

Big Data, Cloud Computing, and Map-Reduce

Session 1: Introduction to MapReduce

Based on Slides from:

<http://lintool.github.io/UMD-courses/bigdata-2015-Spring/>

and

<http://www.mmds.org/>



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

- We will use the book:

“Data-Intensive Text Processing with MapReduce”

by Jimmy Lin and Chris Dyer

(provided in Piazza and on-line)

- We will also concentrate on Hadoop version of MapReduce

Cloud Computing and Big data

- Computing on “big data”
- Focus on applications and algorithm design
- MapReduce... and beyond



processes 20 PB a day (2008)
crawls 20B web pages a day (2012)



>10 PB data, 75B DB calls
per day (6/2012)

>100 PB of user data +
500 TB/day (8/2012)



S3: 449B objects, peak 290k
request/second (7/2011)
1T objects (6/2012)



640K ought to be
enough for
anybody.

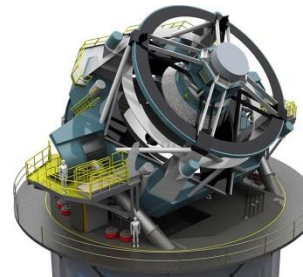


150 PB on 50k+ servers
running 15k apps (6/2011)



Wayback Machine: 240B web
pages archived, 5 PB (1/2013)

LHC: ~15 PB a year



LSST: 6-10 PB a year
(~2015)

SKA: 0.3 – 1.5 EB
per year (~2020)



How much data?

Cloud Computing

- 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

Utility Computing

○ What?

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

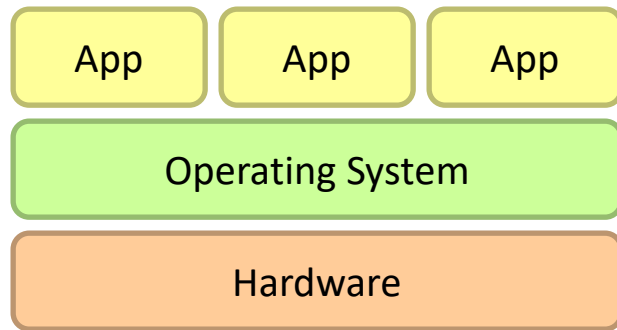
○ Why?

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

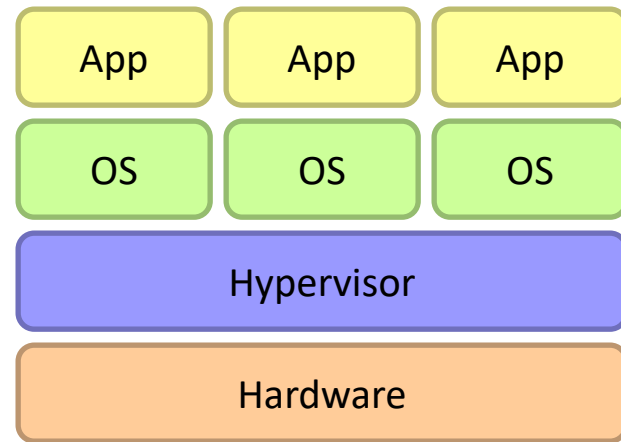
○ Does it make sense?

- Benefits to cloud users
- Business case for cloud providers

Enabling Technology: Virtualization



Traditional Stack

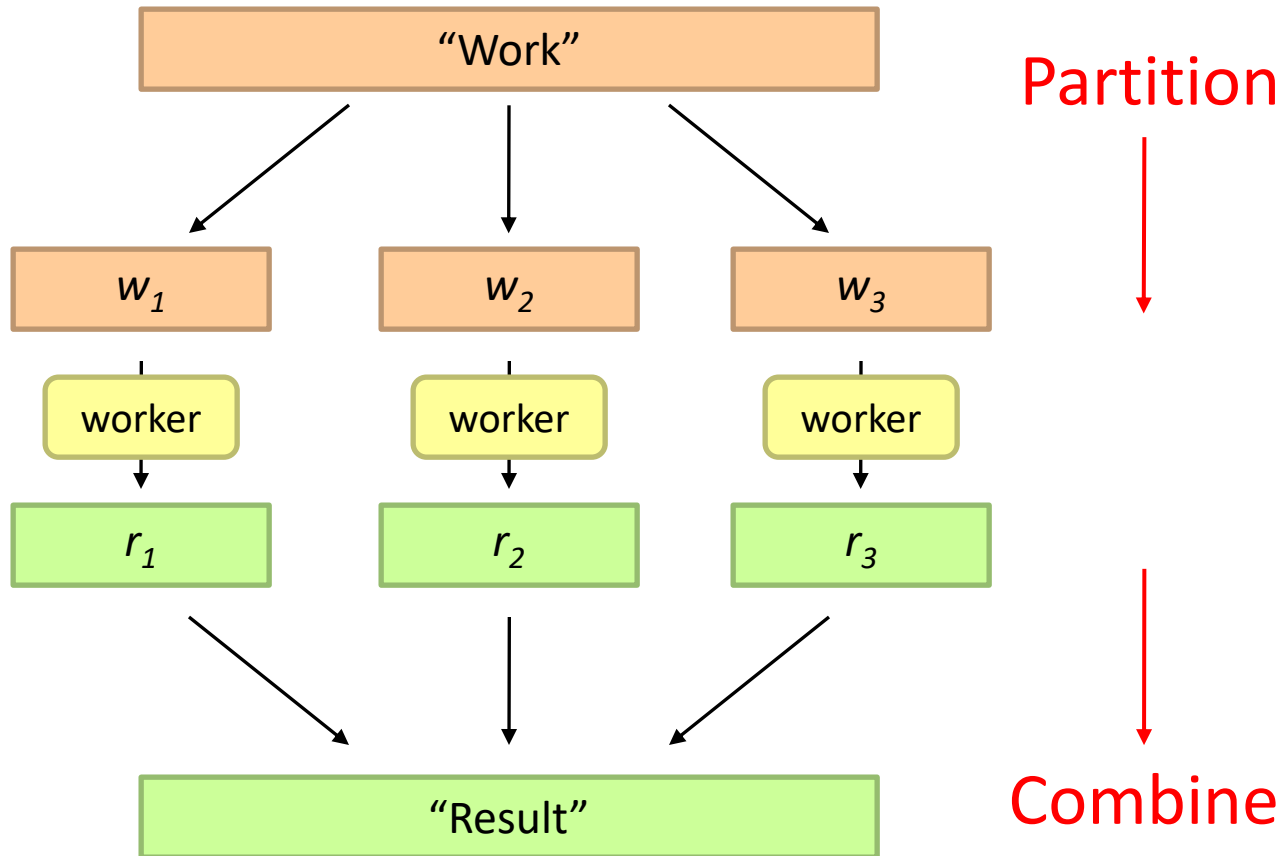


Virtualized Stack

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

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

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!

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



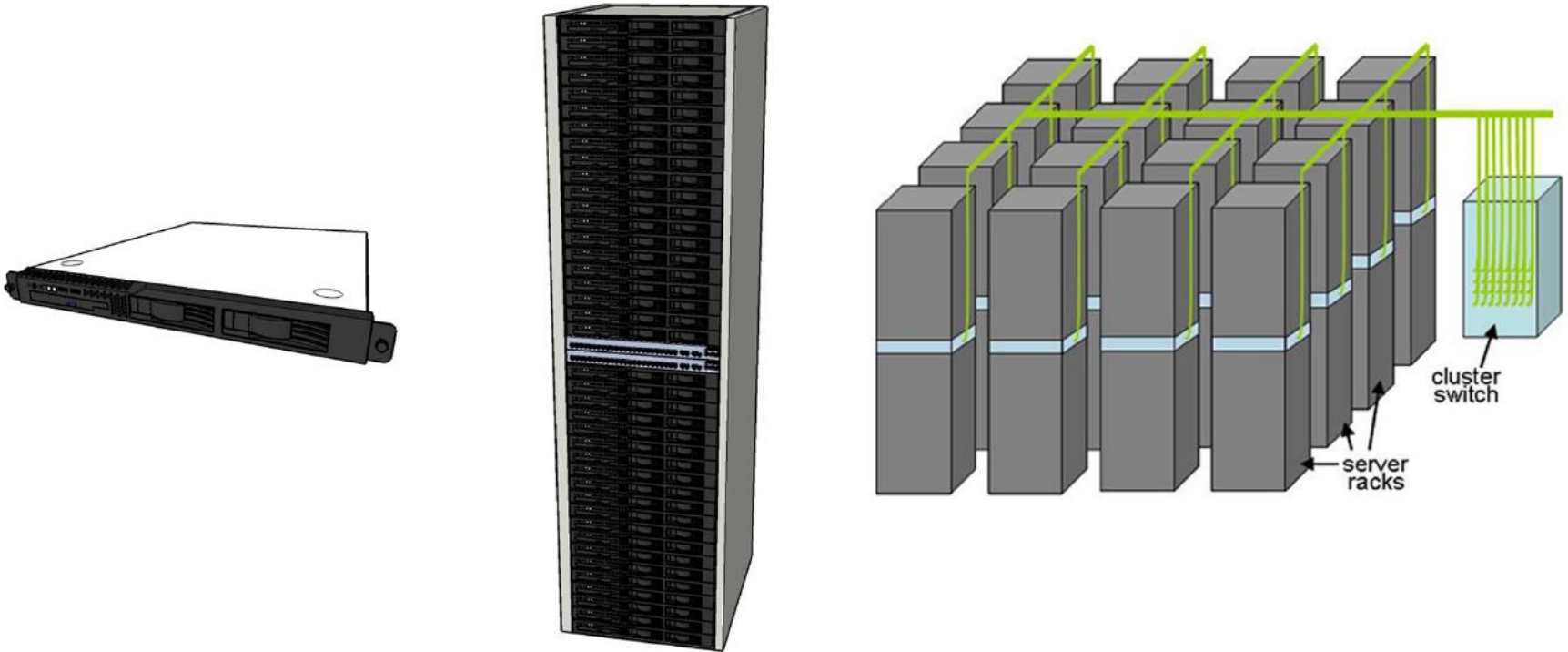
The datacenter *is* the computer!

What's the point?

- It's all about the right level of abstraction
 - Moving beyond the von Neumann architecture
 - We need better programming models
- 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

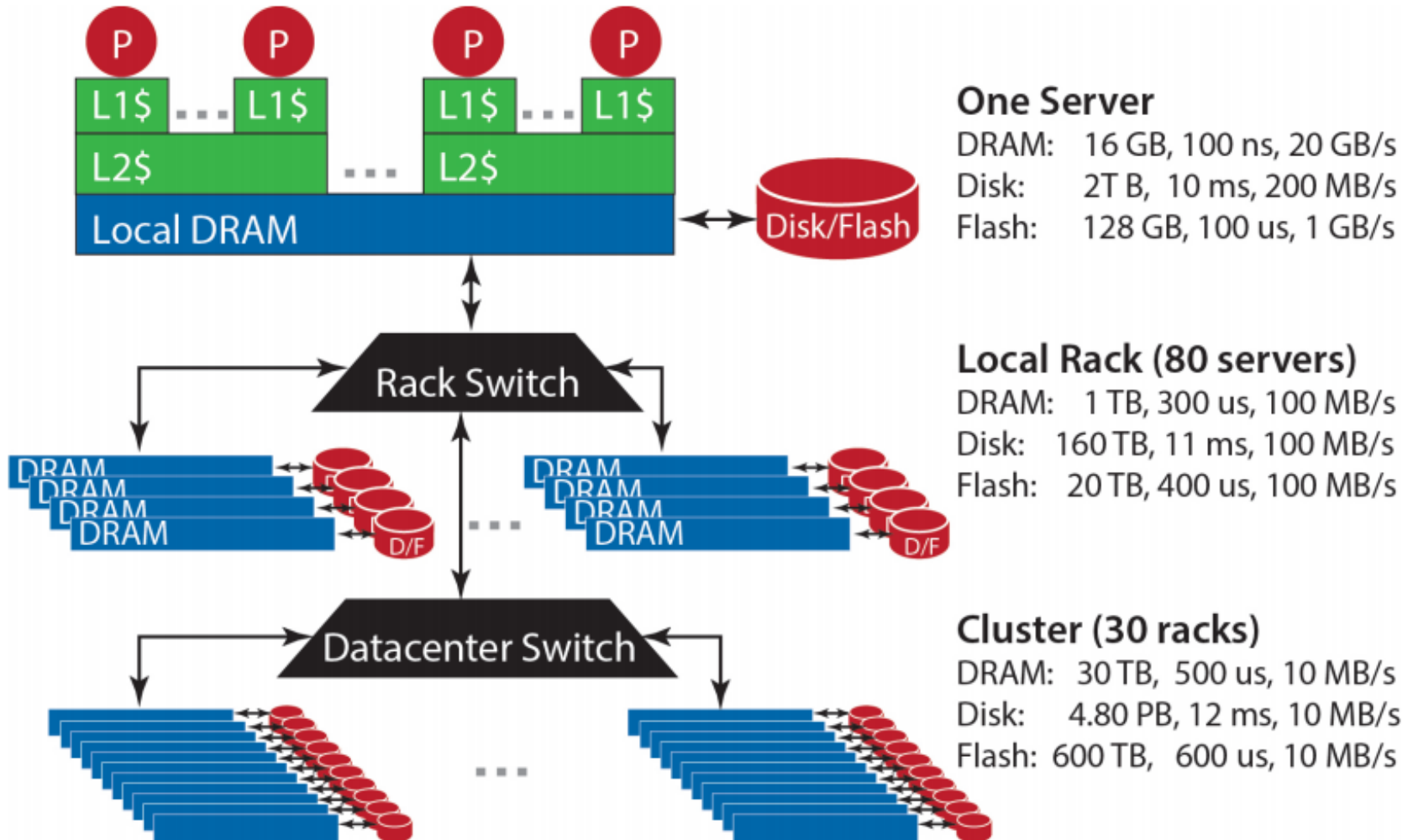
The datacenter *is* the computer!

Building Blocks

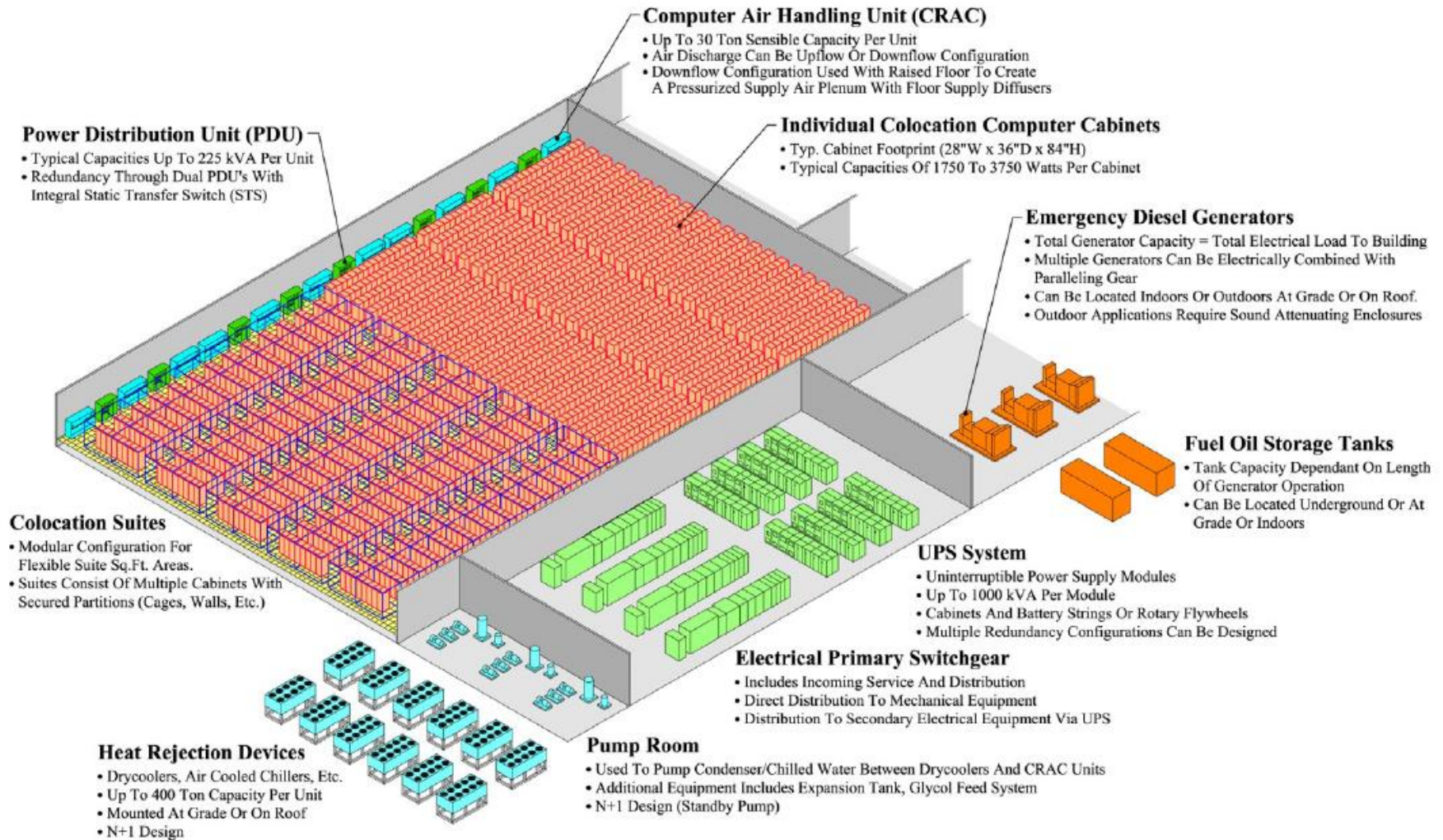




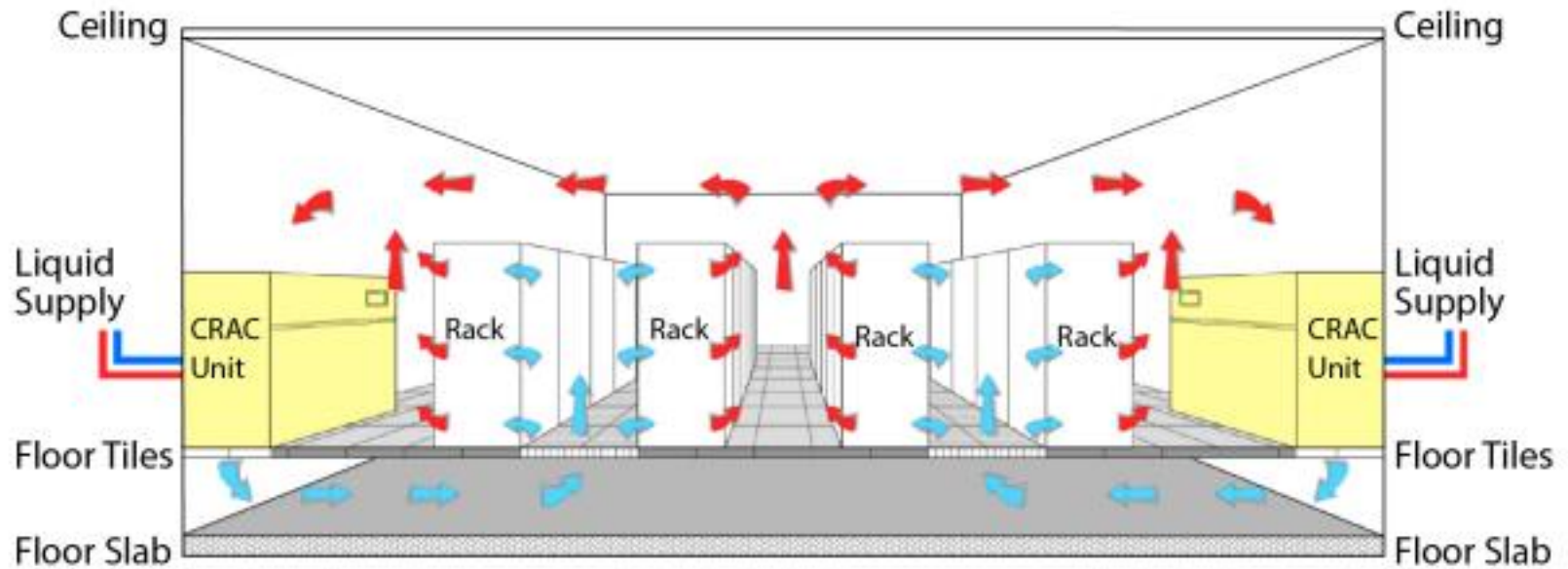
Storage Hierarchy



Anatomy of a Datacenter



Anatomy of a Datacenter



“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, low-perspective shot of a massive data center. The floor is a light-colored square tile. Rows of server racks stretch into the distance, illuminated by a mix of blue and yellow light. A complex network of black overhead cables and metal support structures crisscrosses the upper half of the frame, creating a dense, industrial-looking ceiling. The overall atmosphere is one of high-tech infrastructure.

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 **Reduce**
- Generate final output

Key idea: provide a functional abstraction for these two operations

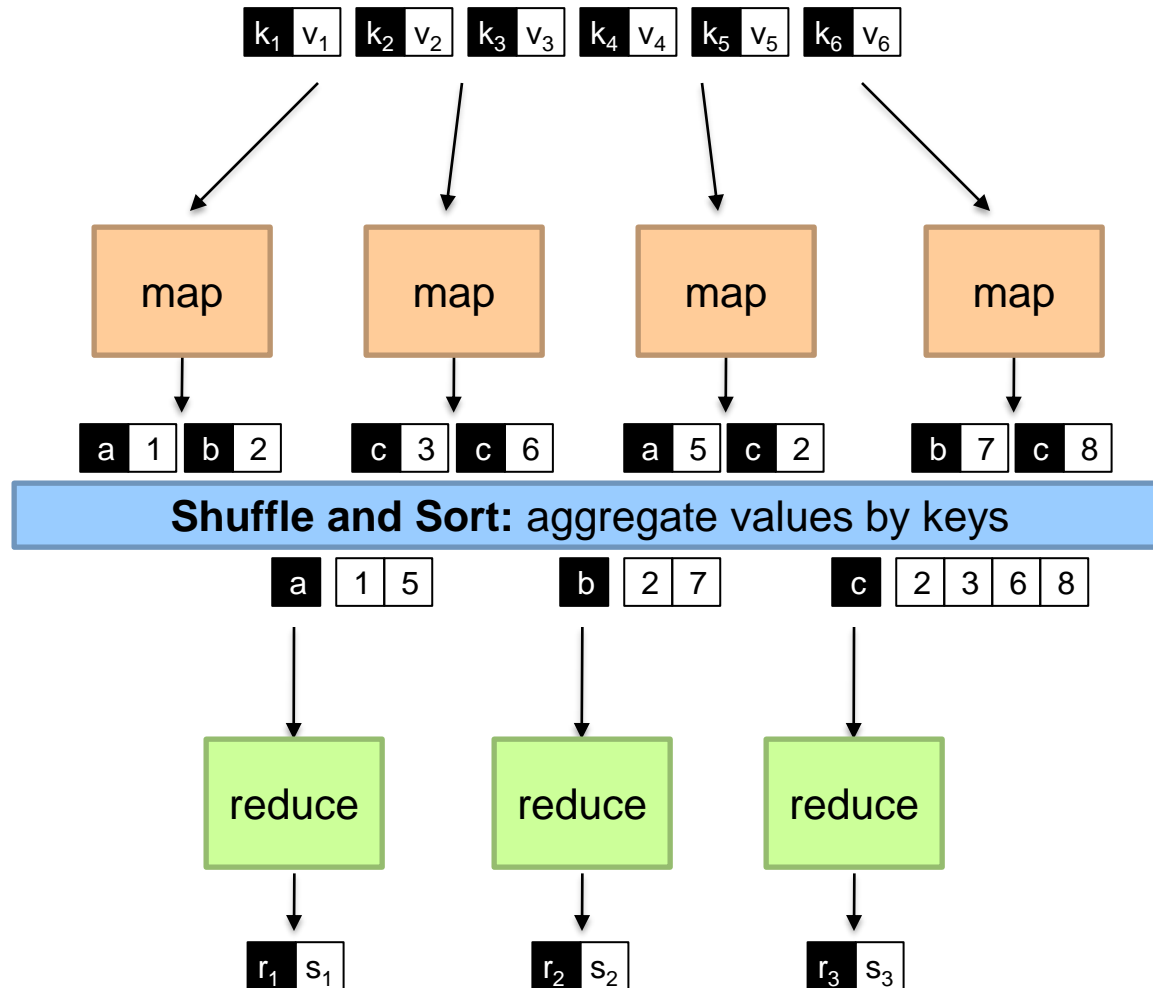
MapReduce

- 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 $(k', v') \rightarrow \langle k', v' \rangle^*$

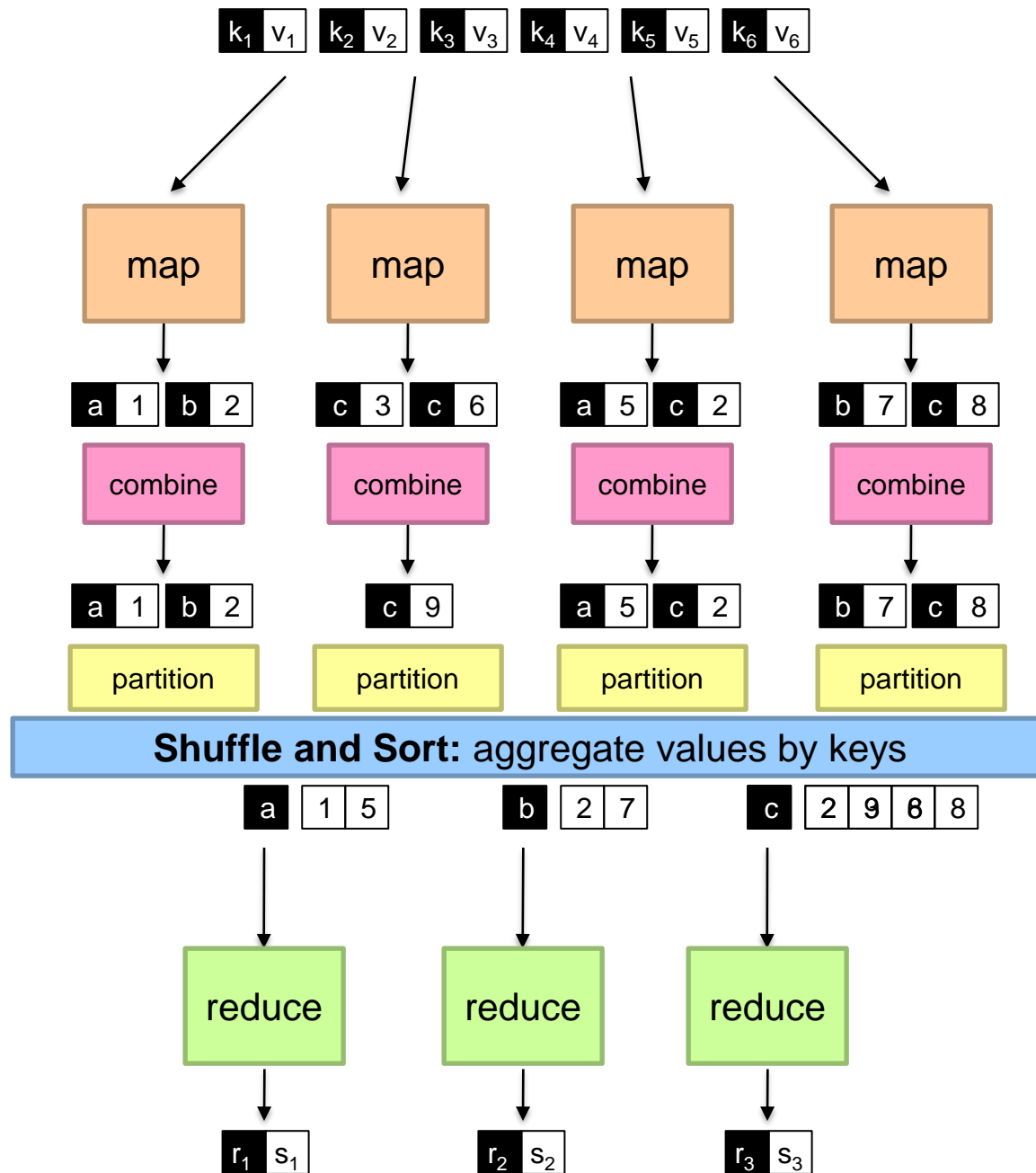
- All values with the same key are reduced together
- The execution framework handles everything else...
- Not quite...usually, programmers also specify:

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...

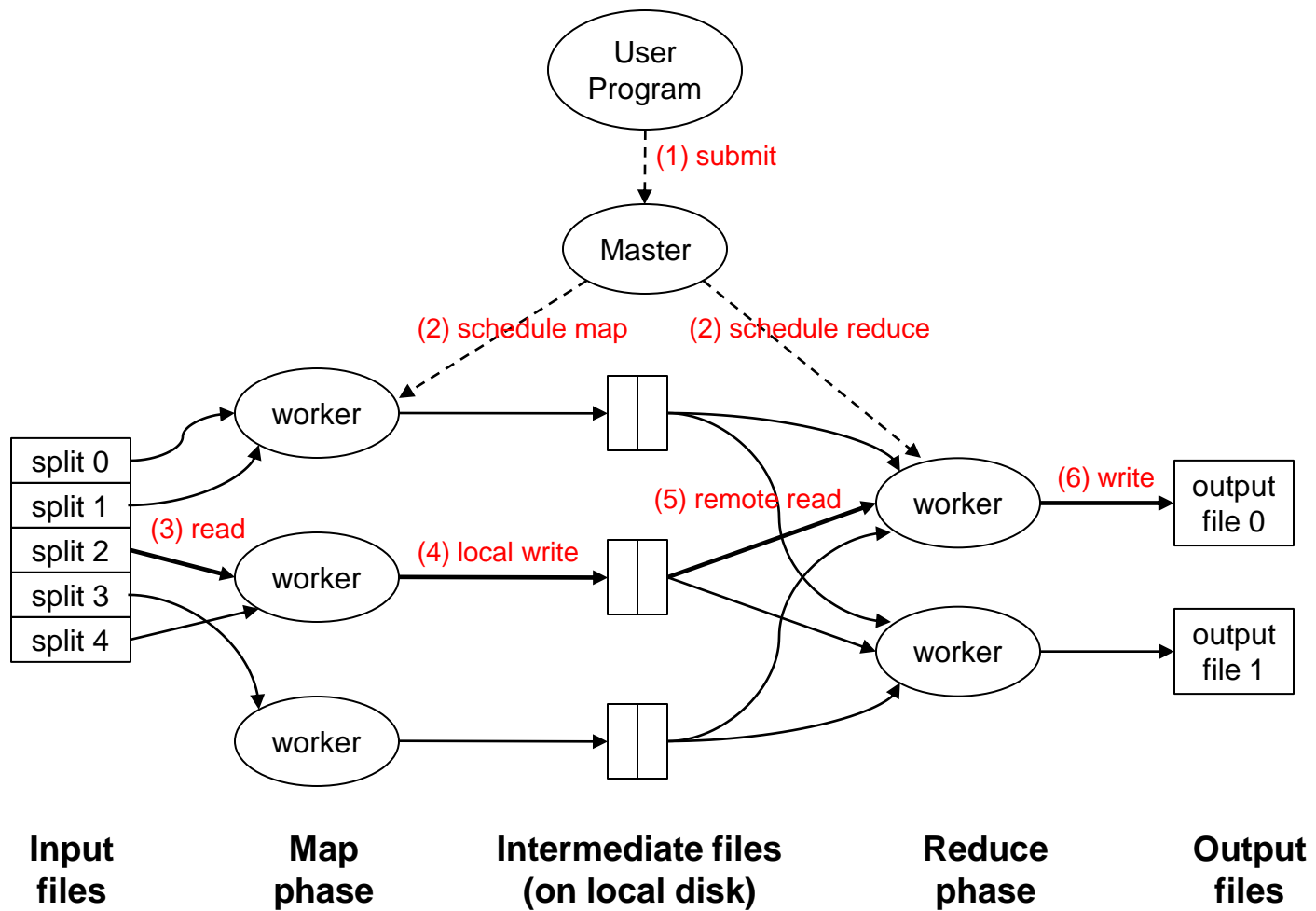
- The programming model
- The execution framework (aka “runtime”)
- The specific implementation

Usage is usually clear from context!

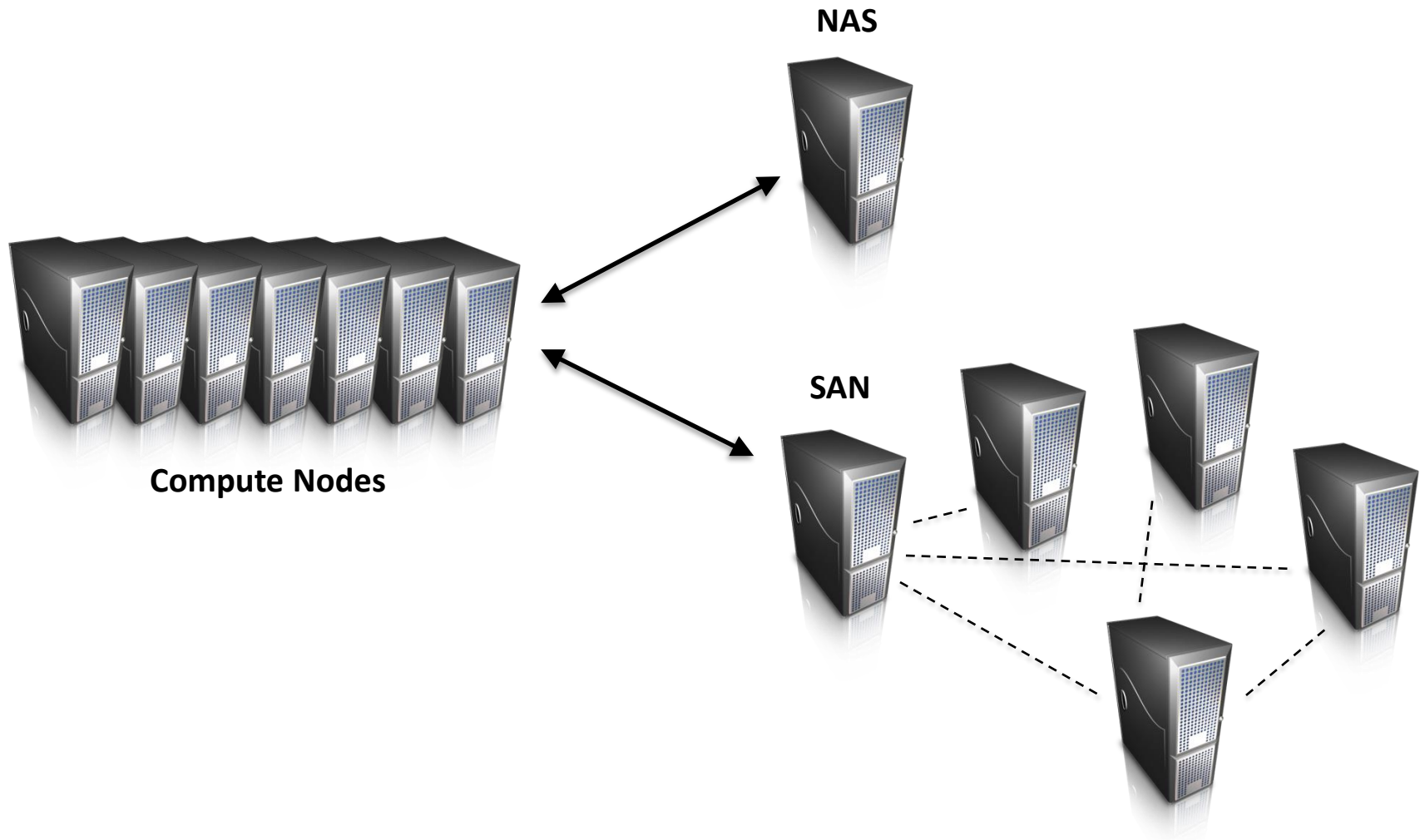
MapReduce Implementations

- 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
 - Rapidly 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?
 - 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)

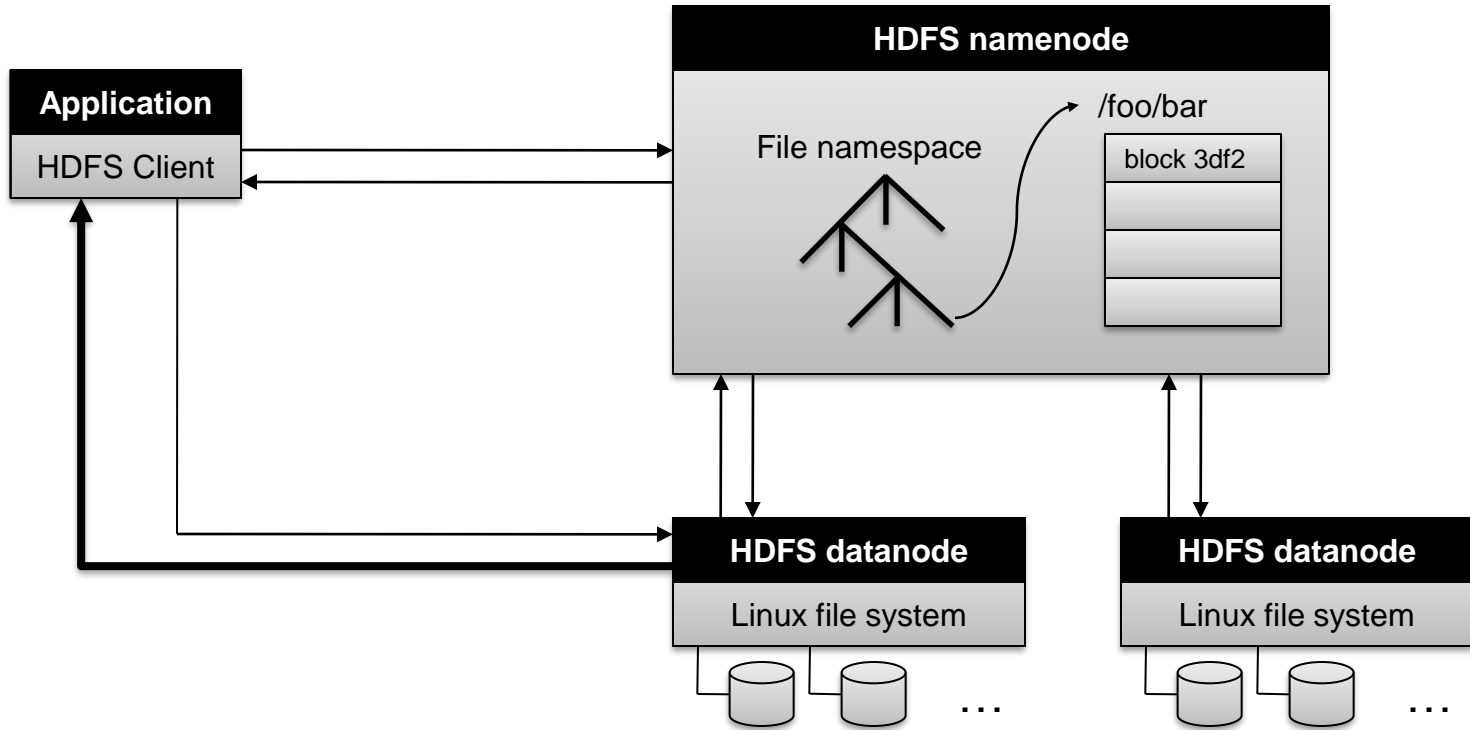
HDFS = GFS clone (same basic ideas)

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



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...

