

# Big Data Infrastructure

## Session I: Introduction

Jimmy Lin  
University of Maryland  
Monday, January 26, 2015



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 this course about?

- What is big data?
- Why big data?
- *Infrastructure* for big data



# From the Ivory Tower...



**... to building sh\*t that works**



**... and back.**



**Big Data**

# Google™

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)

# YAHOO!

Hadoop: 365 PB, 330K nodes (6/2014)

# JPMorganChase

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

# ebay

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

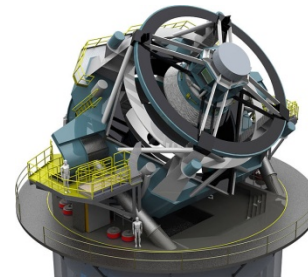
300 PB data in Hive + 600 TB/day (4/2014)

# facebook

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

# amazon web services™

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?

# Why big data?

- Science
- Engineering
- Commerce







# Science

Emergence of the 4<sup>th</sup> Paradigm

Data-intensive e-Science

# Engineering

The unreasonable effectiveness of data

Count and normalize!



Know thy customers

Data → Insights → Competitive advantages

# Commerce





Why big data?  
Infrastructure for big data

# Course Administrivia



# My Expectations

- You're already a good Java programmer
  - This course does *not* teach programming
  - You're expected to pick up Hadoop with minimal help
- You're good at debugging
  - Your own code
  - Compiling, patching, and installing open source software
- You have basic knowledge of:
  - Probability and statistics, discrete math
  - Computer architecture

# How will I actually learn Hadoop?

- Hadoop: The Definitive Guide
- RTFM
- RTFC(!)

# **This course is not for you...**

- If you're not genuinely interested in the topic
- If you can't put in the time
- If you're uncomfortable with the uncertainty, unpredictability, etc. that comes with immature software

**Otherwise, this will be a rewarding and fun course!**



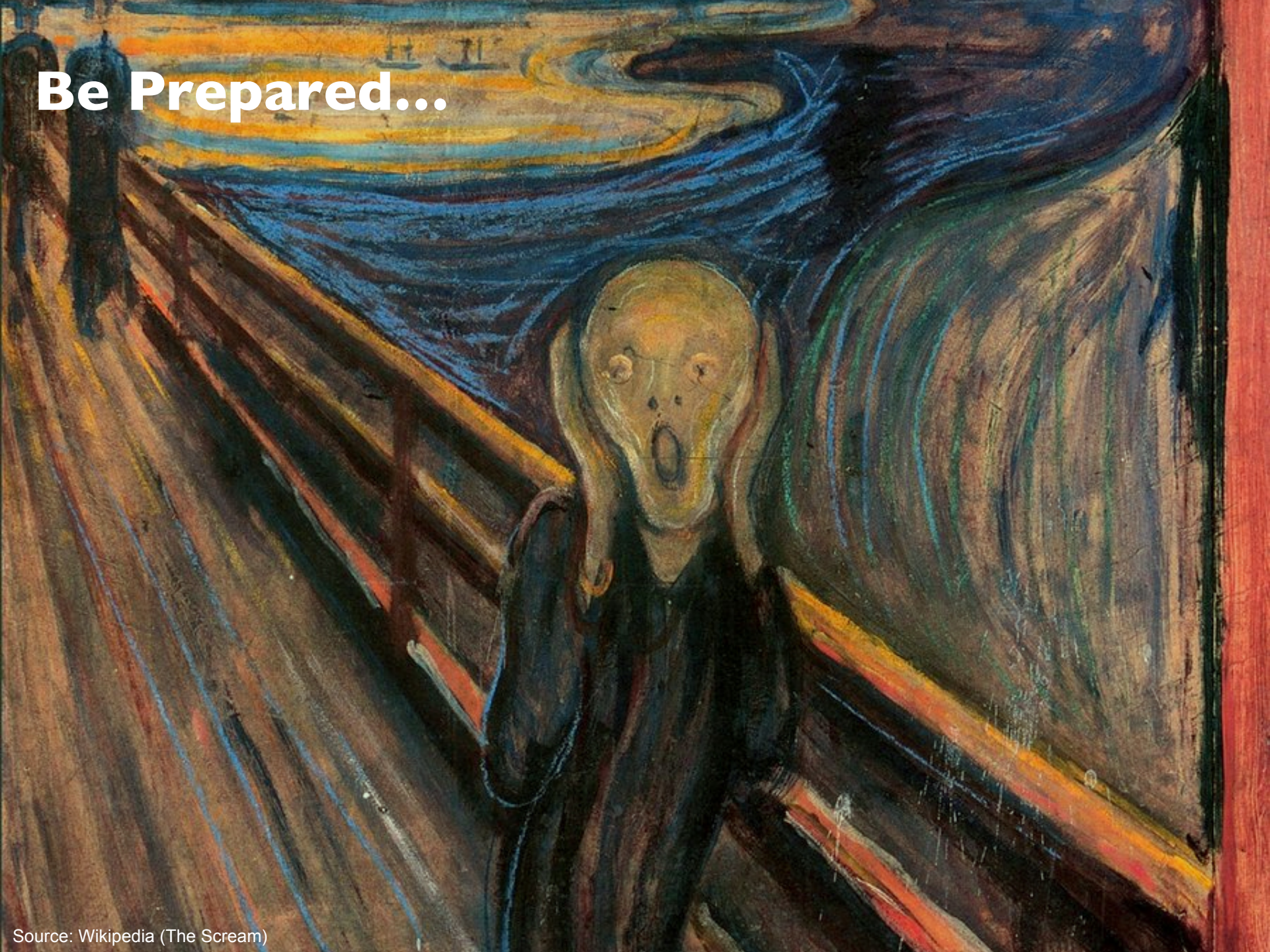
# Details, Details...

- Make sure you're on the mailing list!
- Textbooks
- Components of the final grade:
  - Assignments
  - Final exam
  - Final project
- I am unlikely to accept the following excuses:
  - “Too busy”
  - “It took longer than I thought it would take”
  - “It was harder than I initially thought”
  - “My dog ate my homework” and modern variants thereof

# Hadoop Resources

- Hadoop on your local machine
- Hadoop in a virtual machine on your local machine
- Hadoop on a UMIACS cluster

**Be Prepared...**



# “Hadoop Zen”

- Parts of the ecosystem are still immature
  - We’ve come a long way since 2007, but still far to go...
  - Bugs, undocumented “features”, inexplicable behavior, etc.
- Don’t get frustrated (take a deep breath)...
  - Those W\$\*#T@F! moments
- Be patient...
  - We will inevitably encounter “situations” along the way
- Be flexible...
  - We will have to be creative in workarounds
- Be constructive...
  - Tell me how I can make everyone’s experience better

# “Hadoop Zen”



An aerial photograph showing a vast, dense layer of white, fluffy clouds stretching across the horizon. The clouds are illuminated from the side, creating soft shadows and highlights. The sky above is a clear, deep blue. The overall scene is serene and expansive.

# **Interlude: 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



GENERAL  ELECTRIC

Rr13<sup>8</sup>/<sub>9</sub>



KILOWATTHOURS

CL 200

240V

3W

TYPE I-60-S  
SINGLE STATOR

CAT. NO.



FM 2S  
WATTHOUR METER

720X1G1

TA 30

Kh 7.2

60~

7  
P  
E  
R  
G  
E

397128

•44 617 187•

MADE IN U.S.A.

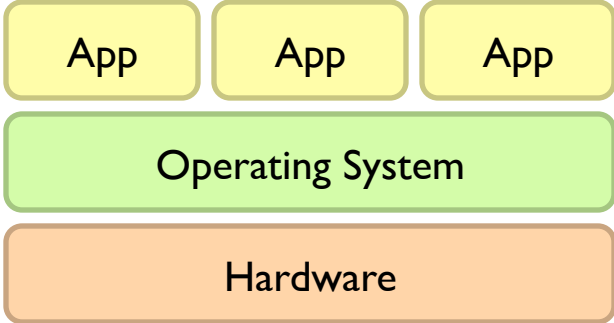
# Utility Computing

- 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

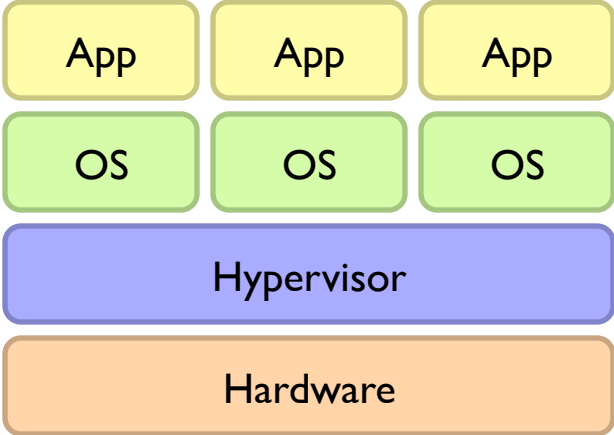
I think there is a world market for about five computers.



# 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

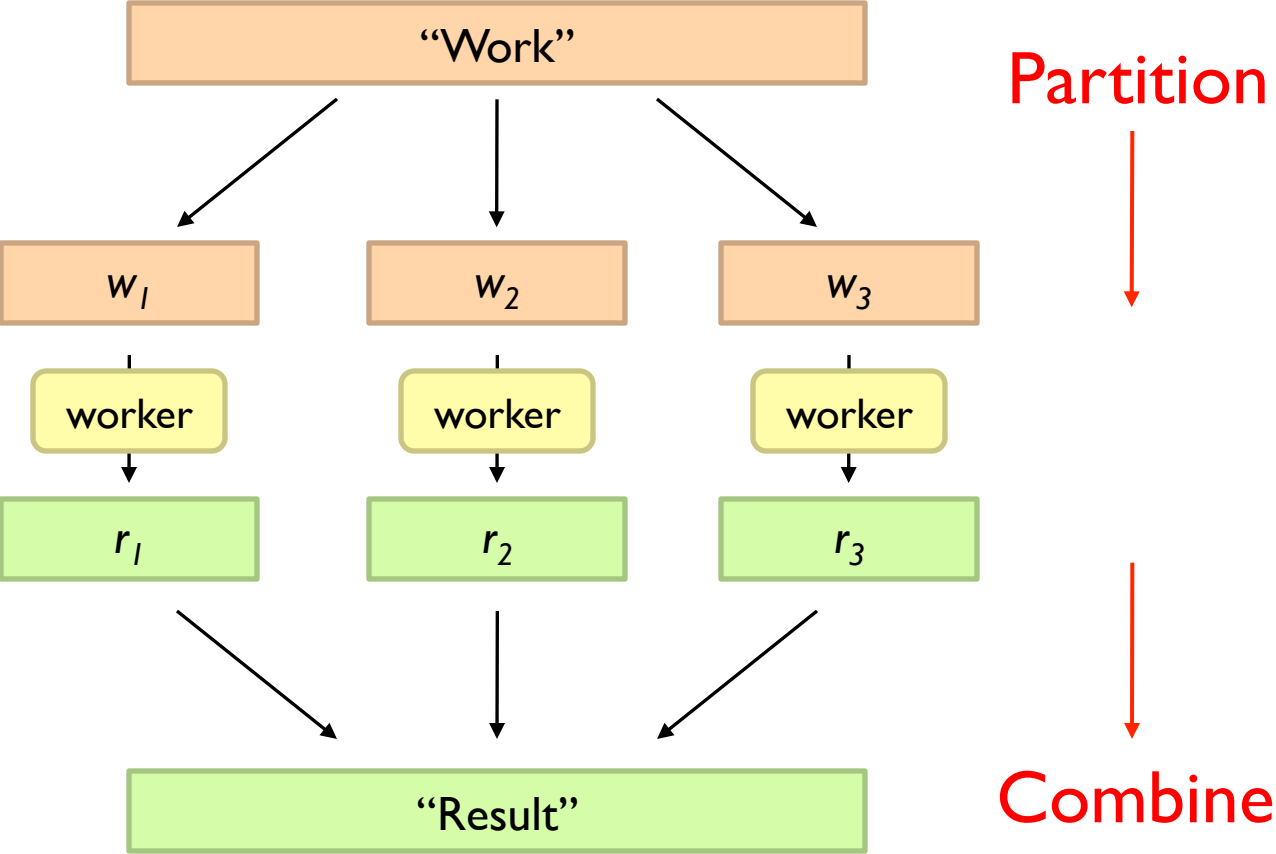
# Who cares?

- 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

# 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. A complex network of metal pipes and cables runs across the ceiling and floor, creating a grid-like pattern. The lighting is a mix of cool blue and warm yellow, creating a futuristic and industrial atmosphere. The ceiling is high and features a complex steel truss structure with various pipes and conduits. The floor is made of large, light-colored tiles. The overall scene conveys a sense of scale and technological complexity.

# 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's 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



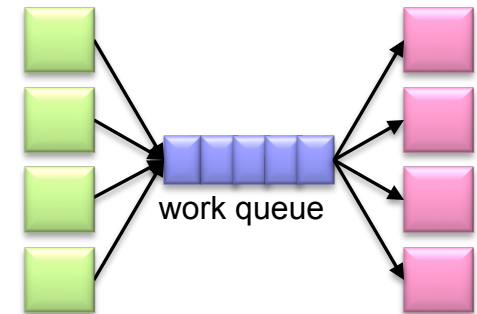
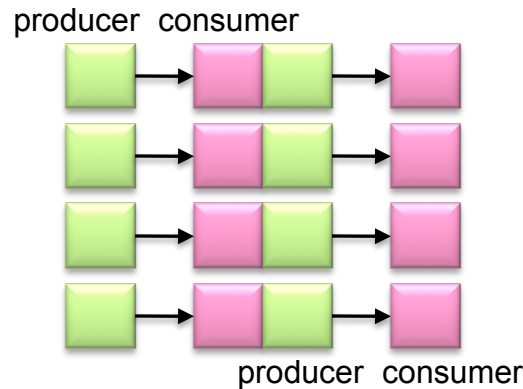
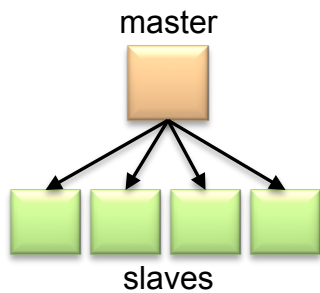
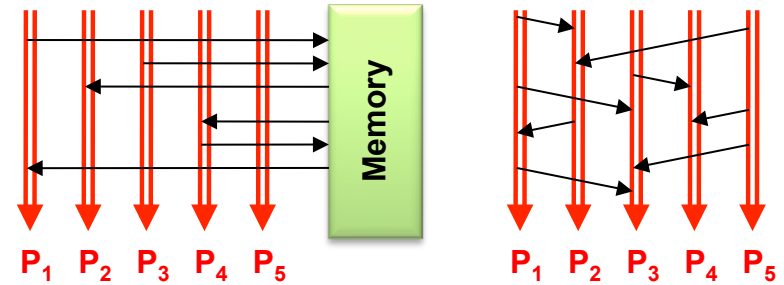
Source: Ricardo Guimarães Herrmann

# 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 when workers need to communicate partial results
  - 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

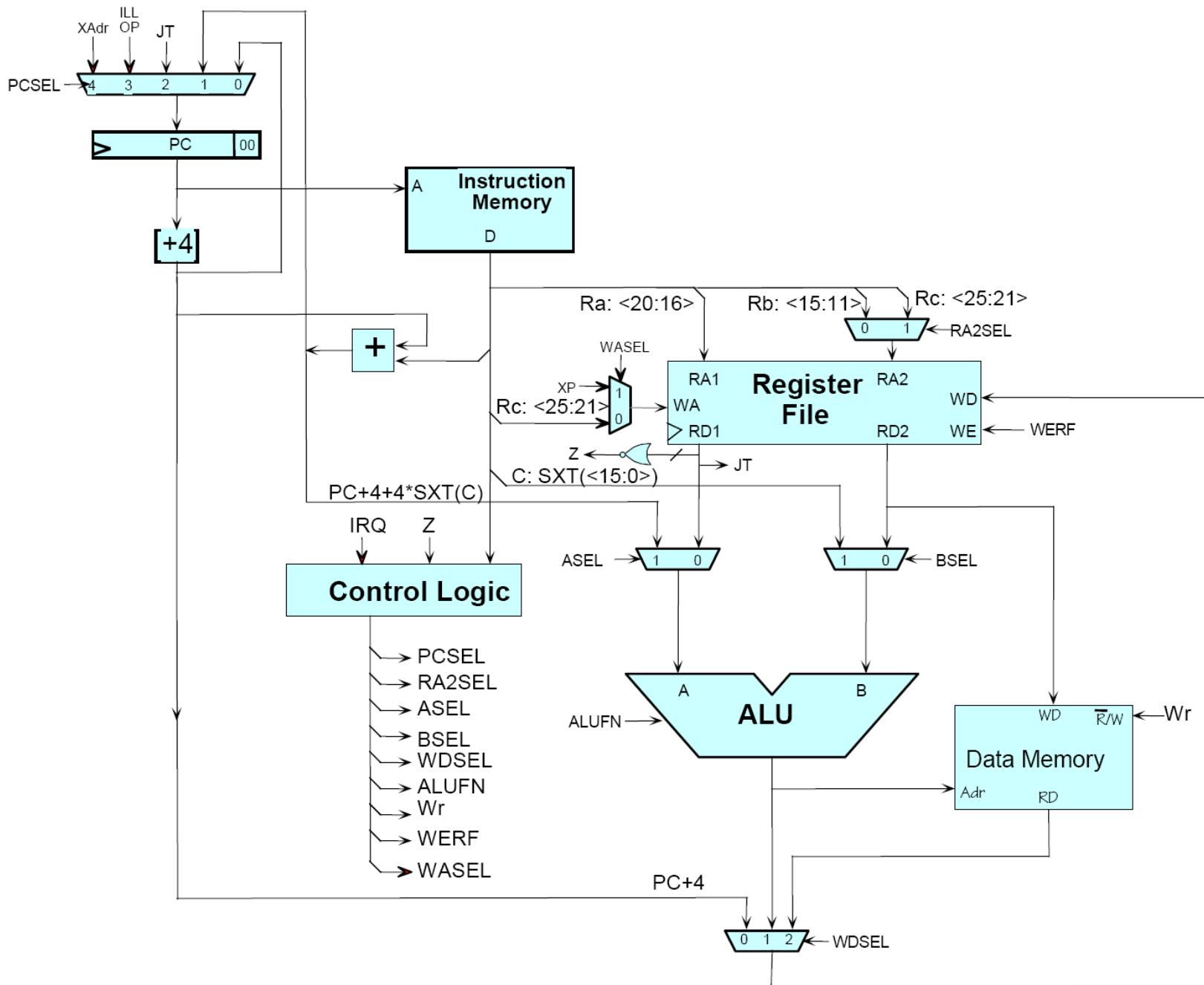


# Where the rubber meets the road

- Concurrency is difficult to reason about
- Concurrency is even more difficult to reason about
  - At the scale of datacenters 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 you own dedicated library, then program with it
  - Burden on the programmer to explicitly manage everything



Source: Wikipedia (Flat Tire)



Source: MIT Open Courseware





An aerial photograph of a large datacenter facility during sunset. The sun is low on the horizon, casting a warm orange glow over the scene. The datacenter consists of several large, white, rectangular buildings arranged in a grid-like pattern. In the foreground, there are several large white storage tanks or containers. The facility is surrounded by green fields and a few scattered trees. The overall atmosphere is serene and industrial.

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



Source: Wikipedia (The Dalles, Oregon)



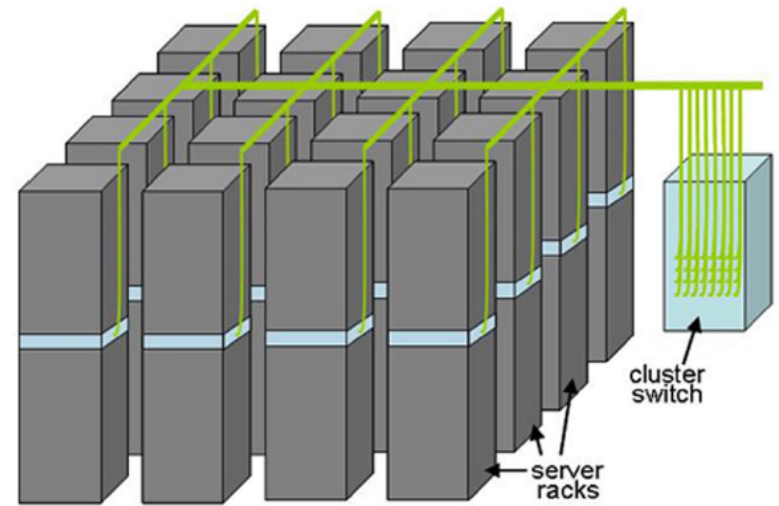
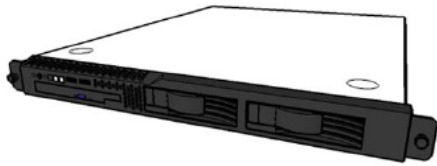


Source: Google



Source: Bonneville Power Administration

# Building Blocks





Source: Google



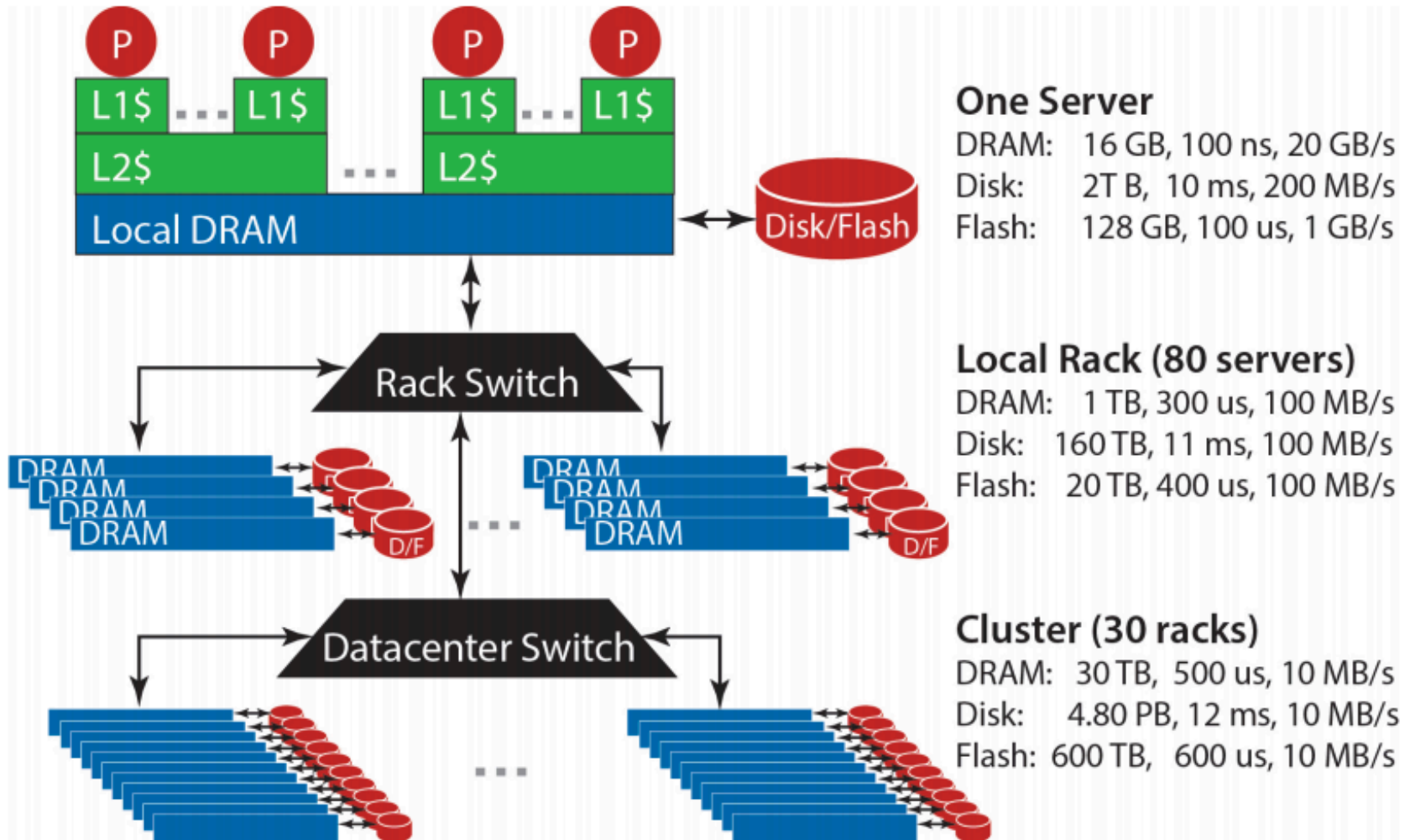
Source: Google



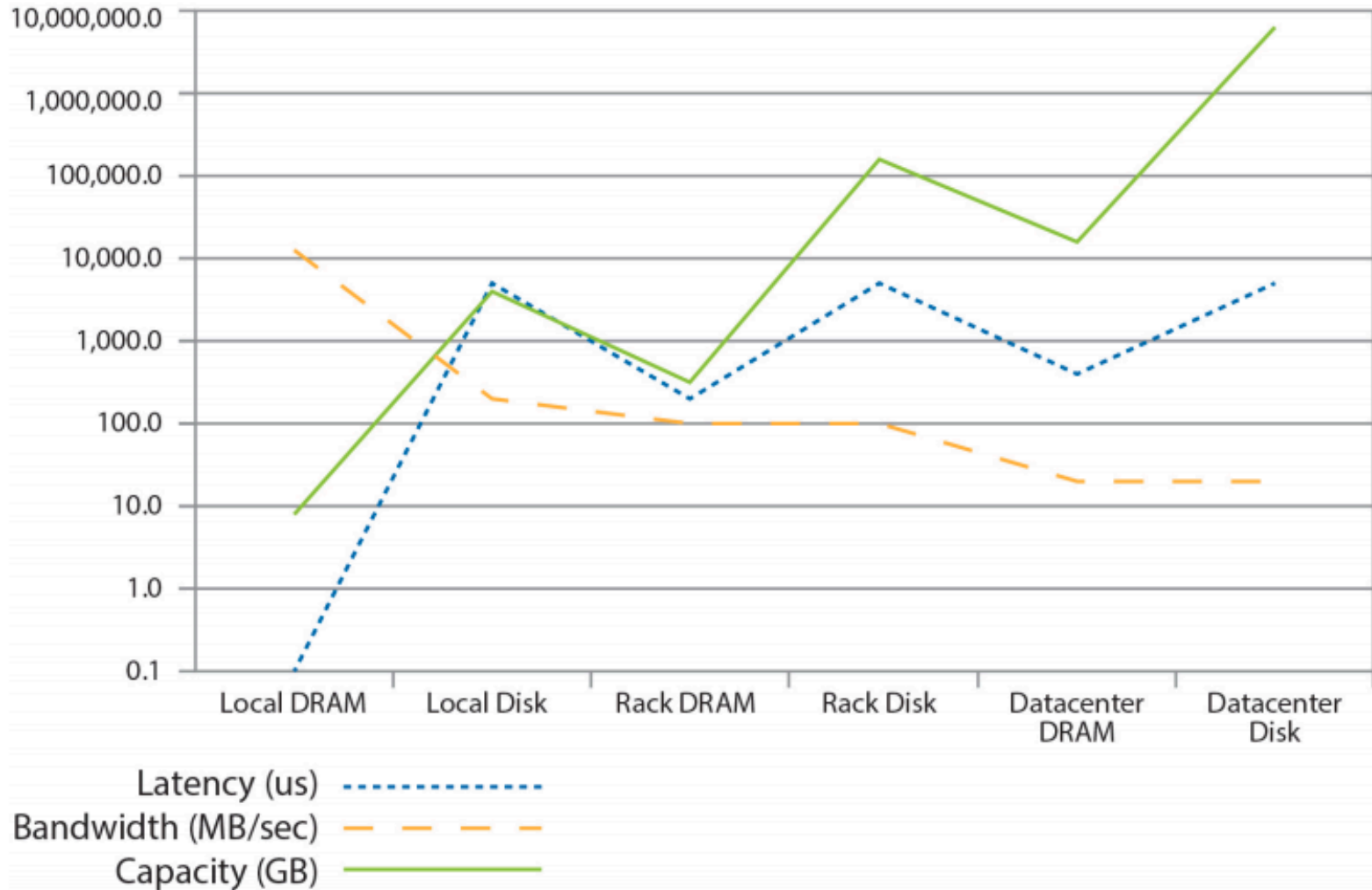


Source: Facebook

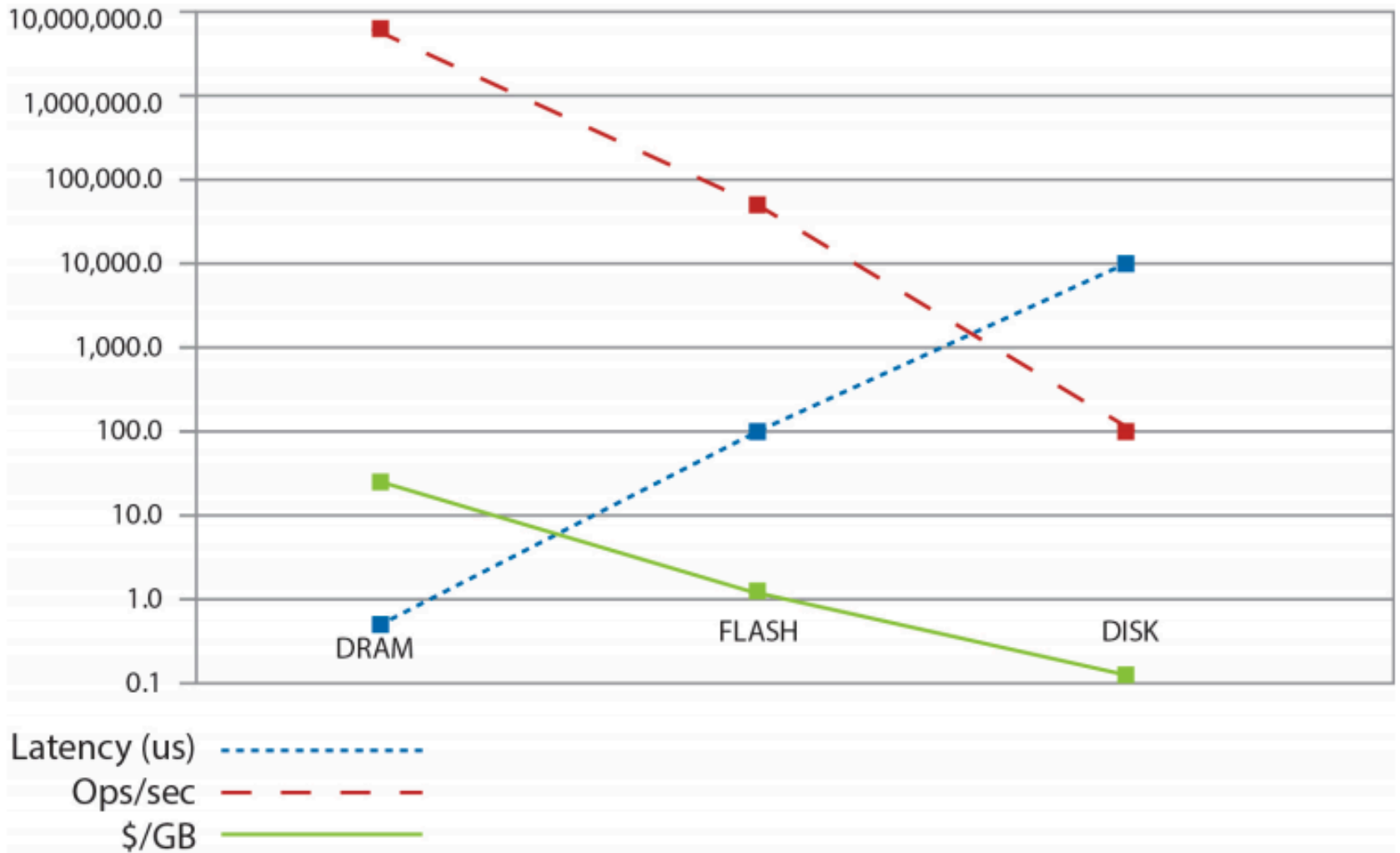
# Storage Hierarchy



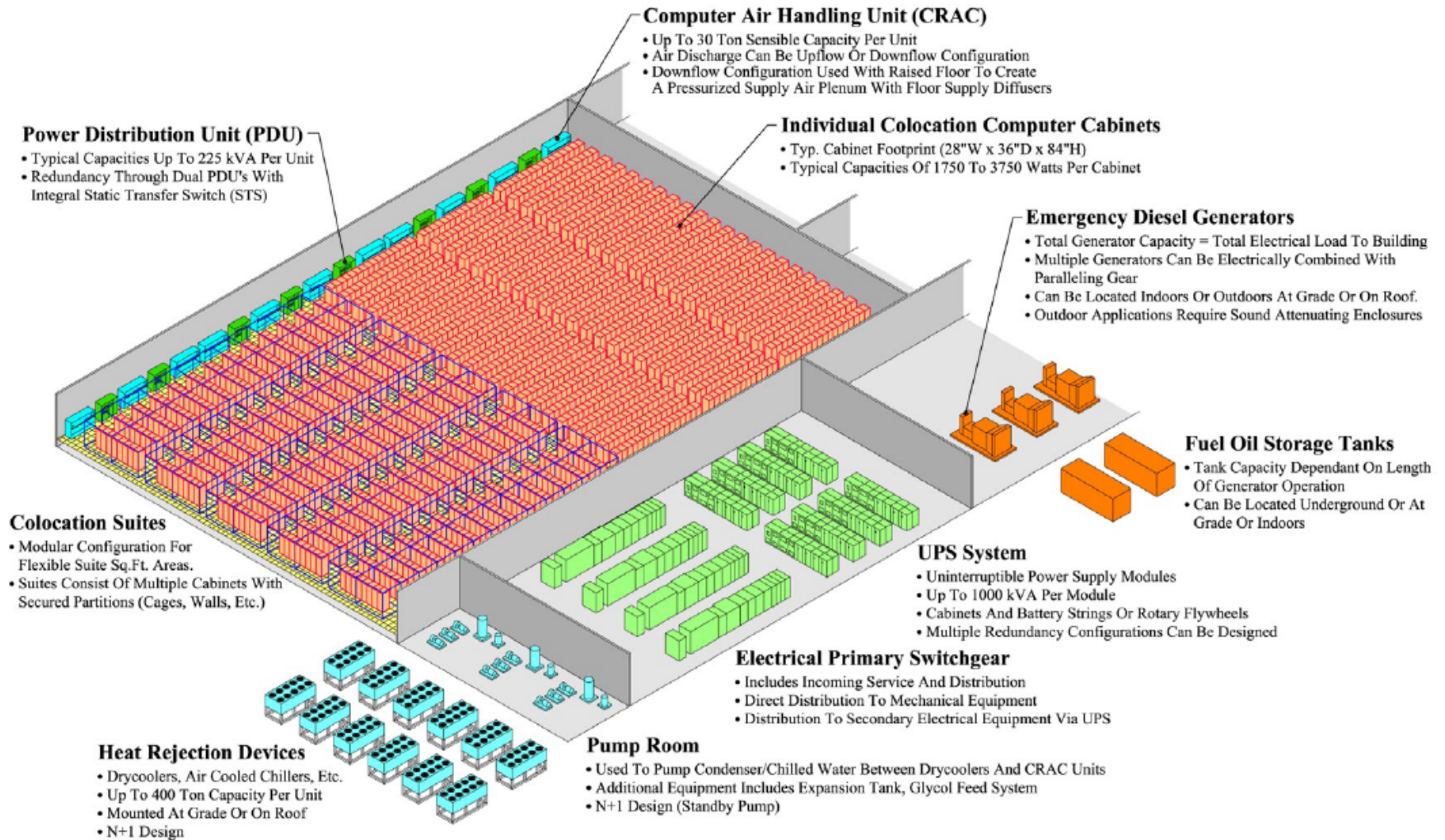
# Storage Hierarchy



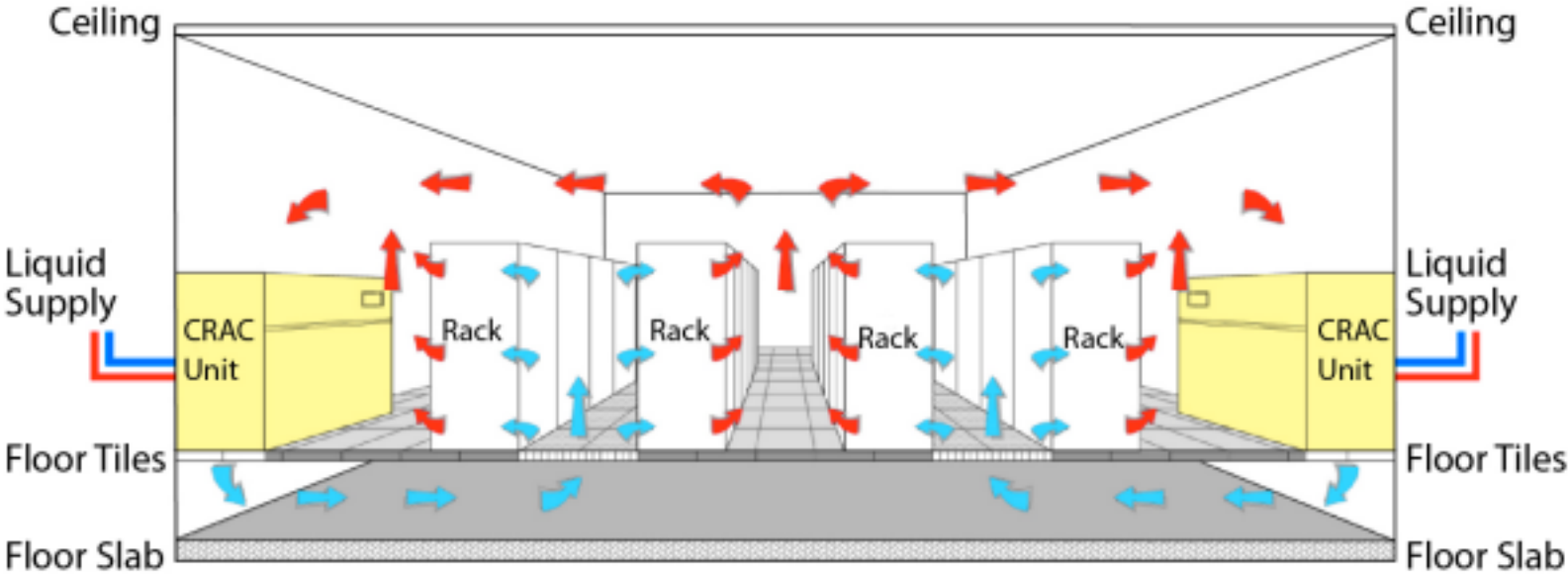
# Storage Hierarchy



# Anatomy of a Datacenter



# Anatomy of a Datacenter



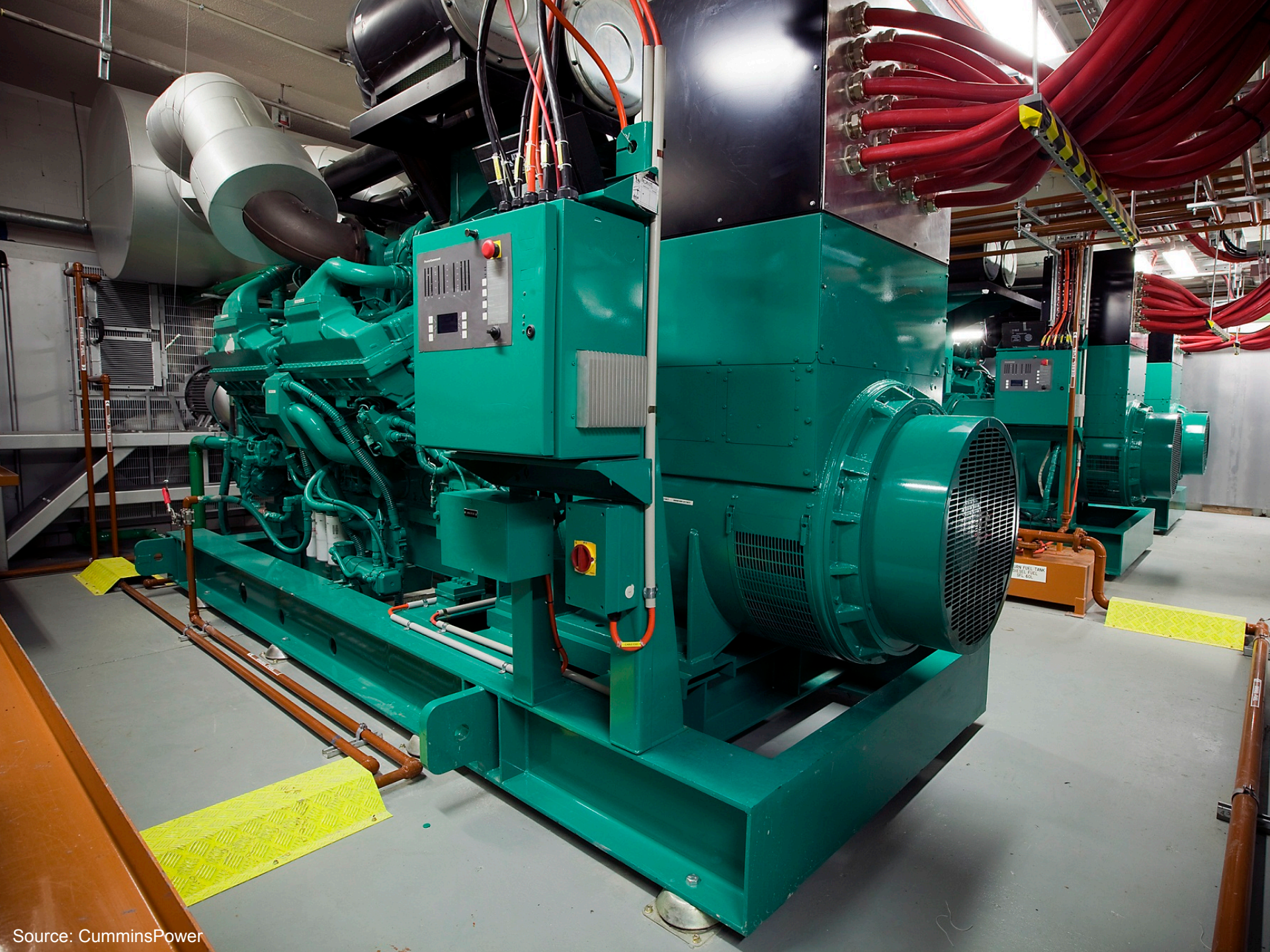
Source: Barroso and Urs Hölzle (2013)





Source: Google







An aerial photograph of an industrial facility, likely a power plant or refinery, during sunset. The sun is low on the horizon, casting a warm orange glow over the scene. The facility consists of several large, white, rectangular buildings with flat roofs, arranged in a grid-like pattern. In the foreground, there are several large, white, cylindrical storage tanks or silos. The surrounding area is a mix of green fields and brown, tilled soil. The sky is a gradient of orange and yellow, with the sun visible as a bright white circle on the left side.

**Aside: How much is 30 MW?**

# 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:  $\sim 100$  ns
  - Inter-node latencies:  $\sim 100$   $\mu$ s
- Let's model communication overhead...

# Modeling Communication Costs

- Simple execution cost model:

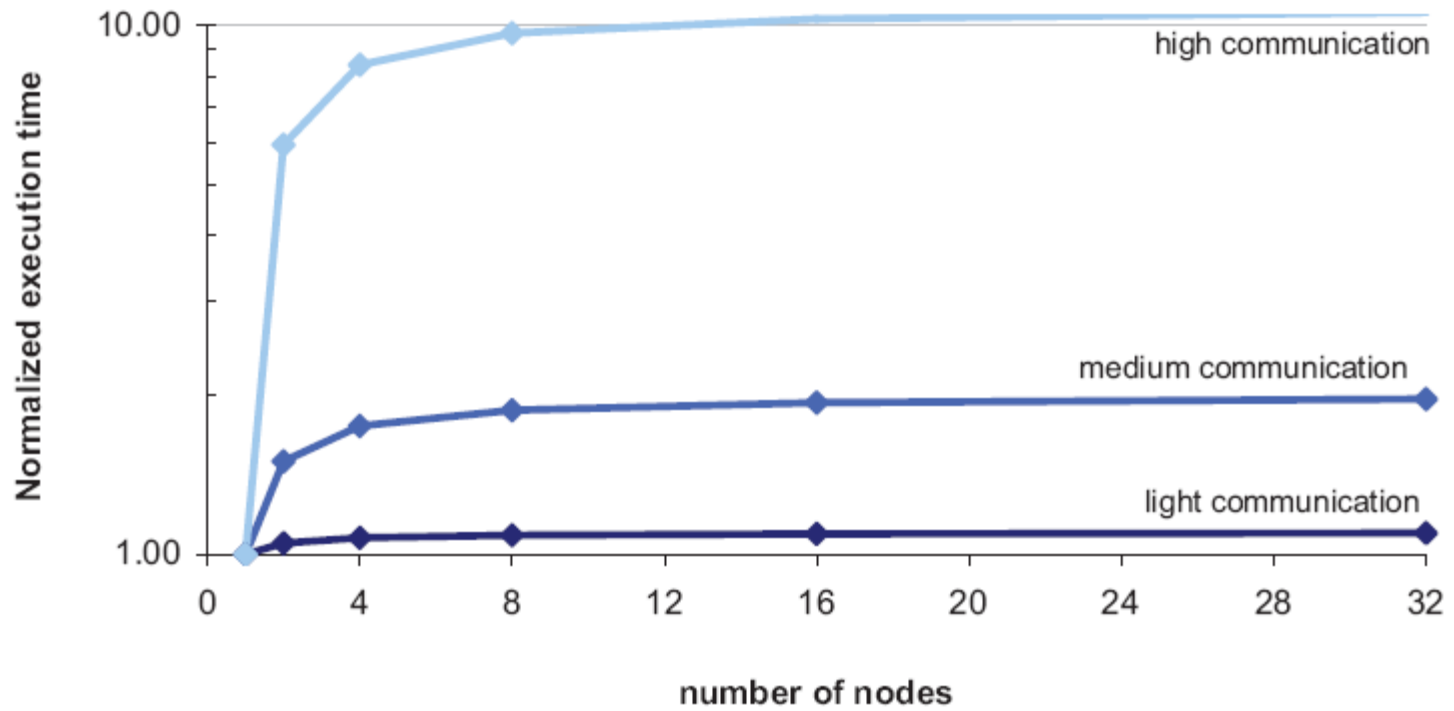
- Total cost = cost of computation + cost to access global data
- Fraction of local access inversely proportional to size of cluster
- $n$  nodes (ignore cores for now)

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

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

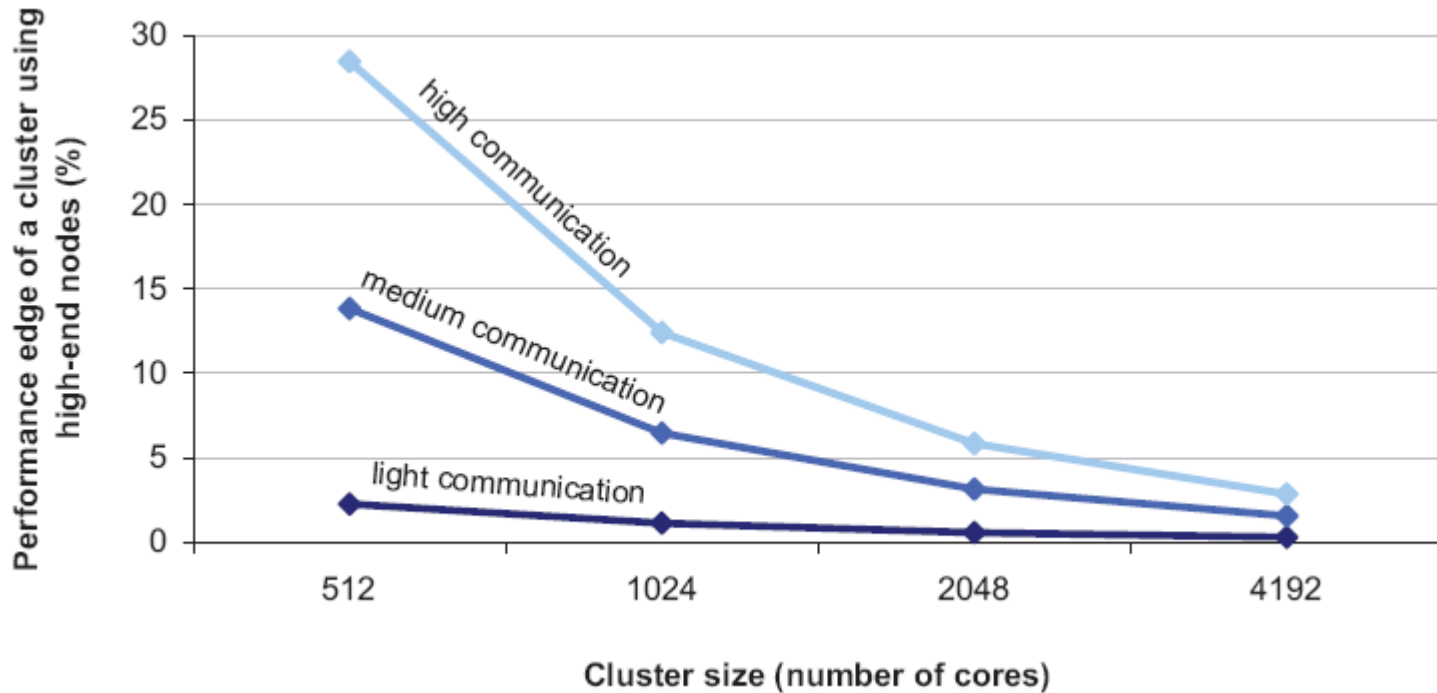
- What are the costs in parallelization?

# Cost of Parallelization



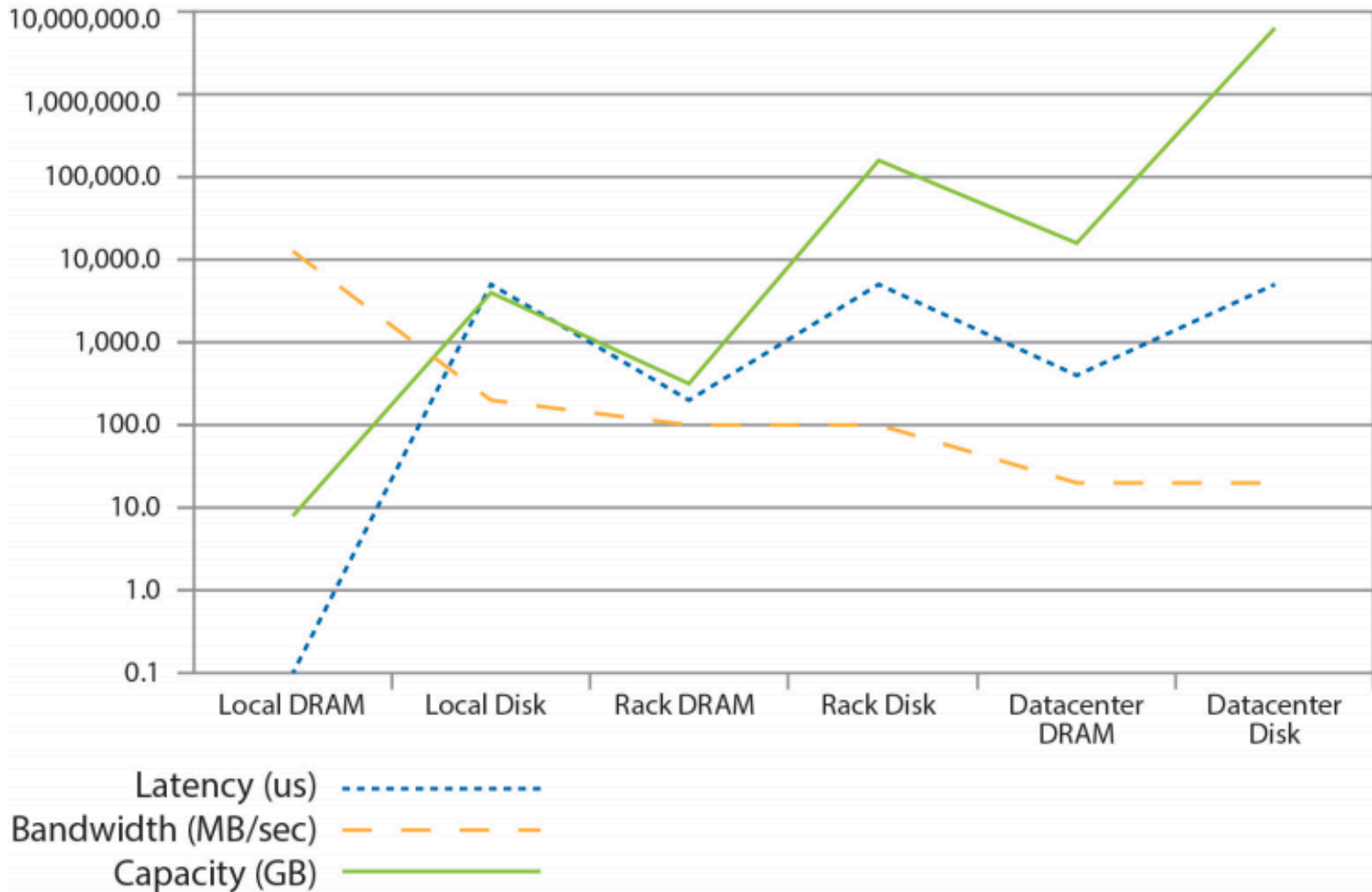


# Advantages of scaling “up”



So why not?  
Why does commodity beat exotic?

# 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

A wide-angle, high-angle photograph of a massive server room. The room is filled with rows of server racks, each with numerous lights glowing. A complex network of metal pipes and cables runs across the ceiling and down the sides of the racks. The lighting is predominantly blue, creating a cool, industrial atmosphere. The ceiling is a high, vaulted structure with a grid of steel beams. The floor is a light-colored, tiled surface. The overall impression is one of a large-scale, high-tech data center.

# Typical Big 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

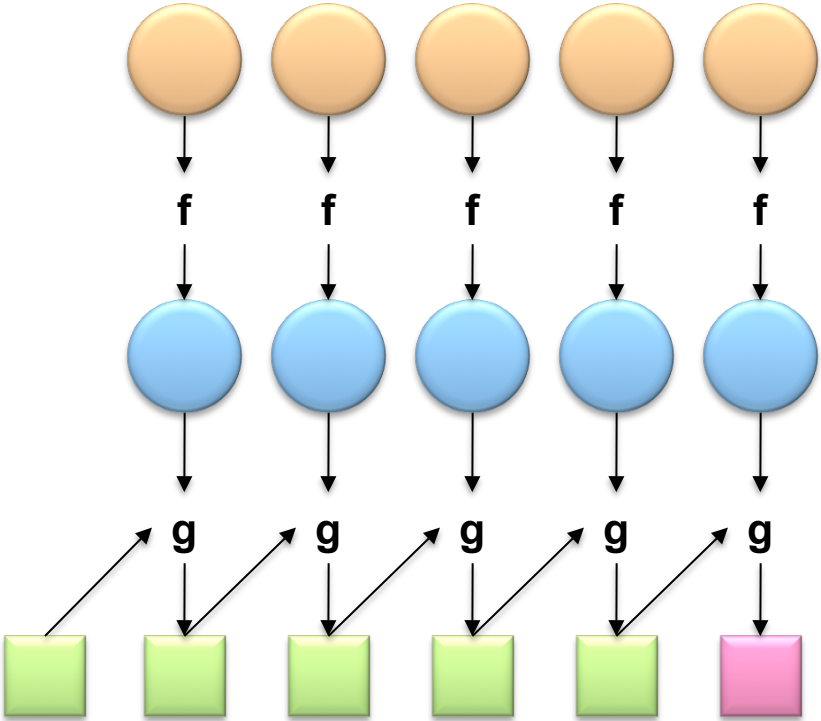
**Reduce**

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

# Roots in Functional Programming

**Map**

**Fold**



# MapReduce

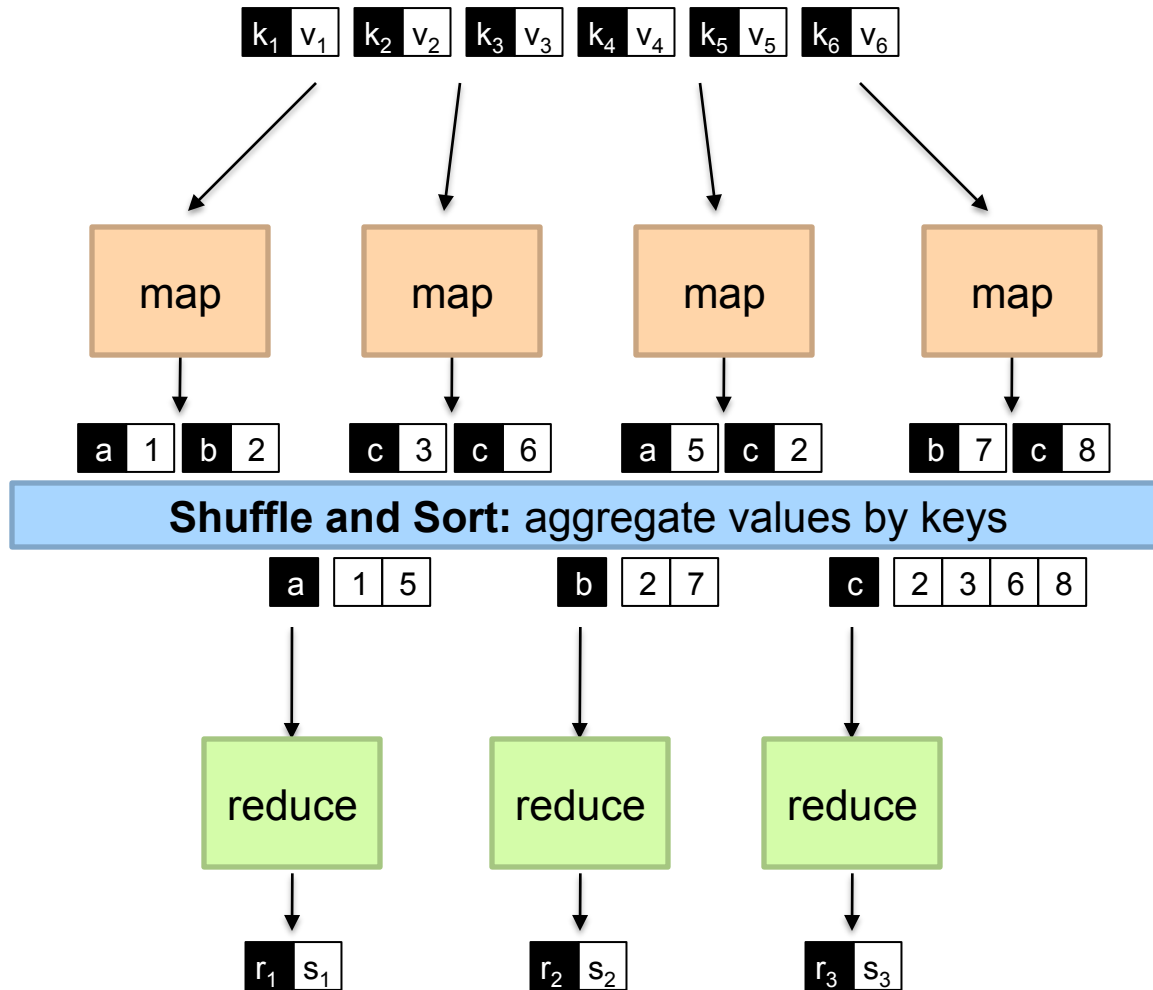
- Programmers specify two functions:

**map**  $(k_1, v_1) \rightarrow [ \langle k_2, v_2 \rangle ]$

**reduce**  $(k_2, [v_2]) \rightarrow [ \langle k_3, v_3 \rangle ]$

- 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...

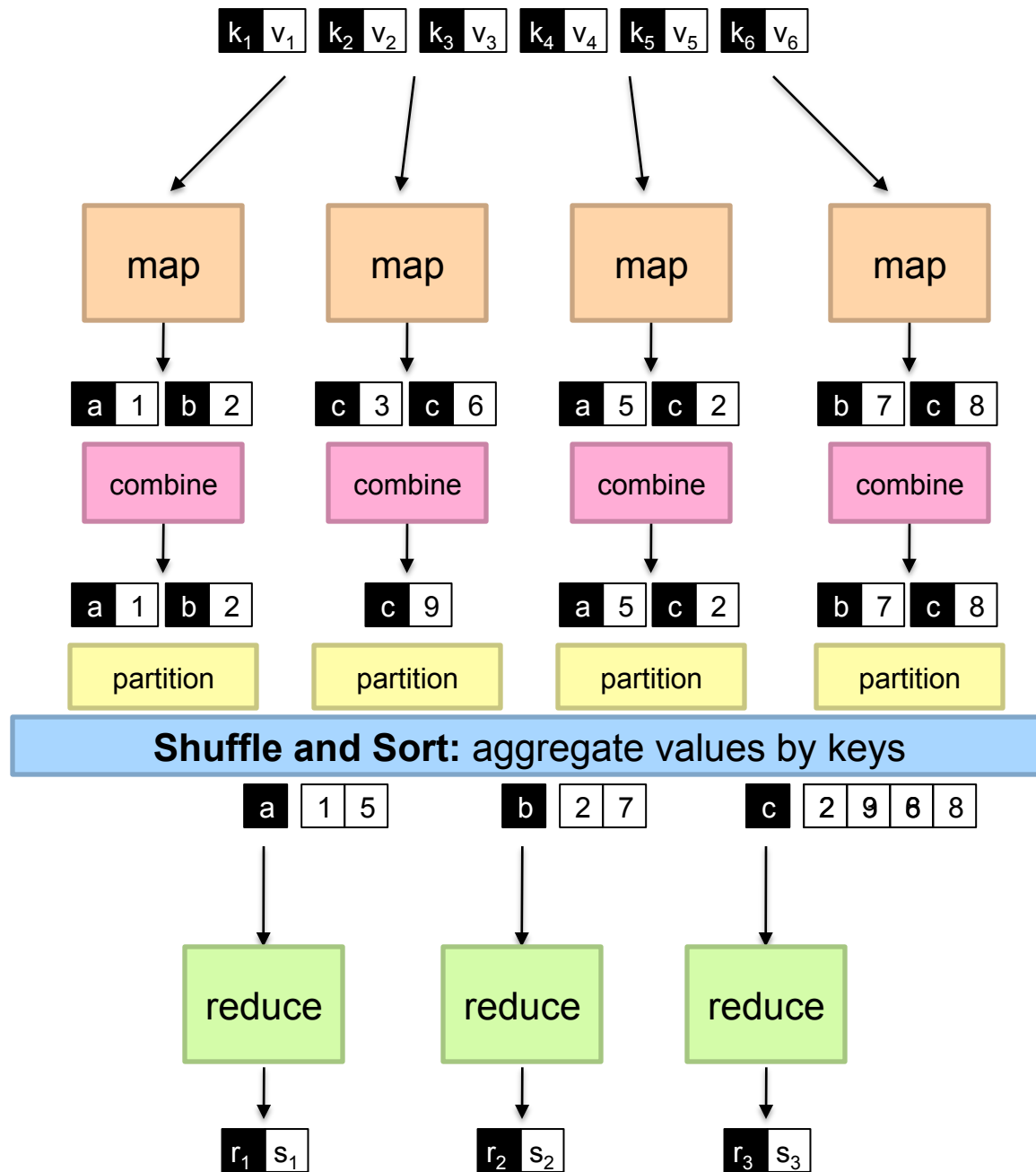
- Not quite...usually, programmers also specify:

**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

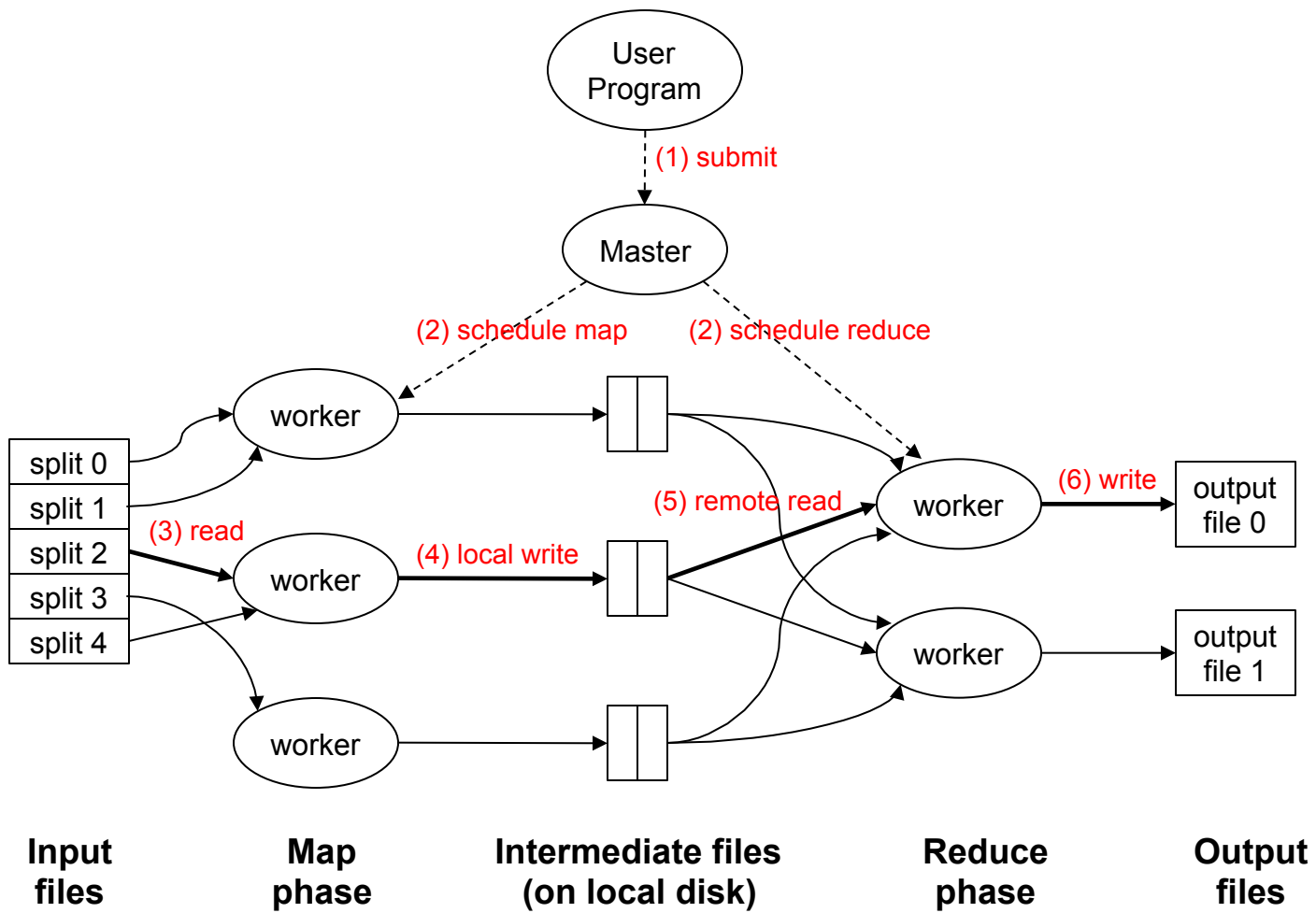
**Usage is usually clear from context!**



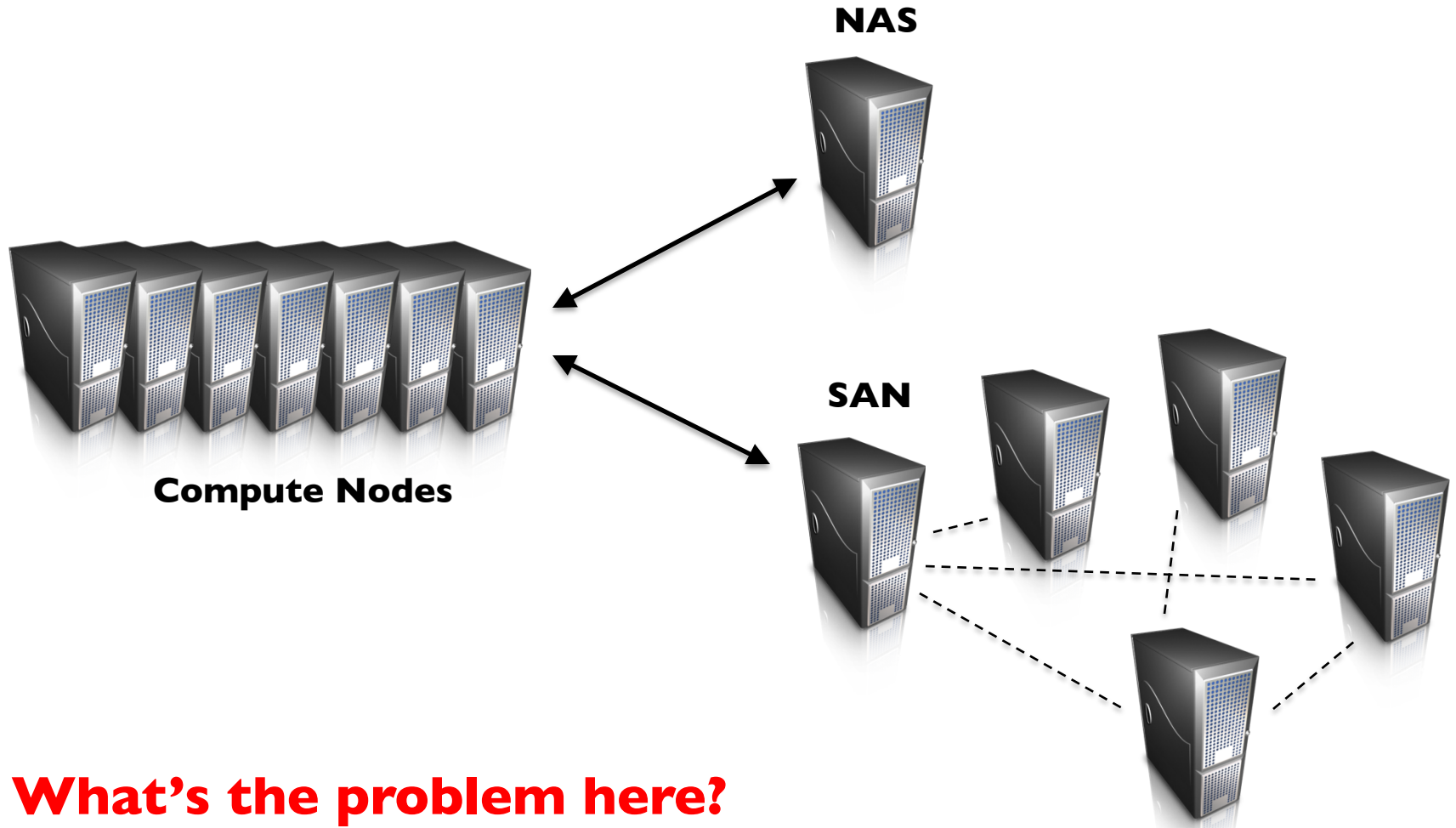
# MapReduce Implementations

- Google has a proprietary implementation in C++
  - Bindings in Java, Python
- Hadoop is an open-source implementation in Java
  - Development led by Yahoo, now an Apache project
  - Used in production at Yahoo, Facebook, Twitter, LinkedIn, Netflix, ...
  - The *de facto* big data processing platform
  - Large and expanding software ecosystem
- Lots of custom research implementations
  - For GPUs, cell processors, etc.





# How do we get data to the workers?



# 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)

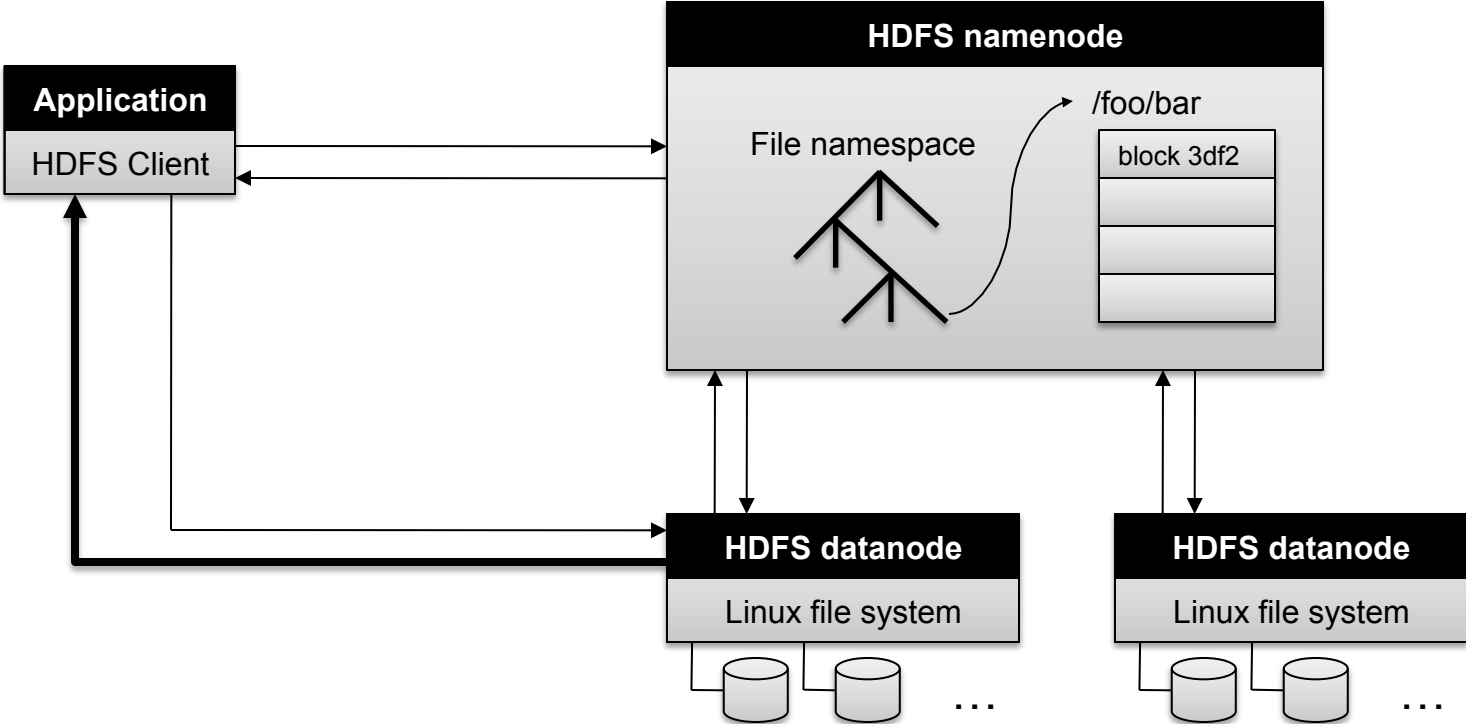
**HDFS = GFS clone (same basic ideas)**

# 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



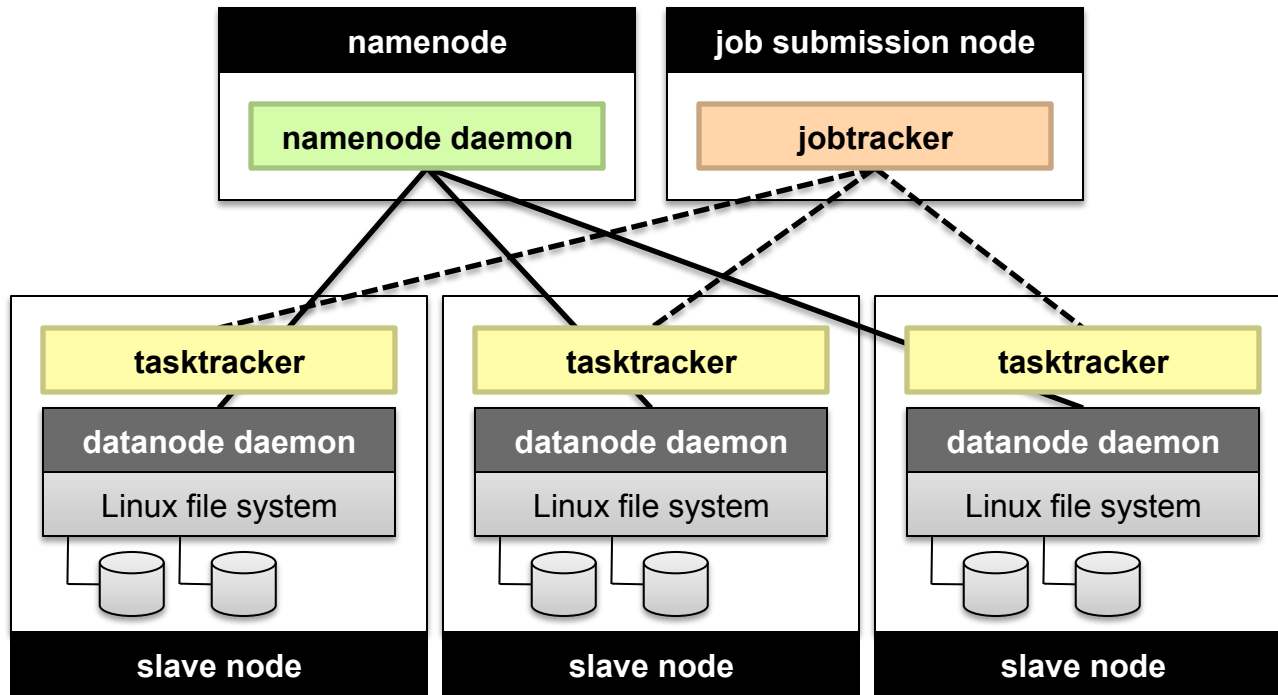
Adapted from (Ghemawat et al., SOSP 2003)



# 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

# Putting everything together...



(Not Quite... We'll come back to YARN later)



# Sequoia

16.32 PFLOPS

98,304 nodes with 1,572,864 million cores

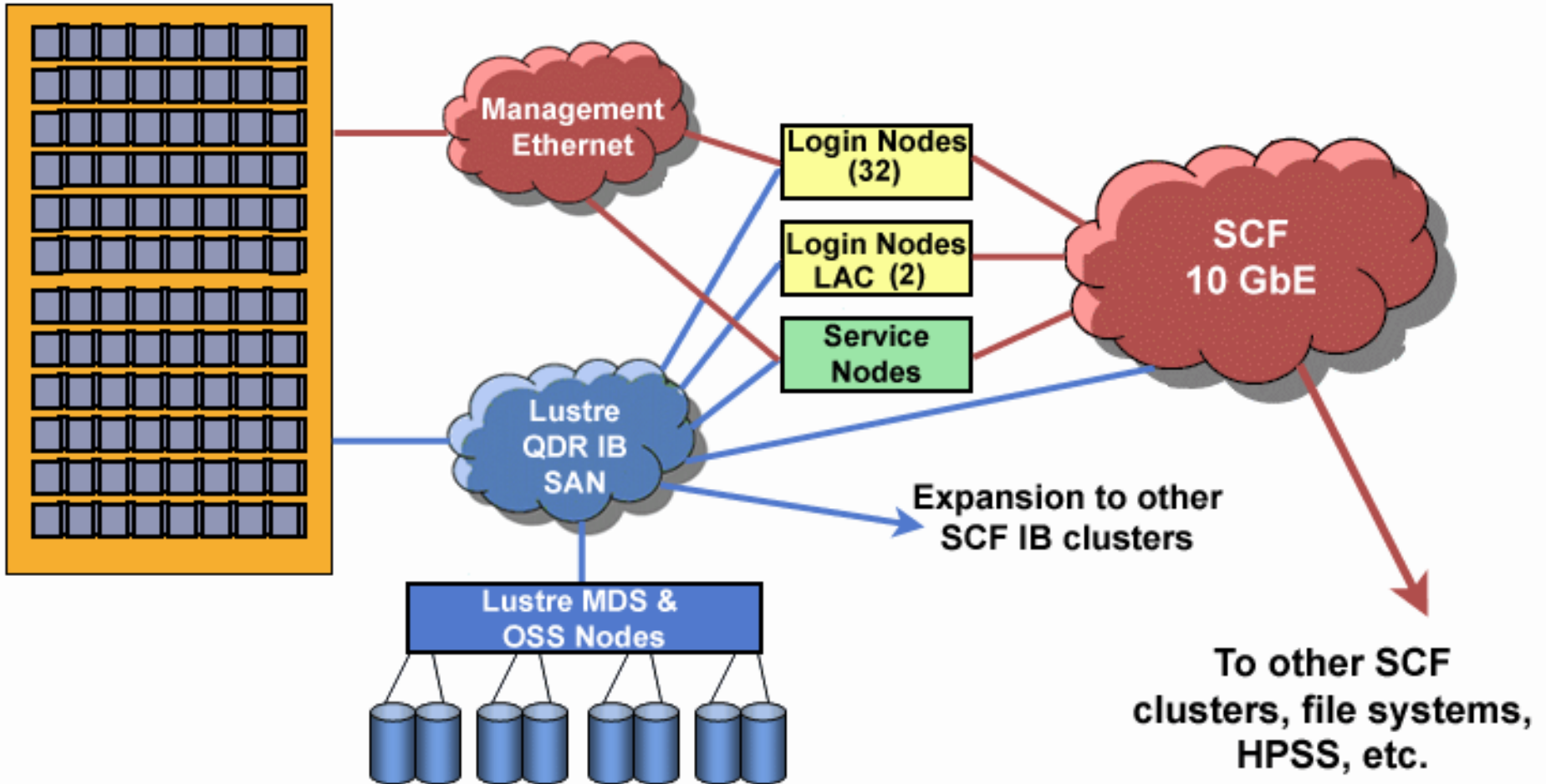
1.6 petabytes of memory

7.9 MWatts total power

# Sequoia

96 racks (12x8)  
98,304 compute nodes  
768 I/O nodes

- BG/Q 5D Torus Fabric
- QDR Infiniband
- Ethernet





# Questions?