$    $r_3$ $s_3$

User Program

(1) submit

Master

(2) schedule map          (2) schedule reduce

worker

split 0
split 1
(3) read
split 2
split 3
split 4

worker

(4) local write

worker

(5) remote read

worker          (6) write          output file 0

worker          output file 1

**Input files**          **Map phase**          **Intermediate files (on local disk)**          **Reduce phase**          **Output files**

Adapted from (Dean and Ghemawat, OSDI 2004)

# 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

The datacenter *is* the computer!

# Building Blocks



cluster
switch

server
racks

# Storage Hierarchy



**One Server**
DRAM: 16 GB, 100 ns, 20 GB/s
Disk: 2T B, 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

# Storage Hierarchy

# Storage Hierarchy



Source: Barroso and Urs Hölzle (2013)

# Anatomy of a Datacenter



**Computer Air Handling Unit (CRAC)**
- Up To 30 Ton Sensible Capacity Per Unit
- Air Discharge Can Be Upflow Or Downflow Configuration
- Downflow Configuration Used With Raised Floor To Create A Pressurized Supply Air Plenum With Floor Supply Diffusers

**Power Distribution Unit (PDU)**
- Typical Capacities Up To 225 kVA Per Unit
- Redundancy Through Dual PDU's With Integral Static Transfer Switch (STS)

**Individual Colocation Computer Cabinets**
- Typ. Cabinet Footprint (28"W x 36"D x 84"H)
- Typical Capacities Of 1750 To 3750 Watts Per Cabinet

**Emergency Diesel Generators**
- Total Generator Capacity = Total Electrical Load To Building
- Multiple Generators Can Be Electrically Combined With Paralleling Gear
- Can Be Located Indoors Or Outdoors At Grade Or On Roof.
- Outdoor Applications Require Sound Attenuating Enclosures

**Fuel Oil Storage Tanks**
- Tank Capacity Dependant On Length Of Generator Operation
- Can Be Located Underground Or At Grade Or Indoors

**Colocation Suites**
- Modular Configuration For Flexible Suite Sq.Ft. Areas.
- Suites Consist Of Multiple Cabinets With Secured Partitions (Cages, Walls, Etc.)

**UPS System**
- Uninterruptible Power Supply Modules
- Up To 1000 kVA Per Module
- Cabinets And Battery Strings Or Rotary Flywheels
- Multiple Redundancy Configurations Can Be Designed

**Electrical Primary Switchgear**
- Includes Incoming Service And Distribution
- Direct Distribution To Mechanical Equipment
- Distribution To Secondary Electrical Equipment Via UPS

**Heat Rejection Devices**
- Drycoolers, Air Cooled Chillers, Etc.
- Up To 400 Ton Capacity Per Unit
- Mounted At Grade Or On Roof
- N+1 Design

**Pump Room**
- Used To Pump Condenser/Chilled Water Between Drycoolers And CRAC Units
- Additional Equipment Includes Expansion Tank, Glycol Feed System
- N+1 Design (Standby Pump)

Source: Barroso and Urs Hölzle (2013)

# Anatomy of a Datacenter

Aside: How much is 30 MW?

Source: Google

# 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

# "Big Ideas"

- Scale "out", not "up"
  - Limits of SMP and large shared-memory machines
- Move processing to the data
  - Cluster have limited bandwidth
- Process data sequentially, avoid random access
  - Seeks are expensive, disk throughput is reasonable
- Seamless scalability
  - From the mythical man-month to the tradable machine-hour

# Scaling "up" vs. "out"

- No single machine is large enough
  - Smaller cluster of large SMP machines vs. larger cluster of commodity machines (e.g., 16 128-core machines vs. 128 16-core machines)

- Nodes need to talk to each other!
  - Intra-node latencies: ~100 ns
  - Inter-node latencies: ~100 $\mu$s

- Let's model communication overhead…

# Modeling Communication Costs

○ Simple execution cost model:

- Total cost = cost of computation + cost to access global data
- Fraction of local access inversely proportional to size of cluster
- $n$ nodes (ignore cores for now)

$$1 \text{ ms} + f \times [100 \text{ ns} \times (1/n) + 100 \text{ } \mu\text{s} \times (1 - 1/n)]$$

- Light communication: $f = 1$
- Medium communication: $f = 10$
- Heavy communication: $f = 100$

○ What are the costs in parallelization?

# Cost of Parallelization

# Advantages of scaling "up"



So why not?
Why does commodity beat exotic?

# Counterpoint: Scaling up?

- No single machine is large enough

  - Smaller cluster of large SMP machines vs. larger cluster of commodity machines (e.g., 16 128-core machines vs. 128 16-core machines)

- Is this really true? Modern "commodity" machine:

  - Four 18-core processors: 72 cores total
  - 3TB RAM

Who really has big data problems?

# Numbers Everyone Should Know

According to Jeff Dean

| | |
|---|---:|
| L1 cache reference | 0.5 ns |
| Branch mispredict | 5 ns |
| L2 cache reference | 7 ns |
| Mutex lock/unlock | 100 ns |
| Main memory reference | 100 ns |
| Compress 1K bytes with Zippy | 10,000 ns |
| Send 2K bytes over 1 Gbps network | 20,000 ns |
| Read 1 MB sequentially from memory | 250,000 ns |
| Round trip within same datacenter | 500,000 ns |
| Disk seek | 10,000,000 ns |
| Read 1 MB sequentially from network | 10,000,000 ns |
| Read 1 MB sequentially from disk | 30,000,000 ns |
| Send packet CA->Netherlands->CA | 150,000,000 ns |

Google

# Moving Data Around
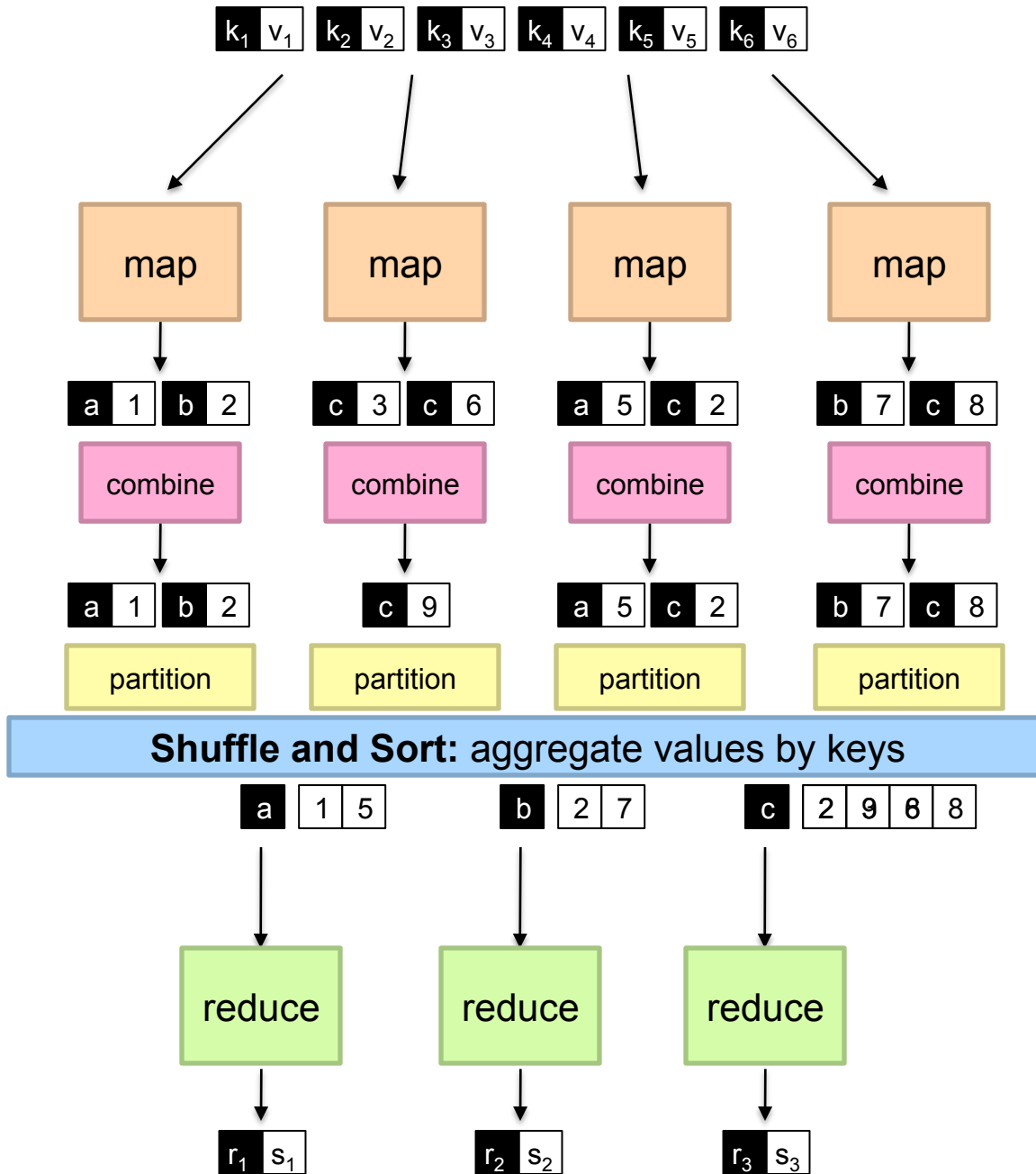
# Seeks vs. Scans

- Consider a 1 TB database with 100 byte records
  - We want to update 1 percent of the records
- Scenario 1: random access
  - Each update takes ~30 ms (seek, read, write)
  - $10^8$ updates = ~35 days
- Scenario 2: rewrite all records
  - Assume 100 MB/s throughput
  - Time = 5.6 hours(!)
- Lesson: avoid random seeks!

# Justifying the "Big Ideas"

- Scale "out", not "up"

  - Limits of SMP and large shared-memory machines

- Move processing to the data

  - Cluster have limited bandwidth

- Process data sequentially, avoid random access

  - Seeks are expensive, disk throughput is reasonable

- Seamless scalability

  - From the mythical man-month to the tradable machine-hour

MapReduce data flow diagram

Input pairs: $k_1\ v_1$, $k_2\ v_2$, $k_3\ v_3$, $k_4\ v_4$, $k_5\ v_5$, $k_6\ v_6$

**map** → a 1 | b 2
**map** → c 3 | c 6
**map** → a 5 | c 2
**map** → b 7 | c 8

**combine** → a 1 | b 2
**combine** → c 9
**combine** → a 5 | c 2
**combine** → b 7 | c 8

**partition** (×4)

**Shuffle and Sort:** aggregate values by keys

a | 1 | 5
b | 2 | 7
c | 2 | 9 | 6 | 8

**reduce** → $r_1\ s_1$
**reduce** → $r_2\ s_2$
**reduce** → $r_3\ s_3$

# How do we get data to the workers?



**NAS**

**SAN**

**Compute Nodes**

**What's the problem here?**

# Distributed File System

- Don't move data to workers… move workers to the data!
  - Store data on the local disks of nodes in the cluster
  - Start up the workers on the node that has the data local
- Why?
  - (Perhaps) not enough RAM to hold all the data in memory
  - Disk access is slow, but disk throughput is reasonable
- A distributed file system is the answer
  - GFS (Google File System) for Google's MapReduce
  - HDFS (Hadoop Distributed File System) for Hadoop

# GFS: Assumptions

- Commodity hardware over "exotic" hardware

  - Scale "out", not "up"

- High component failure rates

  - Inexpensive commodity components fail all the time

- "Modest" number of huge files

  - Multi-gigabyte files are common, if not encouraged

- Files are write-once, mostly appended to

  - Perhaps concurrently

- Large streaming reads over random access

  - High sustained throughput over low latency

# GFS: Design Decisions

- Files stored as chunks

  - Fixed size (64MB)

- Reliability through replication

  - Each chunk replicated across 3+ chunkservers

- Single master to coordinate access, keep metadata

  - Simple centralized management

- No data caching

  - Little benefit due to large datasets, streaming reads

- Simplify the API

  - Push some of the issues onto the client (e.g., data layout)

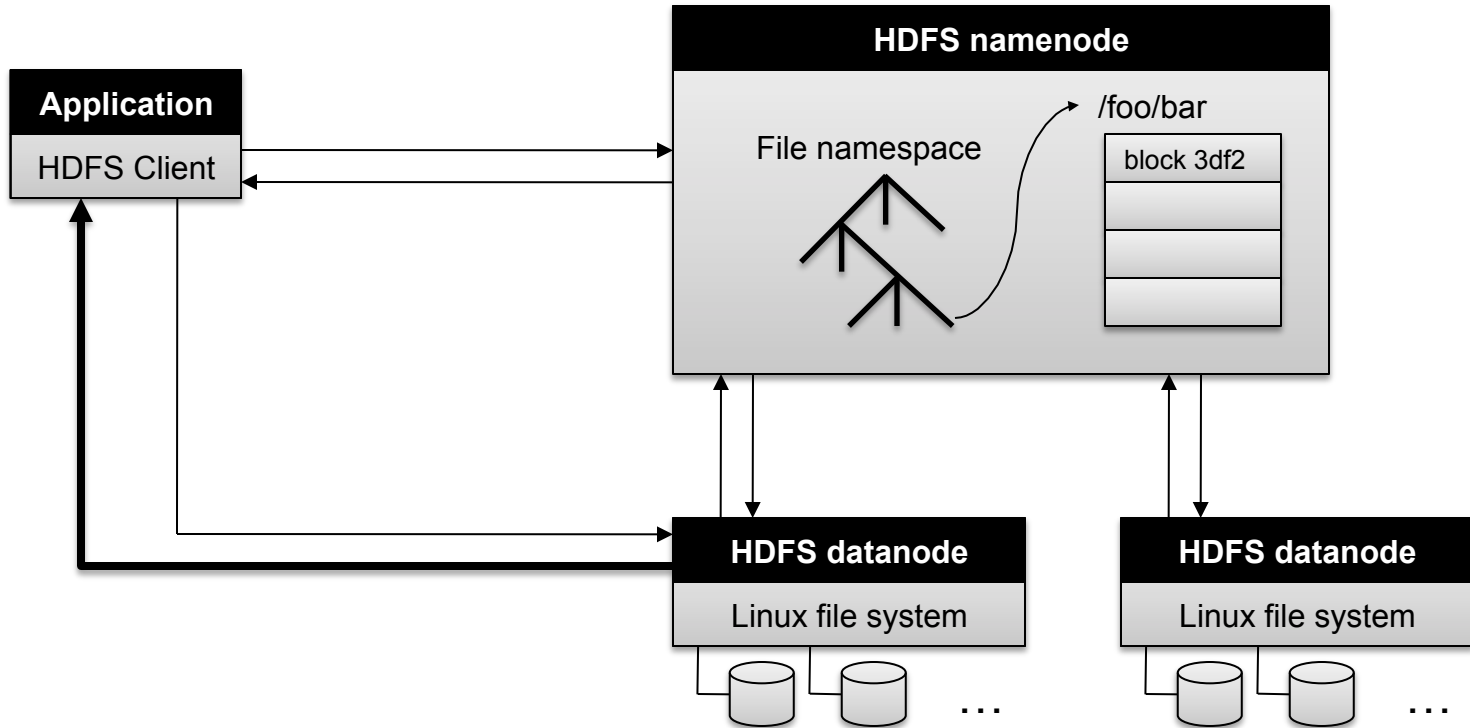**HDFS = GFS clone (same basic ideas)**

# From GFS to HDFS

- Terminology differences:

  - GFS master = Hadoop namenode

  - GFS chunkservers = Hadoop datanodes

- Differences:

  - Different consistency model for file appends
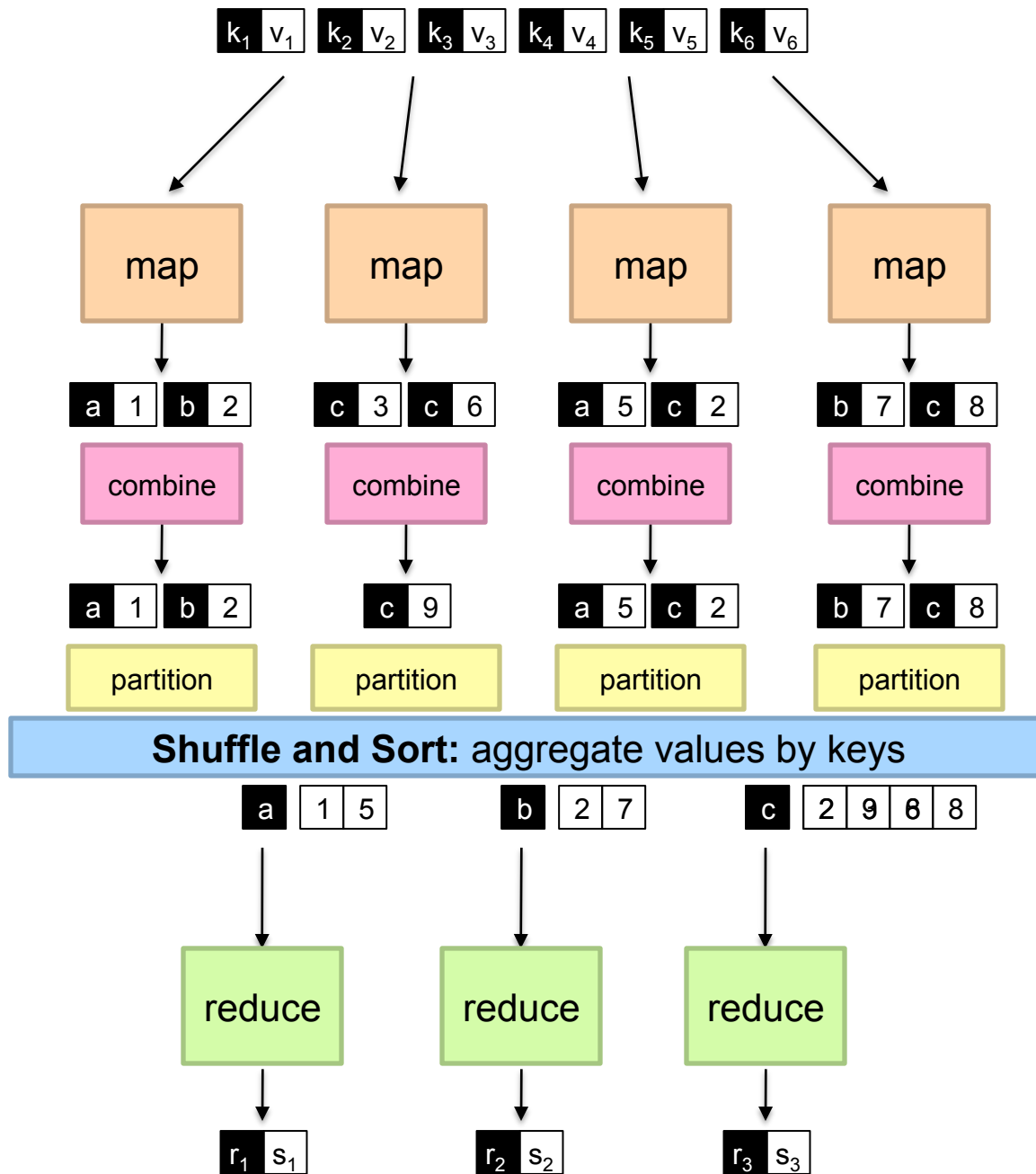
  - Implementation

  - Performance

**For the most part, we'll use Hadoop terminology...**

# HDFS Architecture

**HDFS namenode**

File namespace

/foo/bar

block 3df2

**Application**

HDFS Client

**HDFS datanode**

Linux file system

…

**HDFS datanode**

Linux file system

…

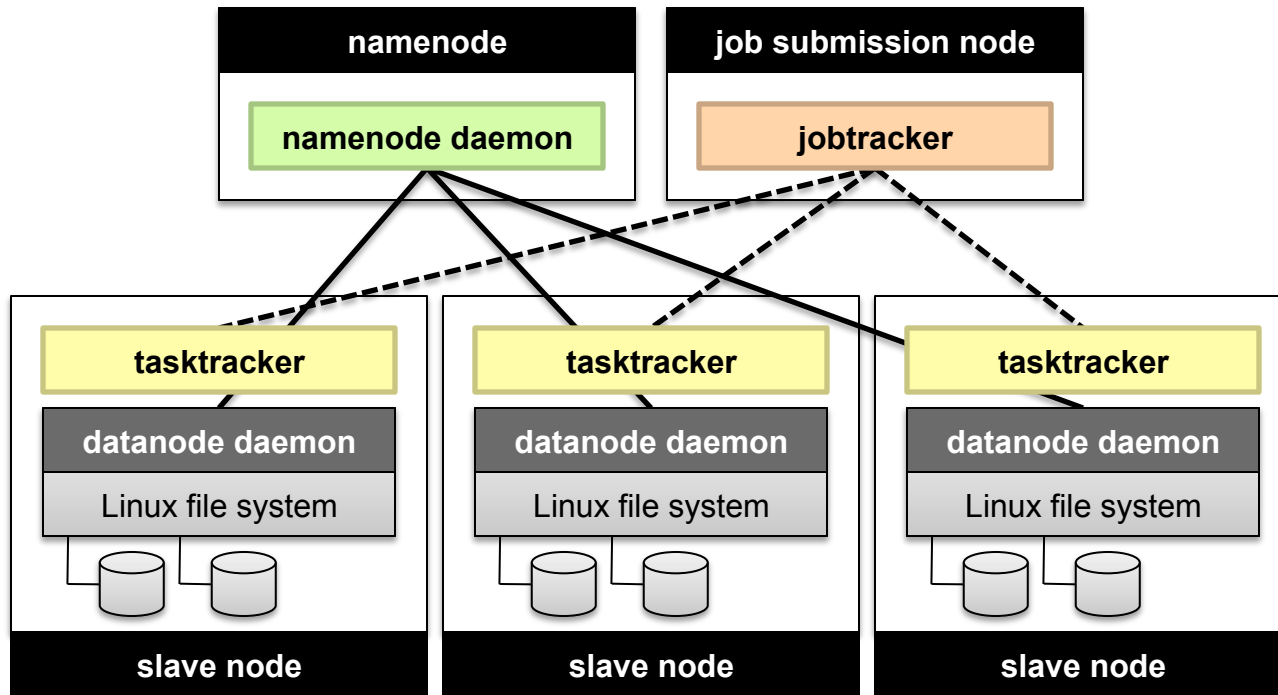Adapted from (Ghemawat et al., SOSP 2003)

# Namenode Responsibilities

- Managing the file system namespace:

  - Holds file/directory structure, metadata, file-to-block mapping, access permissions, etc.

- Coordinating file operations:

  - Directs clients to datanodes for reads and writes
  - No data is moved through the namenode

- Maintaining overall health:

  - Periodic communication with the datanodes
  - Block re-replication and rebalancing
  - Garbage collection

| $k_1$ | $v_1$ | $k_2$ | $v_2$ | $k_3$ | $v_3$ | $k_4$ | $v_4$ | $k_5$ | $v_5$ | $k_6$ | $v_6$ |

map    map    map    map

| a | 1 | b | 2 |    | c | 3 | c | 6 |    | a | 5 | c | 2 |    | b | 7 | c | 8 |

combine    combine    combine    combine

| a | 1 | b | 2 |    | c | 9 |    | a | 5 | c | 2 |    | b | 7 | c | 8 |

partition    partition    partition    partition

**Shuffle and Sort:** aggregate values by keys

| a | 1 | 5 |    | b | 2 | 7 |    | c | 2 | 9 | 6 | 8 |

reduce    reduce    reduce

| $r_1$ | $s_1$ |    | $r_2$ | $s_2$ |    | $r_3$ | $s_3$ |

# Putting everything together…

# Sequoia
16.32 PFLOPS
98,304 nodes with 1,572,864 million cores
1.6 petabytes of memory
7.9 MWatts total power

Source: LLNL

**Aside: Cloud Computing**

# The best thing since sliced bread?

- Before clouds…
  - Grids
  - Connection machine
  - Vector supercomputers
  - …
- Cloud computing means many different things:
  - Big data
  - Rebranding of web 2.0
  - Utility computing
  - Everything as a service

# Rebranding of web 2.0

- Rich, interactive web applications
  - Clouds refer to the servers that run them
  - AJAX as the de facto standard (for better or worse)
  - Examples: Facebook, YouTube, Gmail, …

- "The network is the computer": take two
  - User data is stored "in the clouds"
  - Rise of the netbook, smartphones, etc.
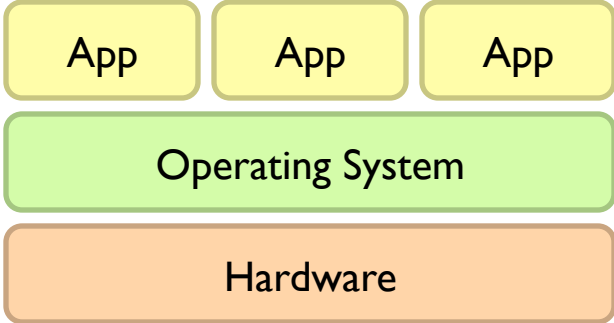  - Browser *is* the OS

# Utility Computing

○ What?

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

○ Why?

- Cost: capital vs. operating expenses
- Scalability: "infinite" capacity
- Elasticity: scale up or down on demand

○ Does it make sense?

- Benefits to cloud users
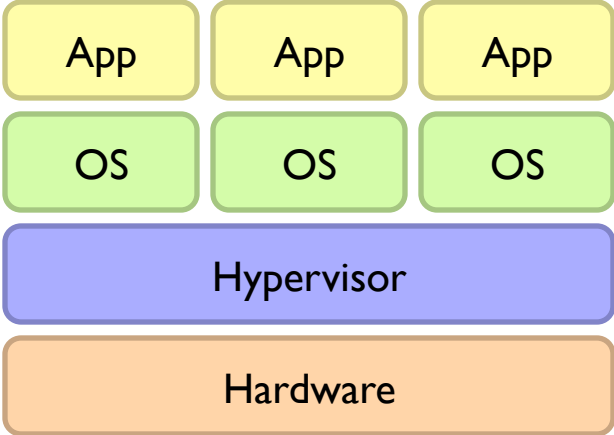- Business case for cloud providers

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

# Enabling Technology: Virtualization



Traditional Stack

Virtualized Stack

# Everything as a Service

- Utility computing = Infrastructure as a Service (IaaS)

  - Why buy machines when you can rent cycles?
  - Examples: Amazon's EC2, Rackspace

- Platform as a Service (PaaS)

  - Give me nice API and take care of the maintenance, upgrades, …
  - Example: Google App Engine

- Software as a Service (SaaS)

  - Just run it for me!
  - Example: Gmail, Salesforce

# Who cares?

○ A source of problems…

- Cloud-based services *generate* big data
- Clouds make it easier to start companies that *generate* big data

○ As well as a solution…

- Ability to provision analytics clusters on-demand in the cloud
- Commoditization and democratization of big data capabilities

# Questions?

Remember: Assignment 0 due next Tuesday at 8:30am