

# Lecture 33: Introduction to MapReduce & Cloud Computing

Dan Andresen  
[dan@k-state.edu](mailto:dan@k-state.edu)



Based on "Introduction to MapReduce" by Jimmy Lin. 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

# What is MapReduce?

- Programming model for expressing distributed computations at a massive scale
- Execution framework for organizing and performing such computations
- Open-source implementation called Hadoop



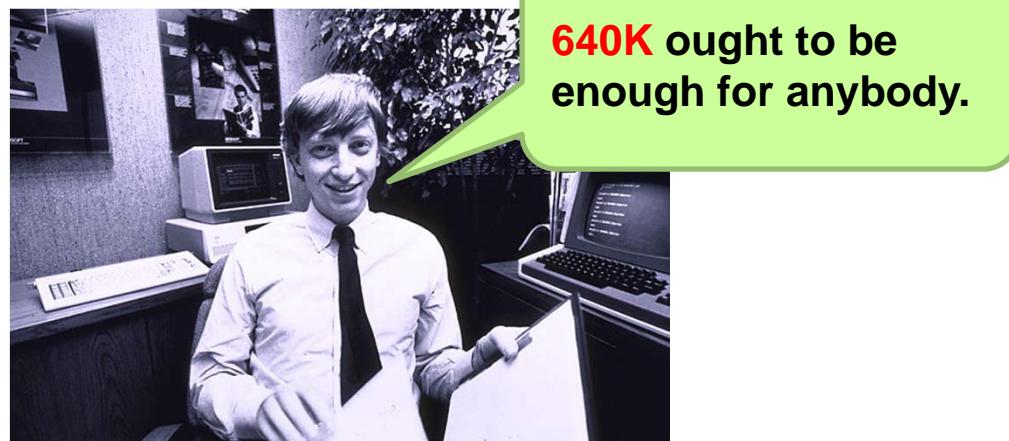
# Why large data?

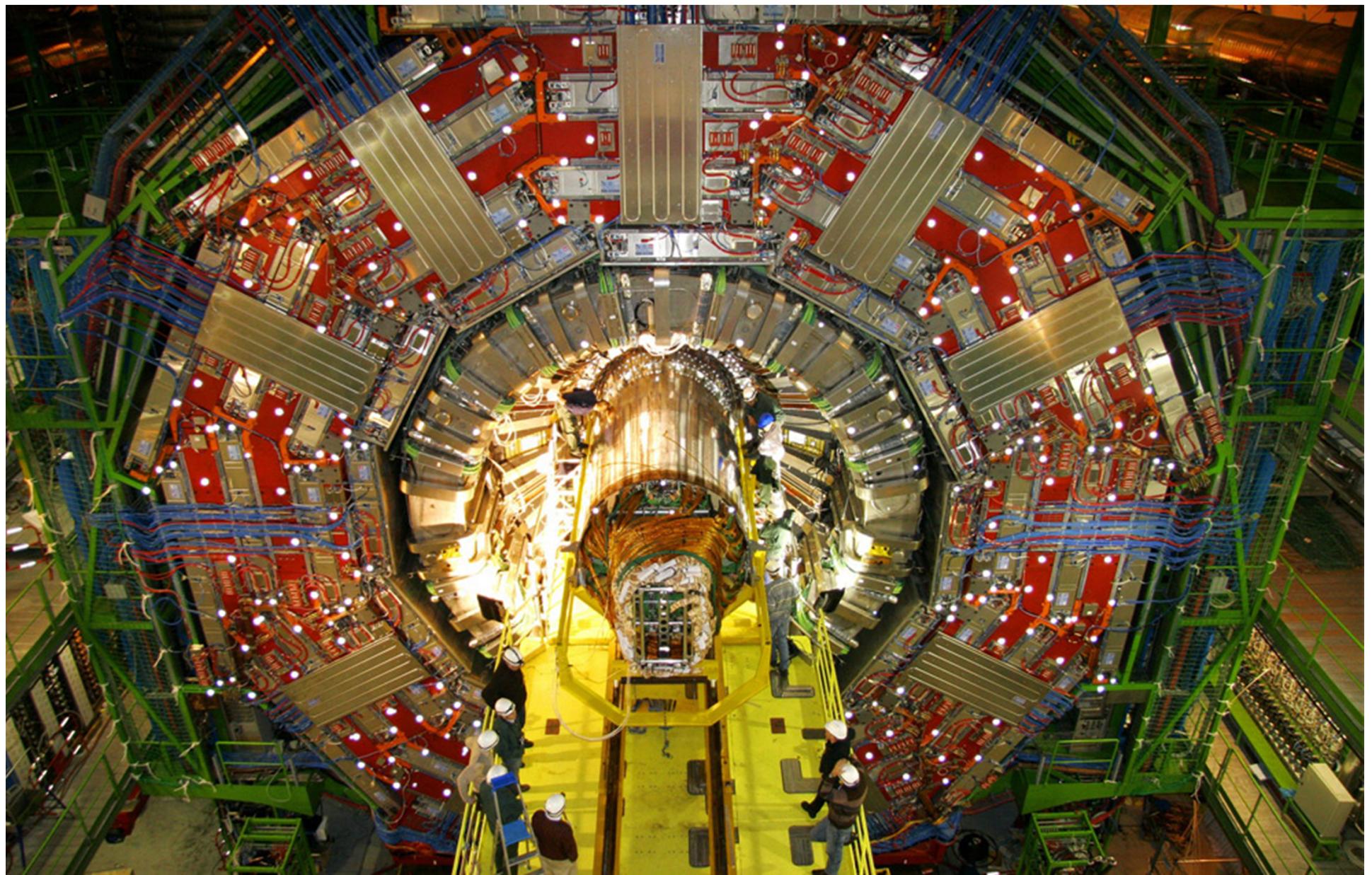


Source: Wikipedia (Everest)

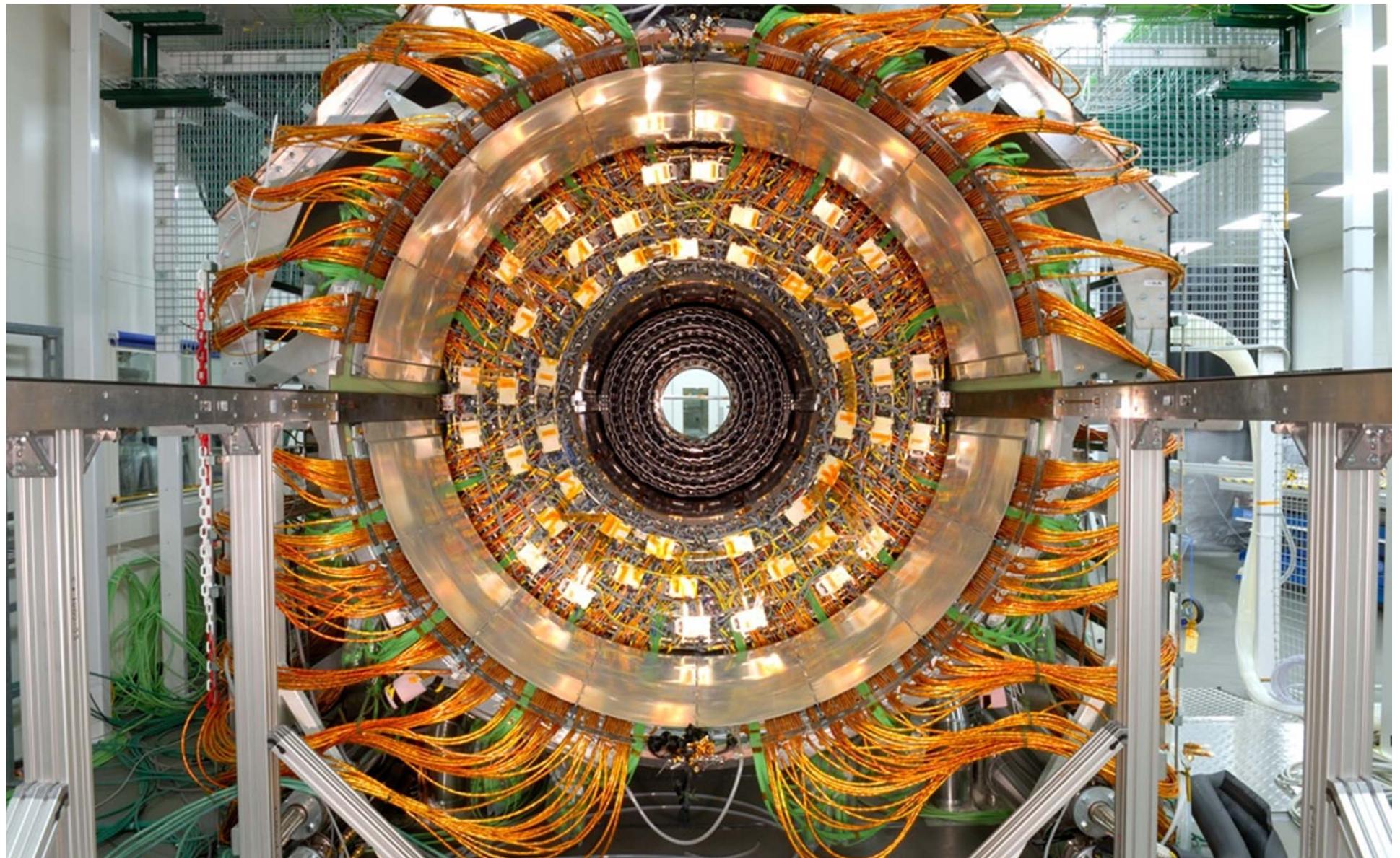
# How much data?

- Google processes 20 PB a day (2008)
- Wayback Machine has 3 PB + 100 TB/month (3/2009)
- Facebook has 2.5 PB of user data + 15 TB/day (4/2009)
- eBay has 6.5 PB of user data + 50 TB/day (5/2009)
- CERN's LHC will generate 15 PB a year (??)





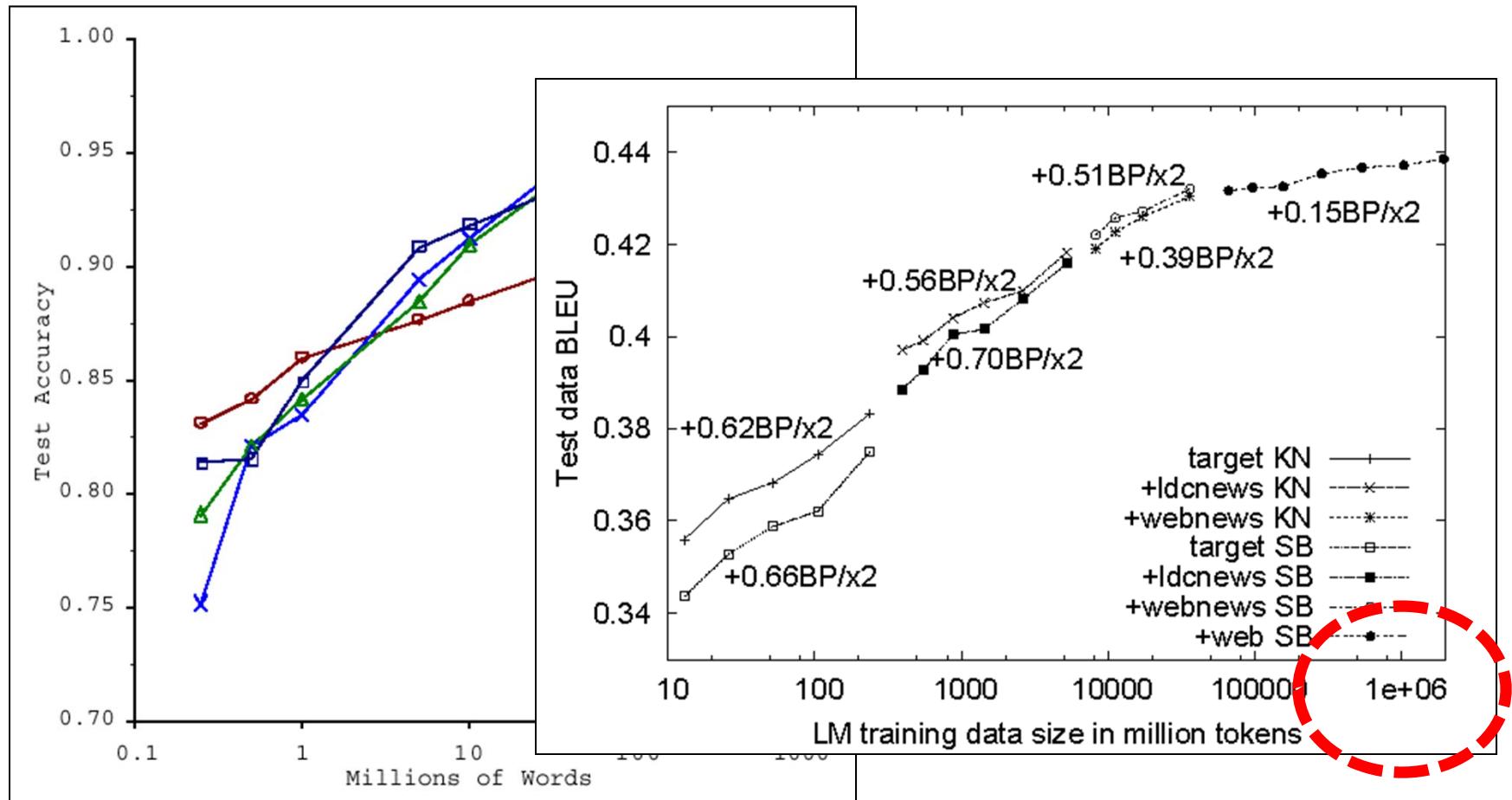
Maximilien Brice, © CERN



Maximilien Brice, © CERN

# No data like more data!

s/knowledge/data/g;



How do we get here if we're not Google?

# What to do with more data?

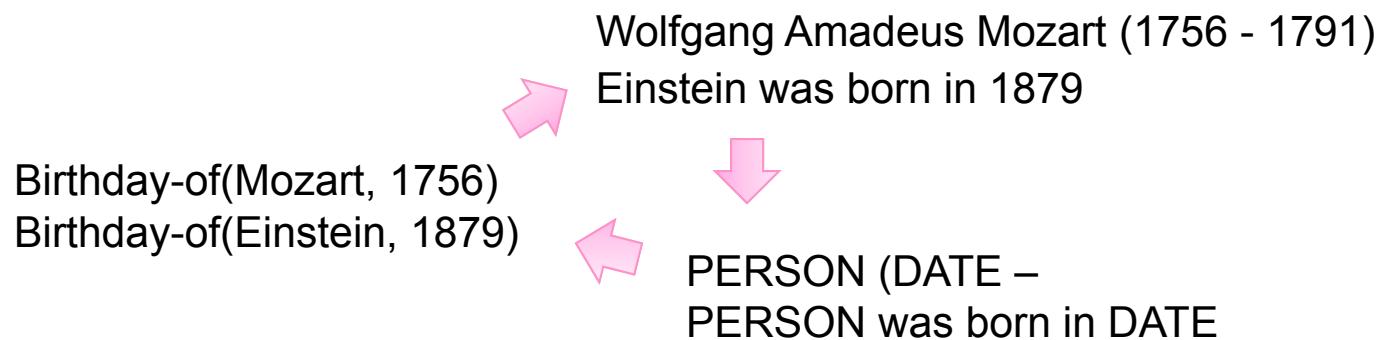
- Answering factoid questions

- Pattern matching on the Web
- Works amazingly well

**Who shot Abraham Lincoln? → X shot Abraham Lincoln**

- Learning relations

- Start with seed instances
- Search for patterns on the Web
- Using patterns to find more instances





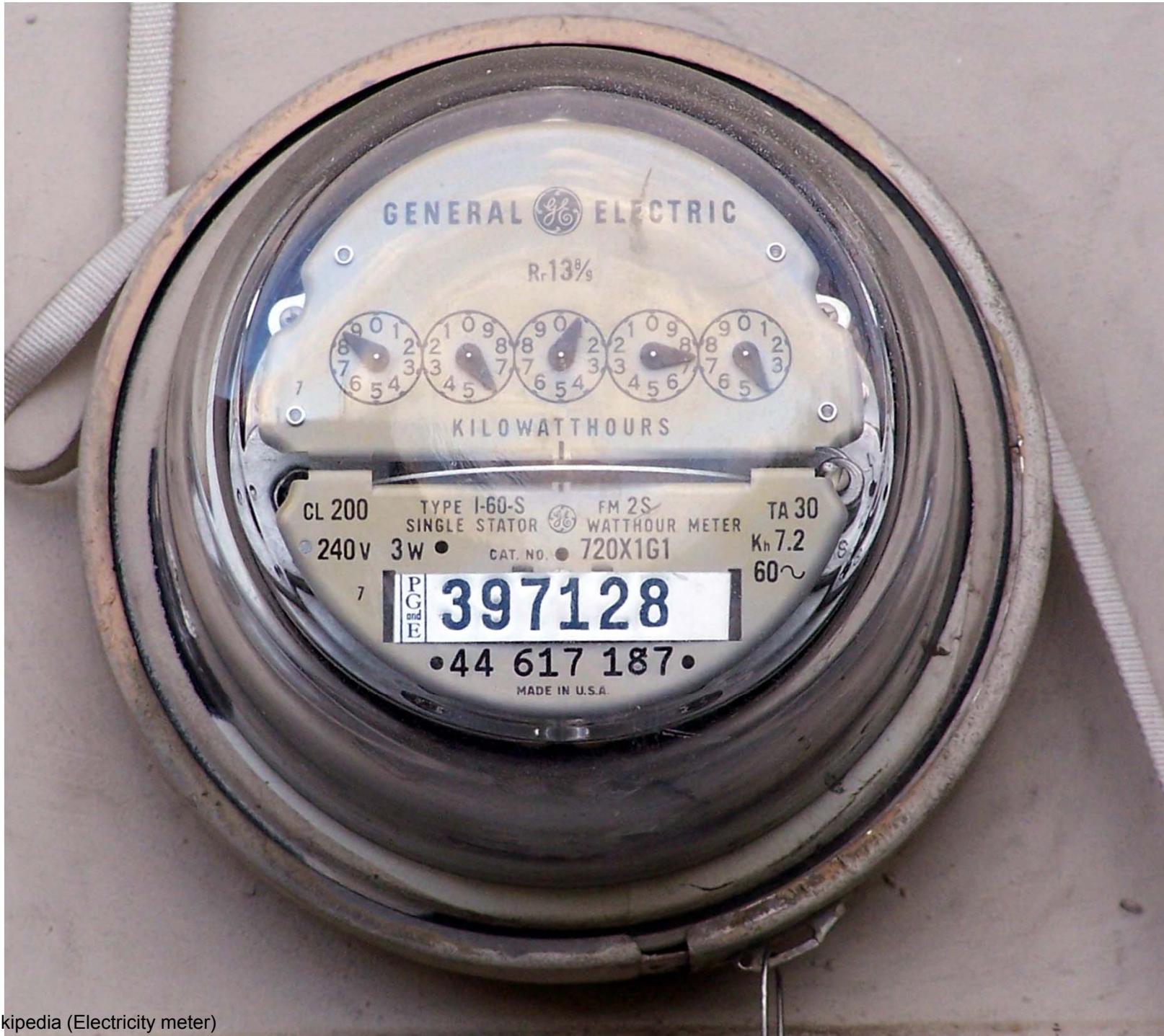
What is cloud computing?

# The best thing since sliced bread?

- Before clouds...
  - Grids
  - Vector supercomputers
  - ...
- Cloud computing means many different things:
  - Large-data processing
  - 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



Source: Wikipedia (Electricity meter)

# Utility Computing

- What?

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

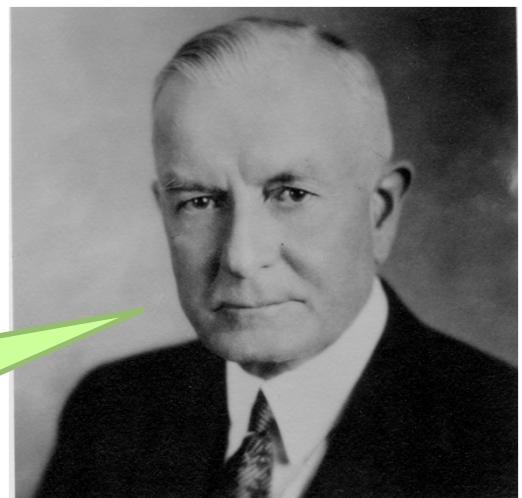
- Why?

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

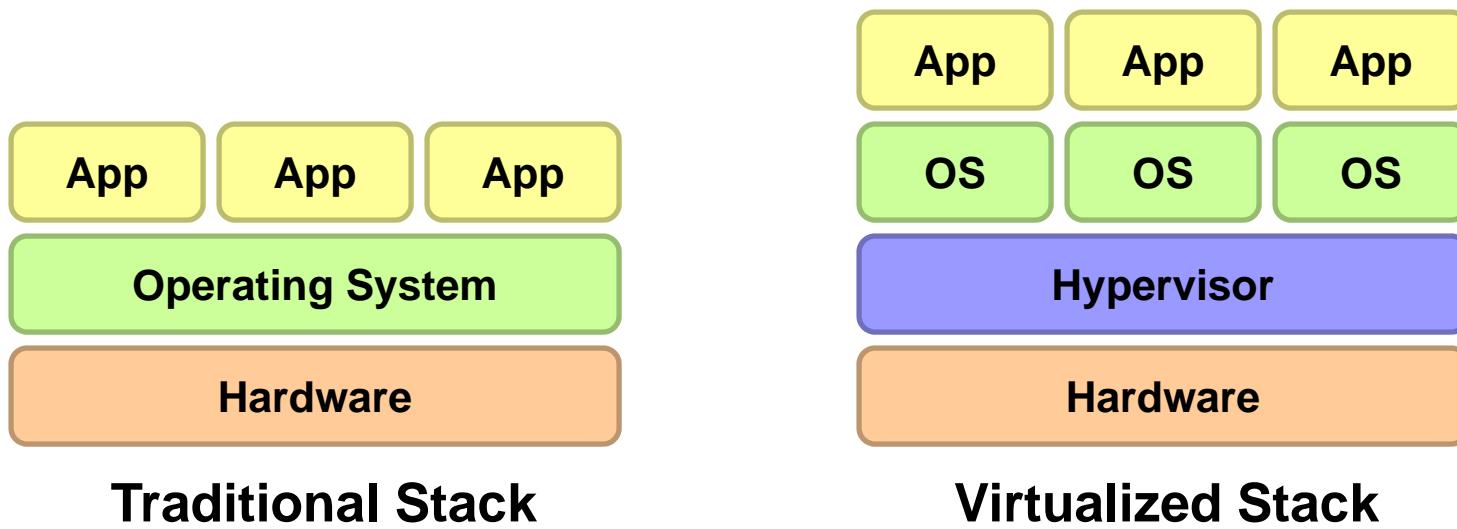
- Does it make sense?

- Benefits to cloud users
- Business case for cloud providers

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



# Enabling Technology: Virtualization



# 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?

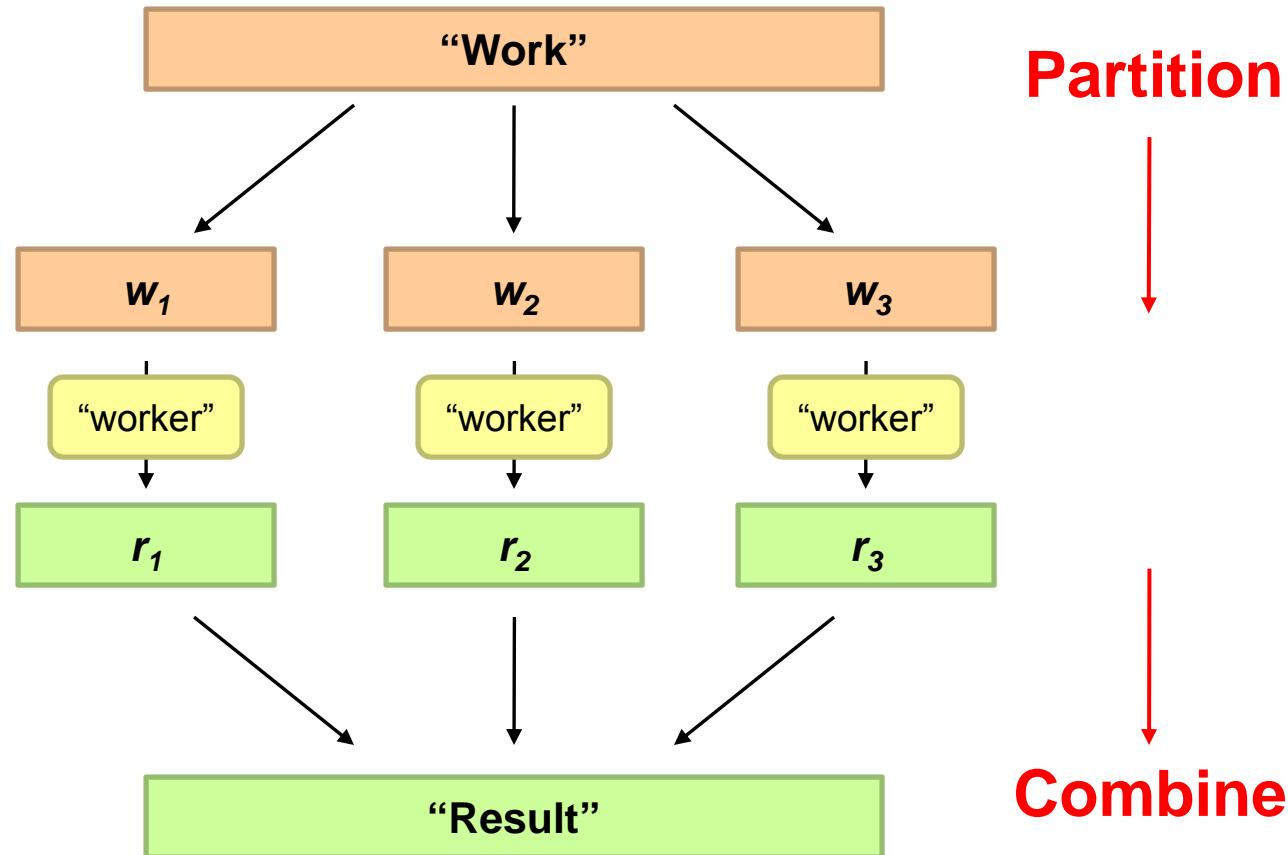
- Ready-made large-data problems
  - Lots of user-generated content
  - Even more user behavior data
  - Examples: Facebook friend suggestions, Google ad placement
  - Business intelligence: gather everything in a data warehouse and run analytics to generate insight
- Utility computing
  - Provision Hadoop clusters on-demand in the cloud
  - Lower barrier to entry for tackling large-data problem
  - Commoditization and democratization of large-data capabilities

# **How do we scale up?**



Source: Wikipedia (IBM Roadrunner)

# 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 share partial results?
- How do we aggregate partial results?
- How do we know all the workers have finished?
- What if workers die?

**What is the common theme of all of these problems?**

# Common Theme?

- Parallelization problems arise from:
  - Communication between workers (e.g., to exchange state)
  - Access to shared resources (e.g., data)
- Thus, we need a synchronization mechanism



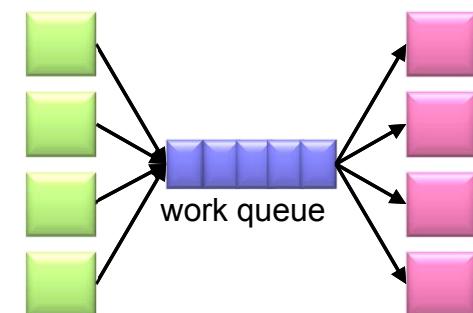
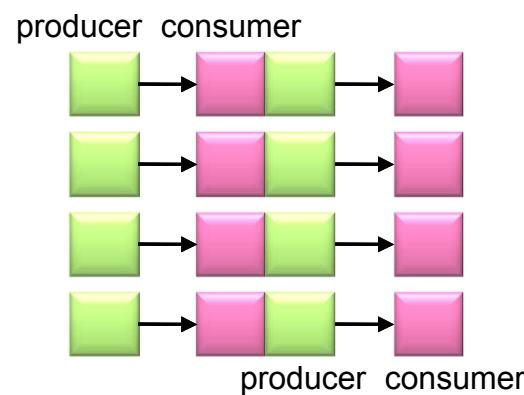
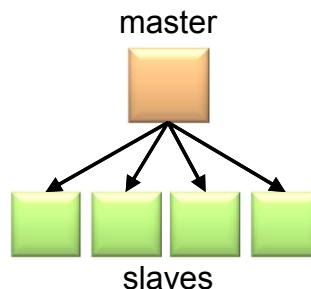
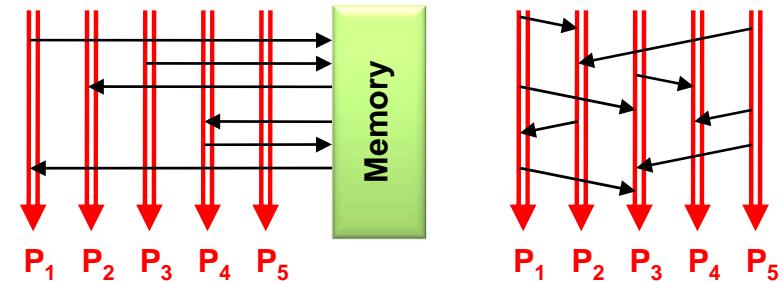
Source: Ricardo Guimarães Herrmann

# Managing Multiple Workers

- 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 the order in which workers access shared data
- Thus, we need:
  - Semaphores (lock, unlock)
  - Conditional variables (wait, notify, broadcast)
  - Barriers
- Still, lots of problems:
  - Deadlock, livelock, race conditions...
  - Dining philosophers, sleeping barbers, cigarette smokers...
- Moral of the story: be careful!

# Current Tools

- Programming models
  - Shared memory (pthreads)
  - Message passing (MPI)
- Design Patterns
  - Master-slaves
  - Producer-consumer flows
  - Shared work queues

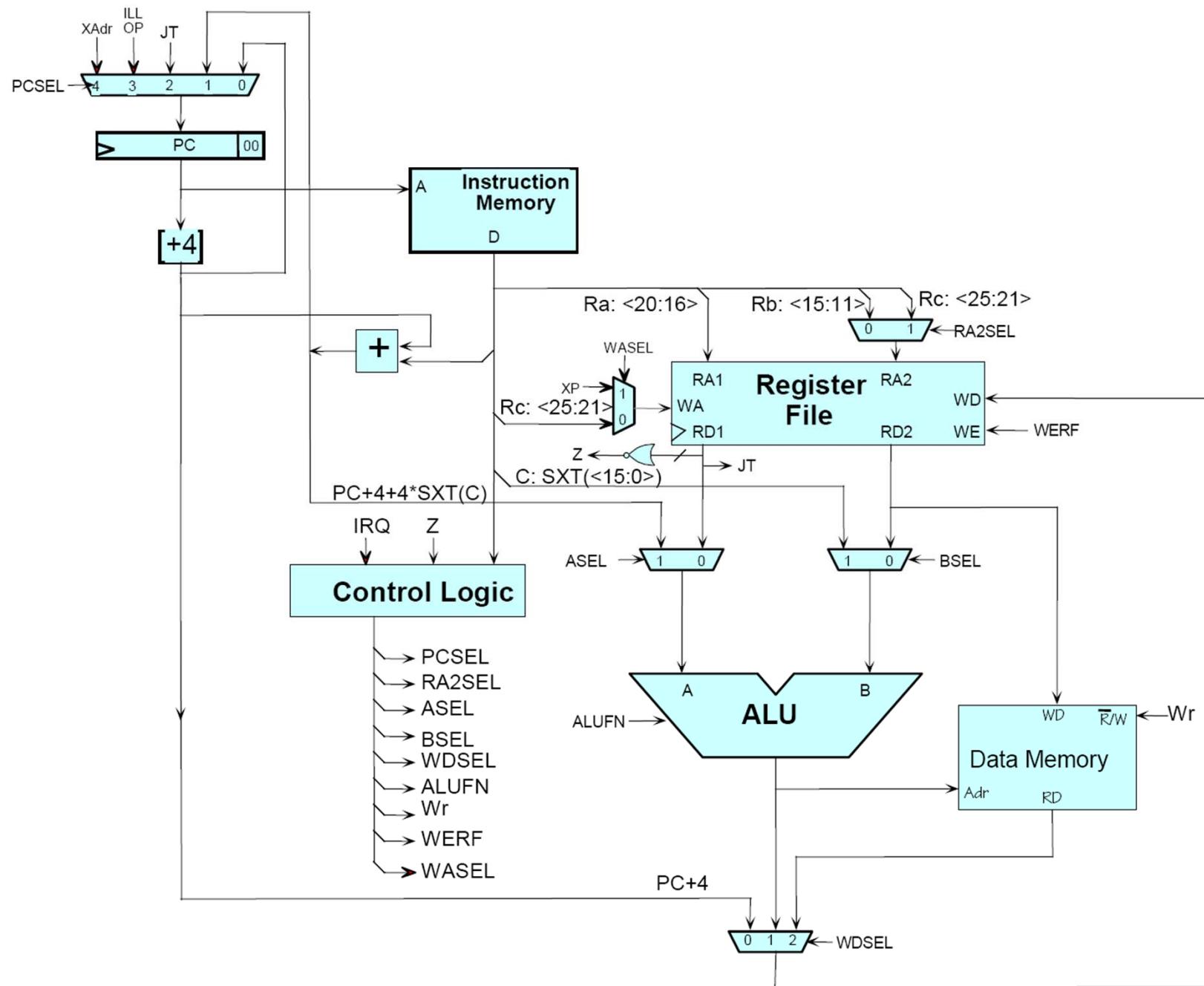


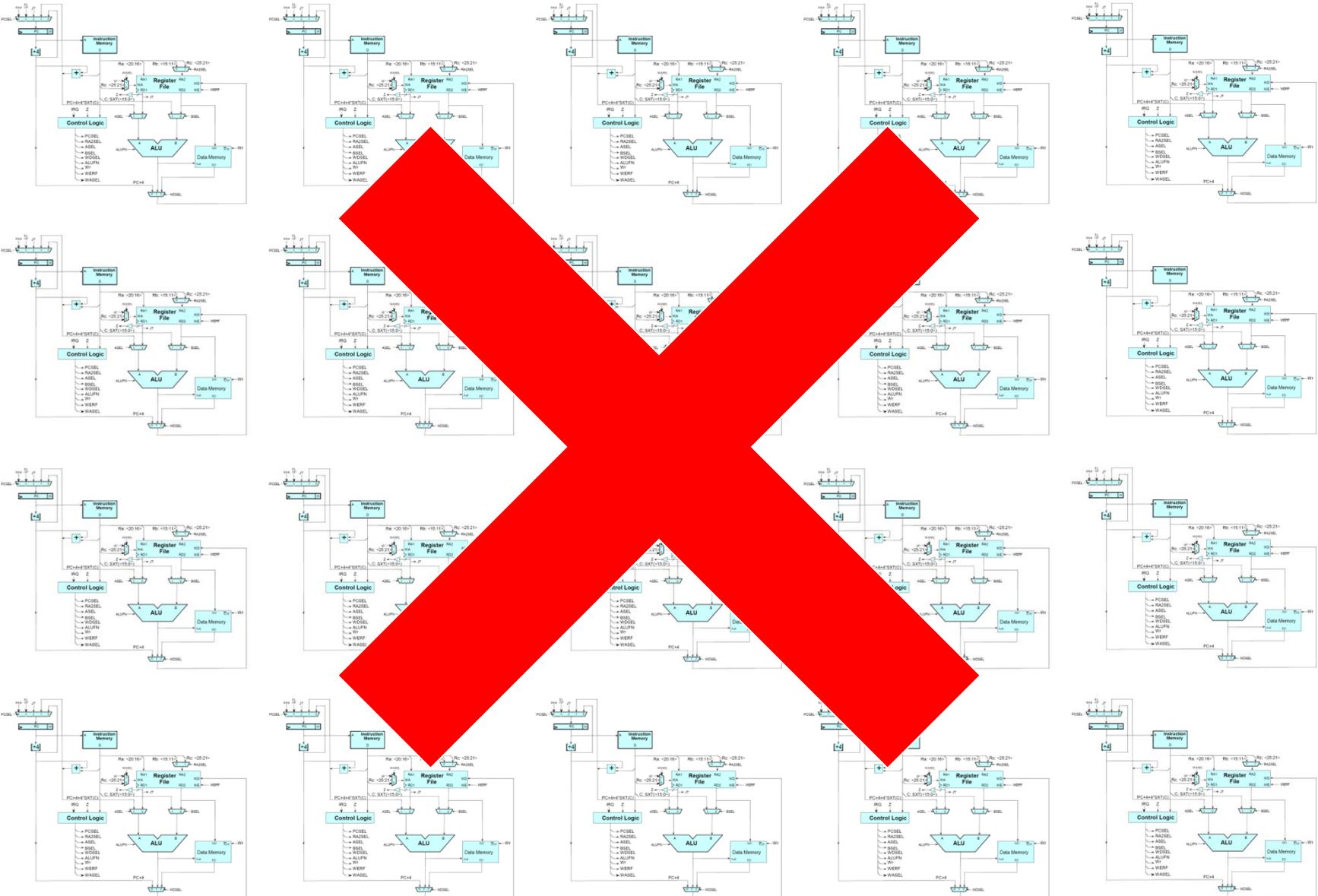
# Where the rubber meets the road

- Concurrency is difficult to reason about
- Concurrency is even more difficult to reason about
  - At the scale of datacenters (even 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



Source: Wikipedia (Flat Tire)

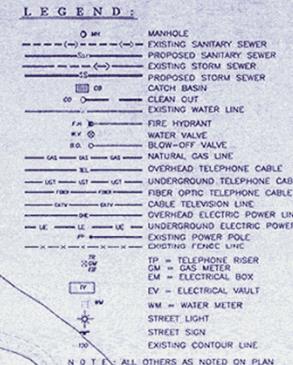




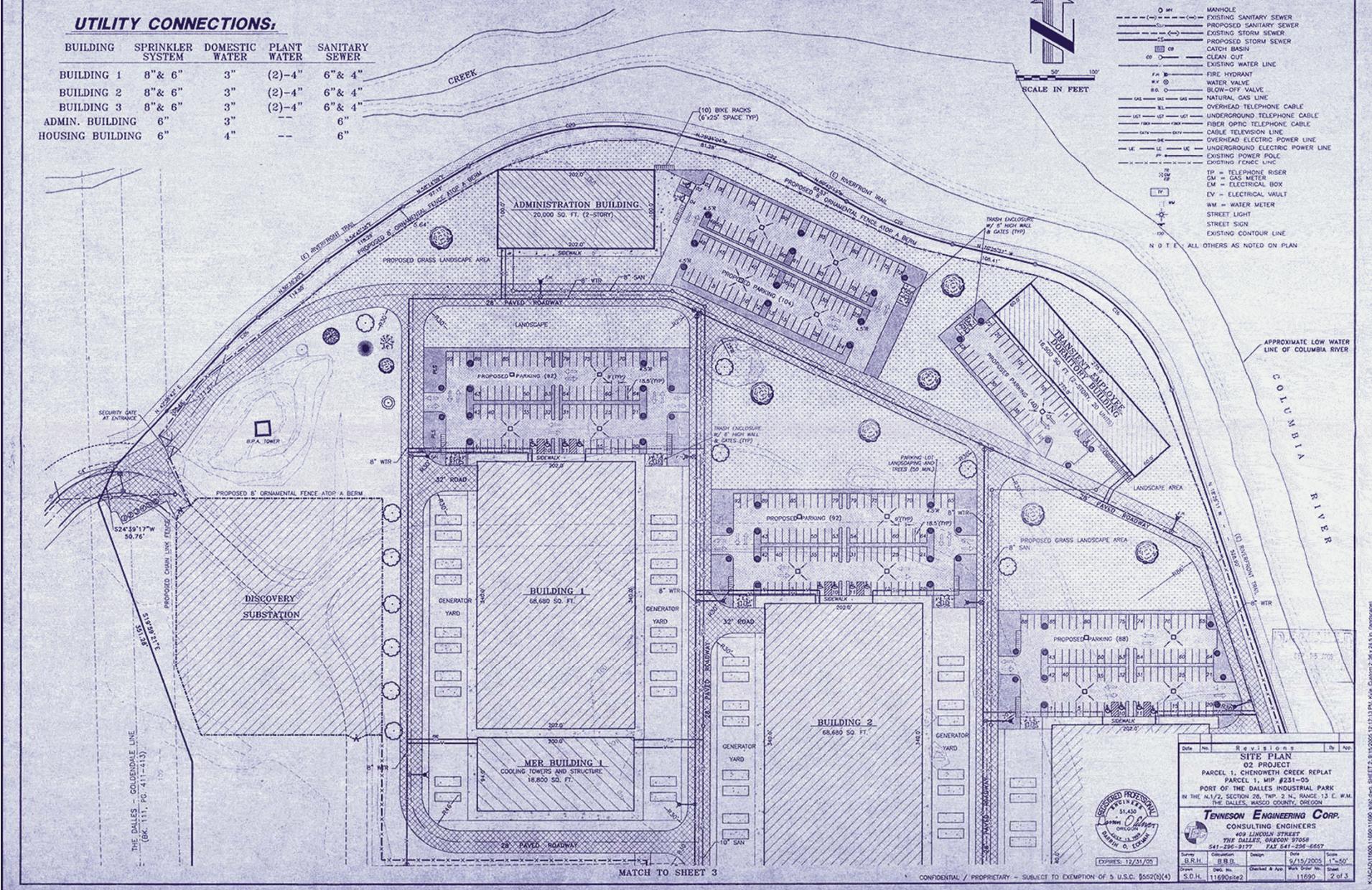
Source: MIT Open Courseware

## UTILITY CONNECTIONS:

BUILDING	SPRINKLER SYSTEM	DOMESTIC WATER	PLANT WATER	SANITARY SEWER
BUILDING 1	8" & 6"	3"	(2)-4"	6" & 4"
BUILDING 2	8" & 6"	3"	(2)-4"	6" & 4"
BUILDING 3	8" & 6"	3"	(2)-4"	6" & 4"
ADMIN. BUILDING	6"	3"	--	6"
HOUSING BUILDING	6"	4"	--	6"



~~NOTE: ALL OTHERS AS NOTED ON PLAN~~



Source: Harper's (Feb, 2008)

# What's the point?

- It's all about the right level of abstraction
  - The von Neumann architecture has served us well, but is no longer appropriate for the multi-core/cluster environment
- Hide system-level details from the developers
  - No more race conditions, lock contention, etc.
- Separating the *what* from *how*
  - Developer specifies the computation that needs to be performed
  - Execution framework (“runtime”) handles actual execution

**The datacenter *is* the computer!**

# “Big Ideas”

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

# MapReduce

# Typical Large-Data Problem

- Iterate over a large number of records

**Map** Extract something of interest from each

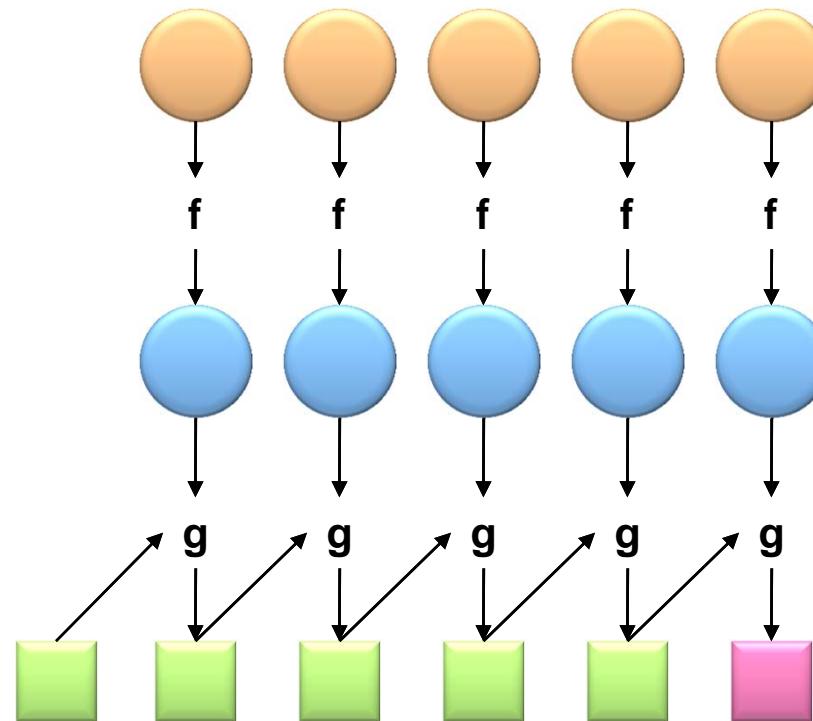
- Shuffle and sort intermediate results
- Aggregate intermediate results
- Generate final output

**Reduce**

**Key idea: provide a functional abstraction for these two operations**

# Roots in Functional Programming

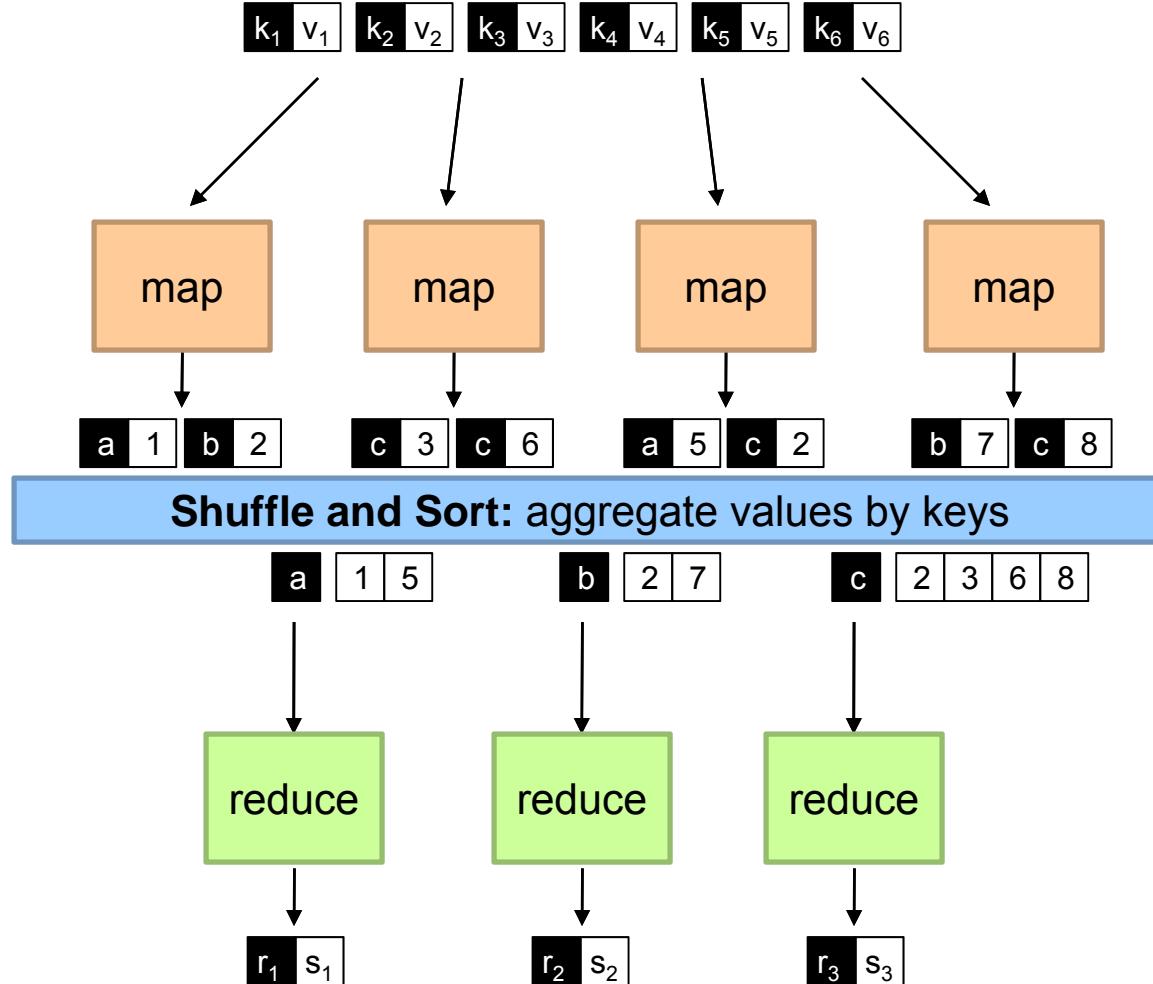
**Map**



**Fold**

# MapReduce

- Programmers specify two functions:
  - map** ( $k, v$ )  $\rightarrow \langle k', v' \rangle^*$
  - reduce** ( $k', v'$ )  $\rightarrow \langle k', v' \rangle^*$ 
    - All values with the same key are sent to the same reducer
- The execution framework handles everything else...



# MapReduce

- Programmers specify two functions:
  - map** ( $k, v$ )  $\rightarrow \langle k', v' \rangle^*$
  - reduce** ( $k', v'$ )  $\rightarrow \langle k', v' \rangle^*$ 
    - All values with the same key are sent to the same reducer
- The execution framework handles everything else...

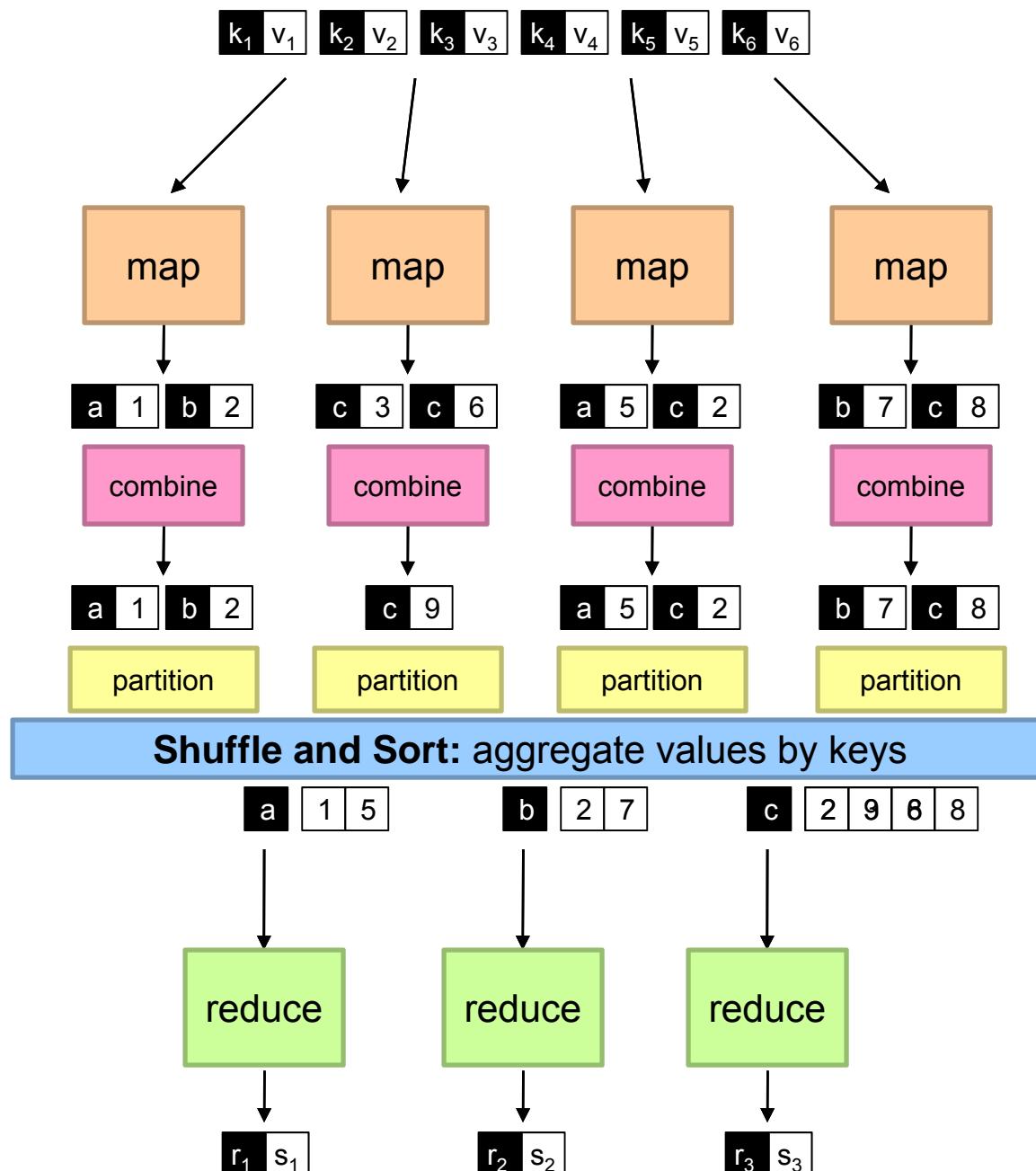
**What's “everything else”?**

# MapReduce “Runtime”

- 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 FS (later)

# MapReduce

- Programmers specify two functions:
  - map** ( $k, v \rightarrow \langle k', v' \rangle^*$
  - reduce** ( $k', v' \rightarrow \langle k', v' \rangle^*$ )
    - All values with the same key are reduced together
- The execution framework handles everything else...
- Not quite...usually, programmers also specify:
  - partition** ( $k', \text{number of partitions} \rightarrow \text{partition for } k'$ )
    - Often a simple hash of the key, e.g.,  $\text{hash}(k') \bmod n$
    - Divides up key space for parallel reduce operations
  - combine** ( $k', v' \rightarrow \langle k', v' \rangle^*$ )
    - Mini-reducers that run in memory after the map phase
    - Used as an optimization to reduce network traffic



## Two more details...

- Barrier between map and reduce phases
  - But we can begin copying intermediate data earlier
- Keys arrive at each reducer in sorted order
  - No enforced ordering across reducers

# “Hello World”: Word Count

**Map(String docid, String text):**

for each word w in text:

    Emit(w, 1);

**Reduce(String term, Iterator<Int> values):**

    int sum = 0;

    for each v in values:

        sum += v;

    Emit(term, value);

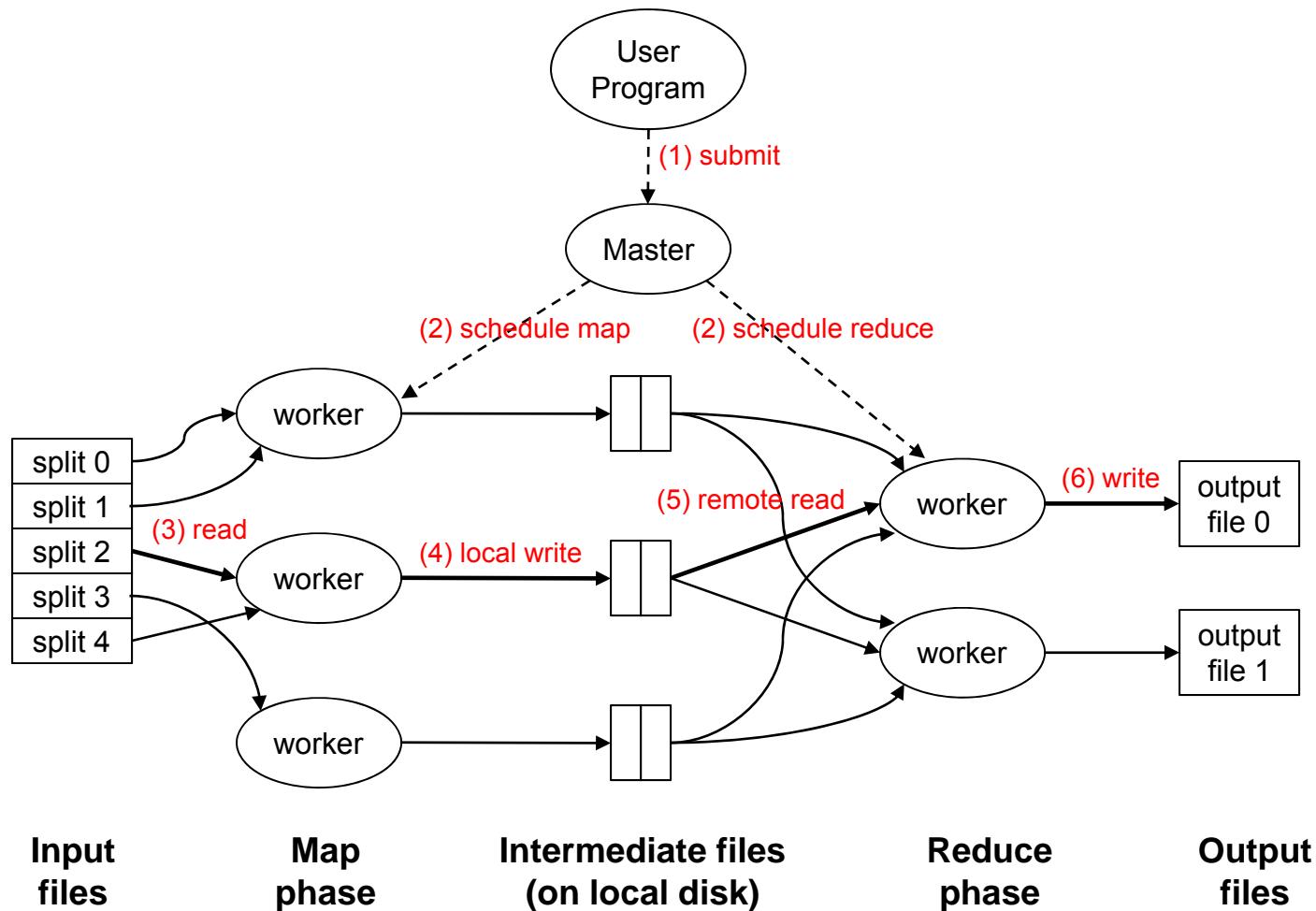
# MapReduce can refer to...

- The programming model
- The execution framework (aka “runtime”)
- The specific implementation

**Usage is usually clear from context!**

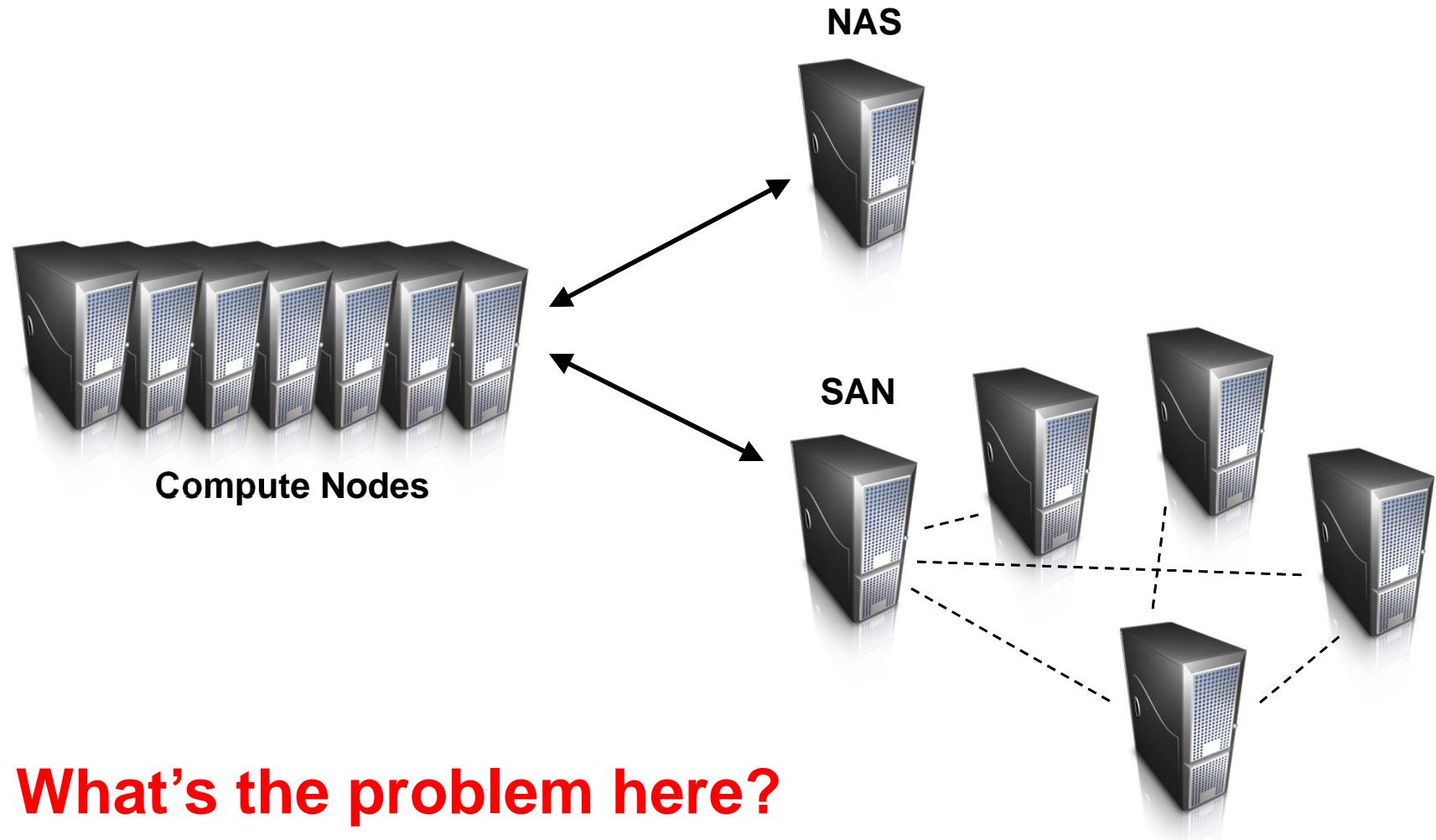
# MapReduce Implementations

- Google has a proprietary implementation in C++
  - Bindings in Java, Python
- Hadoop is an open-source implementation in Java
  - Development led by Yahoo, used in production
  - Now an Apache project
  - Rapidly expanding software ecosystem
- Lots of custom research implementations
  - For GPUs, cell processors, etc.



Adapted from (Dean and Ghemawat, OSDI 2004)

# How do we get data to the workers?



**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?
  - 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)

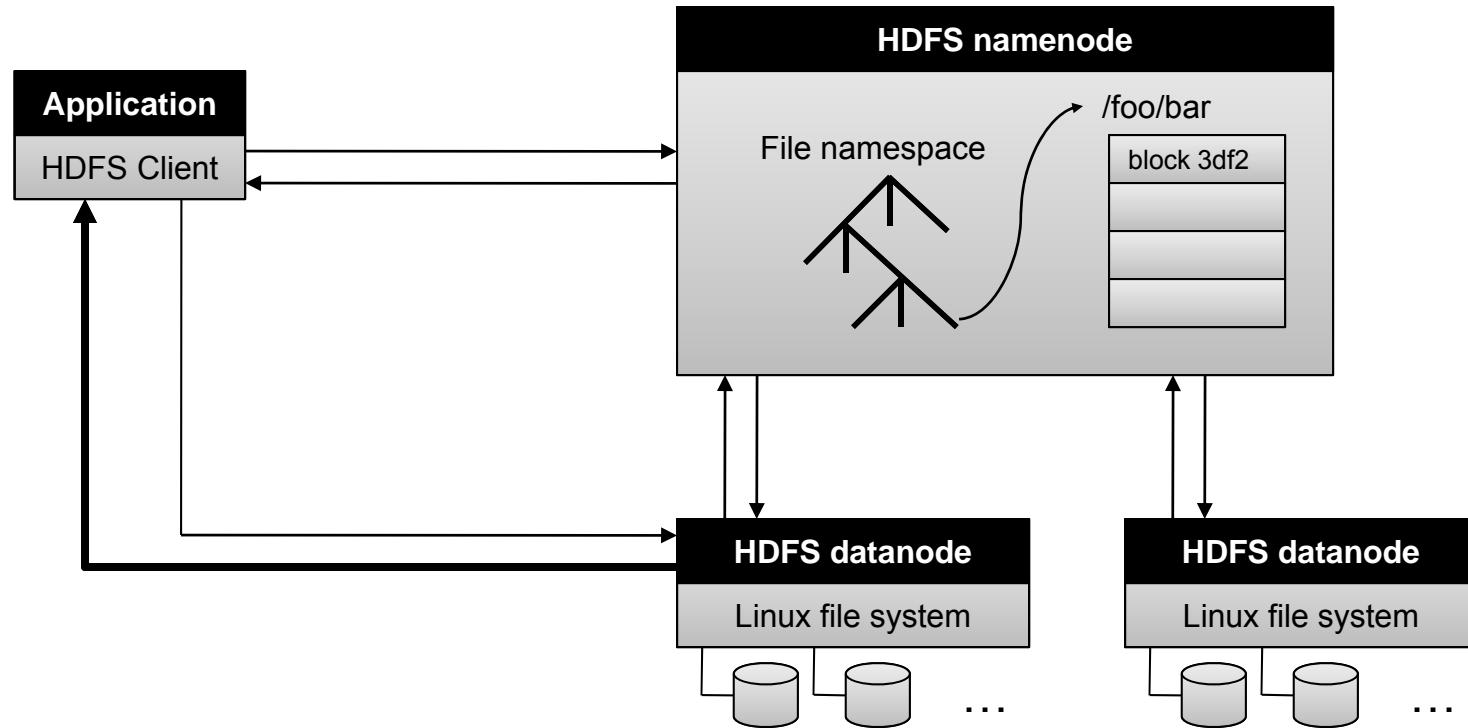
**HDFS = GFS clone (same basic ideas)**

# From GFS to HDFS

- Terminology differences:
  - GFS master = Hadoop namenode
  - GFS chunkservers = Hadoop datanodes
- Functional differences:
  - No file appends in HDFS (planned feature)
  - HDFS performance is (likely) slower

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

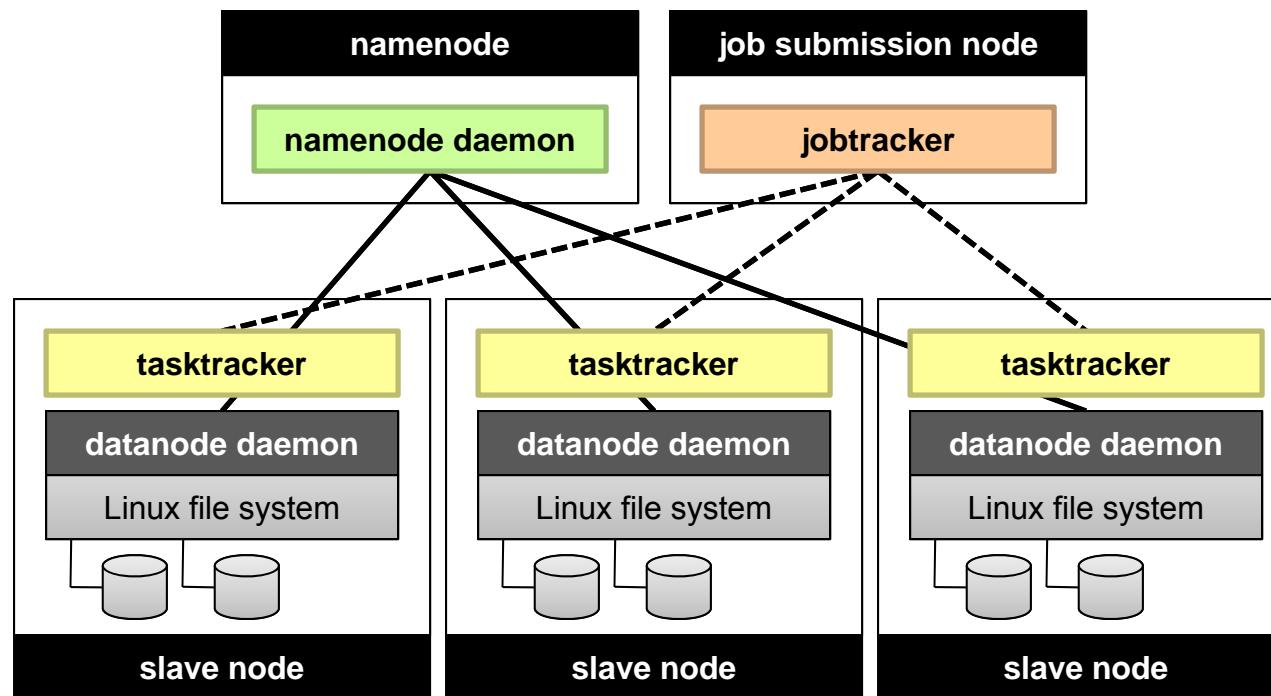
# HDFS Architecture



# 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

# Putting everything together...



# Recap

- Why large data?
- Cloud computing and MapReduce
- Large-data processing: “big ideas”
- What is MapReduce?
- Importance of the underlying distributed file system

A photograph of a tropical sunset. The sky is filled with clouds that are lit from below by the setting sun, creating a warm orange and yellow glow against a darker blue and purple background. In the foreground, the dark silhouette of palm trees stands on a beach. A few small figures of people are visible on the sand near the water's edge. The ocean is calm with some small waves. The overall atmosphere is peaceful and scenic.

# Questions?

Photo credit: Jimmy Lin