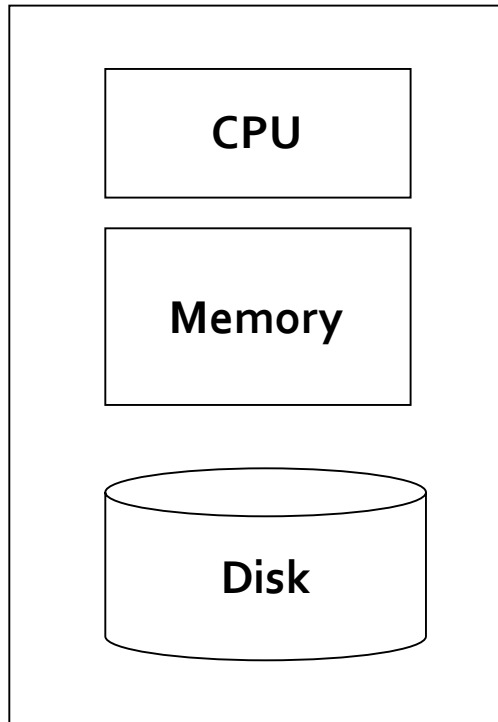


# Map Reduce

David Wemhoener

Acknowledgement: Majority of the slides are taken from Mining of Massive Datasets  
Jure Leskovec, Anand Rajaraman, Jeff Ullman

# Single Node Architecture



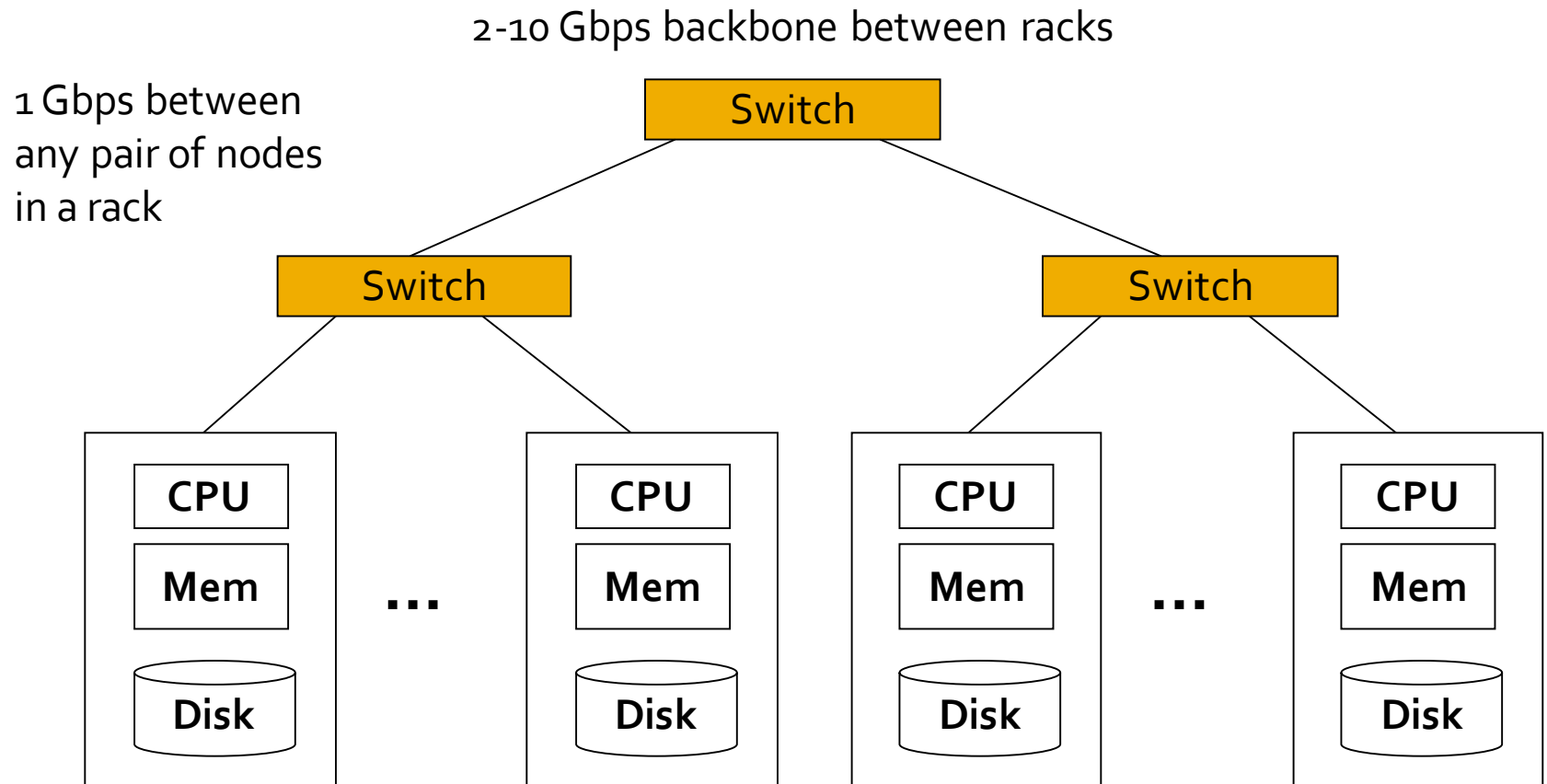
**Machine Learning, Statistics**

**“Classical” Data Mining**

# Motivation: Google Example

- 20+ billion web pages x 20KB = 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!**
- **Today, a standard architecture for such problems is emerging:**
  - Cluster of commodity Linux nodes
  - Commodity network (ethernet) to connect them

# Cluster Architecture



Each rack contains 16-64 nodes

In 2011 it was guestimated that Google had 1M machines, <http://bit.ly/Shh0RO>



# Large-scale Computing

- **Large-scale computing for data mining problems on commodity hardware**
- **Challenges:**
  - **How do you distribute computation?**
  - **How can we make it easy to write distributed programs?**
  - **Machines fail:**
    - One server may stay up 3 years (1,000 days)
    - If you have 1,000 servers, expect to loose 1/day
    - People estimated Google had ~1M machines in 2011
      - 1,000 machines fail every day!

# Idea and Solution

- **Issue:** Copying data over a network takes time
- **Idea:**
  - Bring computation close to the data
  - Store files multiple times for reliability
- **Map-reduce** addresses these problems
  - Google's computational/data manipulation model
  - Elegant way to work with big data
  - **Storage Infrastructure – File system**
    - Google: GFS. Hadoop: HDFS
  - **Programming model**
    - Map-Reduce

# Storage Infrastructure

- **Problem:**

- If nodes fail, how to store data persistently?

- **Answer:**

- **Distributed File System:**

- Provides global file namespace
    - Google GFS; Hadoop HDFS;

- **Typical usage pattern**

- Huge files (100s of GB to TB)
  - Data is rarely updated in place
  - Reads and appends are common



# Distributed File System

## ■ Chunk servers

- File is split into contiguous chunks
- Typically each chunk is 16-64MB
- Each chunk replicated (usually 2x or 3x)
- Try to keep replicas in different racks

## ■ Master node

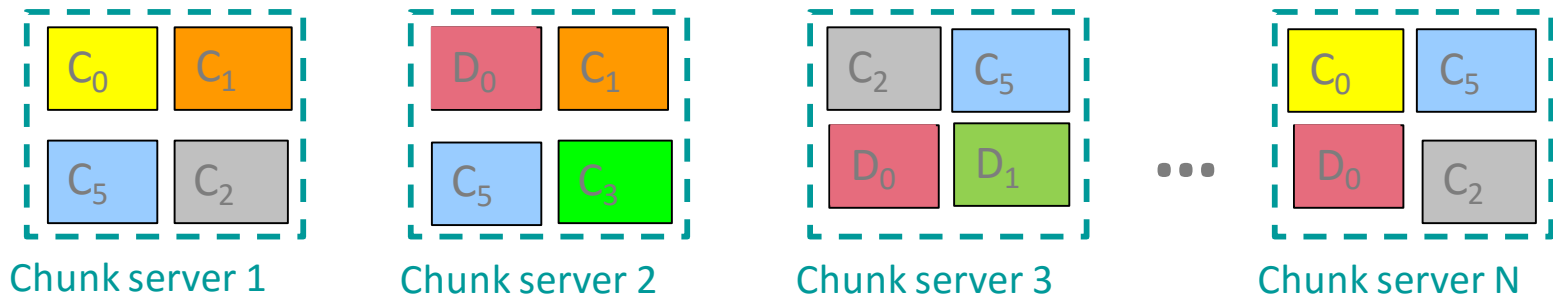
- a.k.a. Name Node in Hadoop's HDFS
- Stores metadata about where files are stored
- Might be replicated

## ■ Client library for file access

- Talks to master to find chunk servers
- Connects directly to chunk servers to access data

# Distributed File System

- **Reliable distributed file system**
- Data kept in “chunks” spread across machines
- Each chunk **replicated** on different machines
  - Seamless recovery from disk or machine failure



Bring computation directly to the data!

Chunk servers also serve as compute servers

# Map-Reduce: Environment

## Map-Reduce environment takes care of:

- Partitioning the input data
- Scheduling the program's execution across a set of machines
- Performing the **group by key** step
- Handling machine failures
- Managing required inter-machine communication

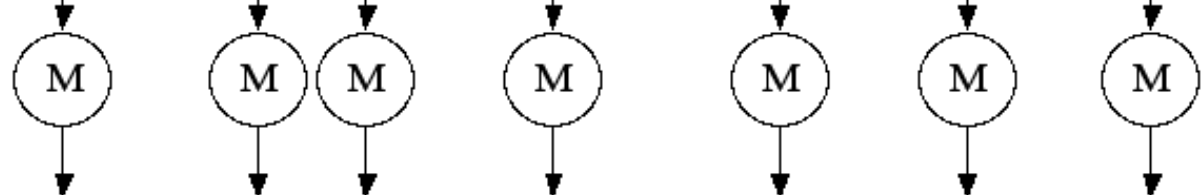
# Map-Reduce: A diagram

Input

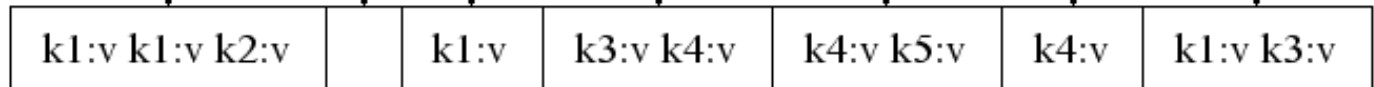


## MAP:

Read input and produces a set of key-value pairs

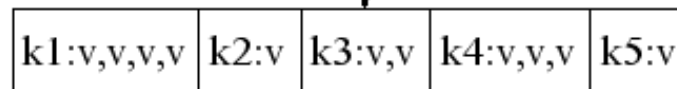


Intermediate

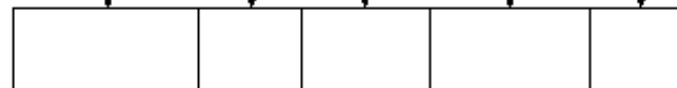


Group by Key

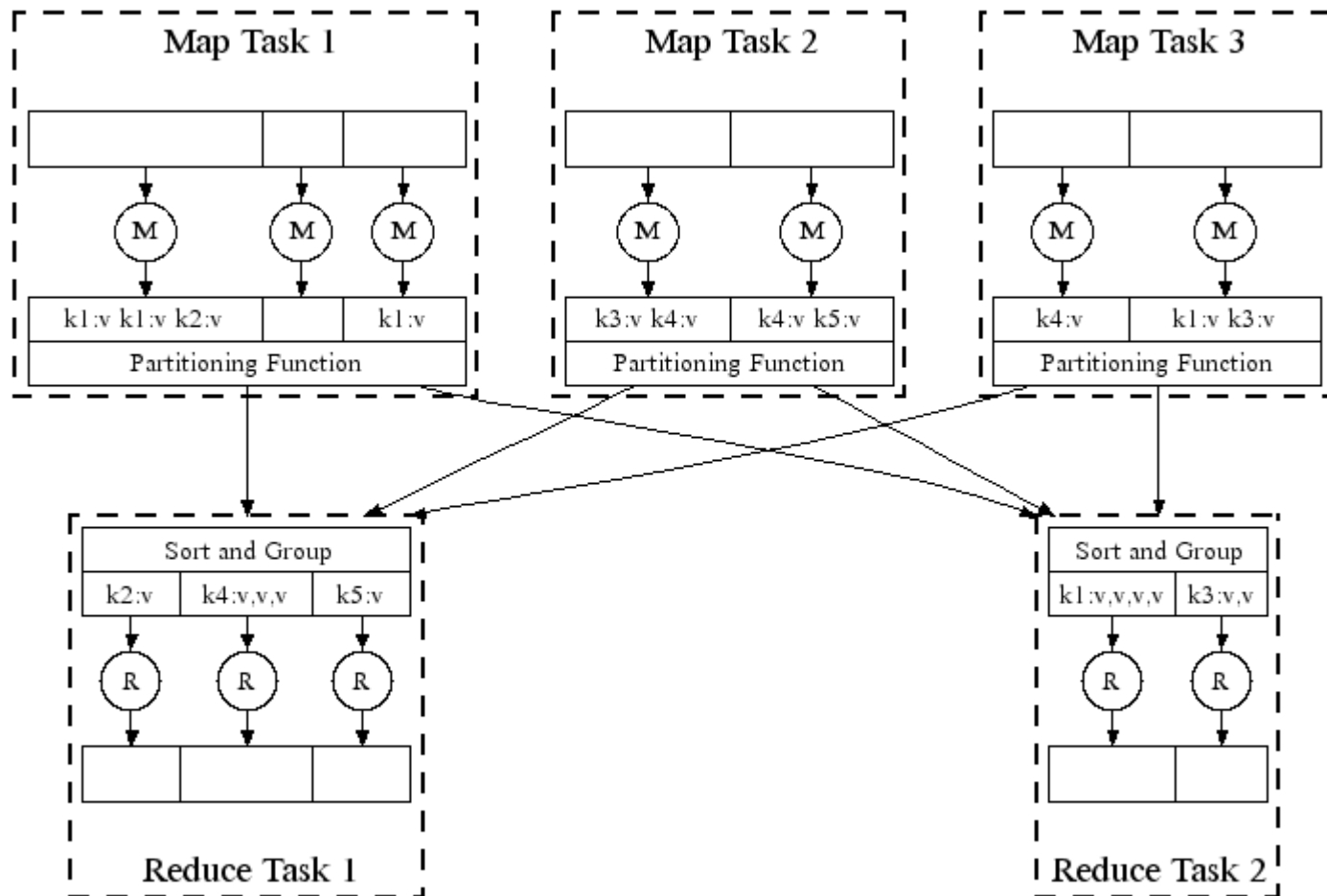
Grouped



Output



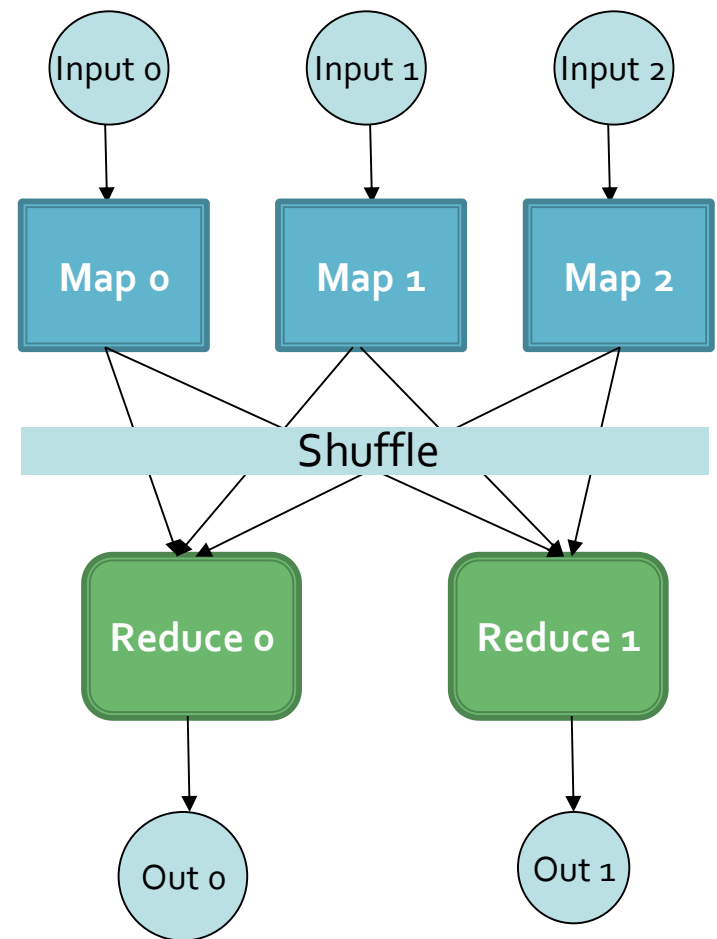
# Map-Reduce: In Parallel



All phases are distributed with many tasks doing the work

# Map-Reduce

- **Programmer specifies:**
  - Map and Reduce and input files
- **Workflow:**
  - Read inputs as a set of key-value-pairs
  - **Map** transforms input kv-pairs into a new set of k'v'-pairs
  - Sorts & Shuffles the k'v'-pairs to output nodes
  - All k'v'-pairs with a given k' are sent to the same **reduce**
  - **Reduce** processes all k'v'-pairs grouped by key into new k''v''-pairs
  - Write the resulting pairs to files
- **All phases are distributed with many tasks doing the work**



# Data Flow

- **Input and final output are stored on a distributed file system (FS):**
  - Scheduler tries to schedule map tasks “close” to physical storage location of input data
- **Intermediate results are stored on local FS of Map and Reduce workers**
- **Output is often input to another MapReduce task**

# Coordination: Master

- **Master node takes care of coordination:**
  - **Task status:** (idle, in-progress, completed)
  - **Idle tasks** get scheduled as workers become available
  - When a map task completes, it sends the master the location and sizes of its  $R$  intermediate files, one for each reducer
  - Master pushes this info to reducers
- Master pings workers periodically to detect failures



# Dealing with Failures

## ■ Map worker failure

- Map tasks completed or in-progress at worker are reset to idle
- Reduce workers are notified when task is rescheduled on another worker

## ■ Reduce worker failure

- Only in-progress tasks are reset to idle
- Reduce task is restarted

## ■ Master failure

- MapReduce task is aborted and client is notified

# How many Map and Reduce jobs?

- $M$  map tasks,  $R$  reduce tasks
- **Rule of a thumb:**
  - Make  $M$  much larger than the number of nodes in the cluster
  - One DFS chunk per map is common
  - Improves dynamic load balancing and speeds up recovery from worker failures
- **Usually  $R$  is smaller than  $M$** 
  - Because output is spread across  $R$  files

# Refinements: Backup Tasks

## ■ Problem

- Slow workers significantly lengthen the job completion time:
  - Other jobs on the machine
  - Bad disks
  - Weird things

## ■ Solution

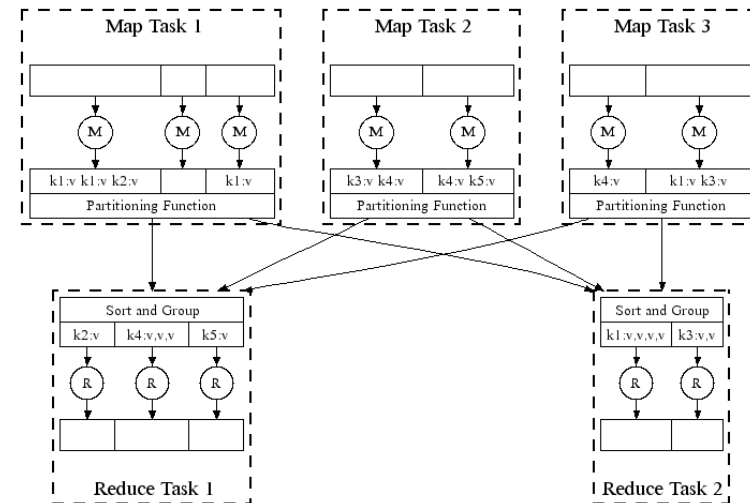
- Near end of phase, spawn backup copies of tasks
  - Whichever one finishes first “wins”

## ■ Effect

- Dramatically shortens job completion time

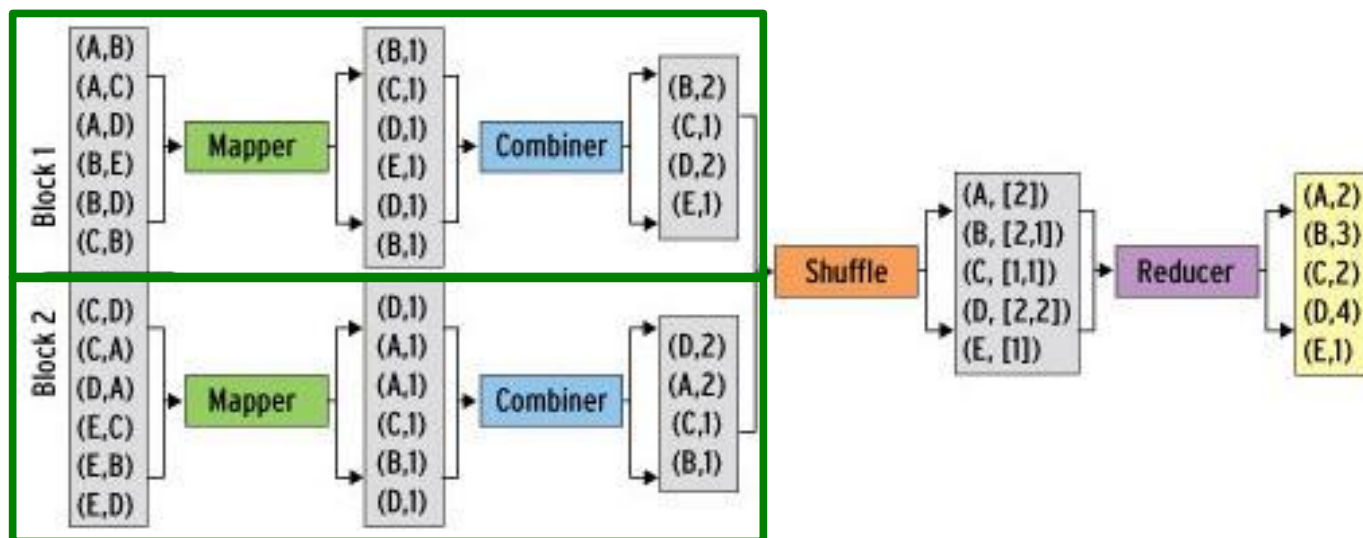
# Refinement: Combiners

- Often a Map task will produce many pairs of the form  $(k, v_1), (k, v_2), \dots$  for the same key  $k$ 
  - E.g., popular words in the word count example
- **Can save network time by pre-aggregating values in the mapper:**
  - $\text{combine}(k, \text{list}(v_1)) \rightarrow v_2$
  - Combiner is usually same as the reduce function
- Works only if reduce function is commutative and associative



# Refinement: Combiners

- **Back to our word counting example:**
  - Combiner combines the values of all keys of a single mapper (single machine):



- Much less data needs to be copied and shuffled!

# Refinement: Partition Function

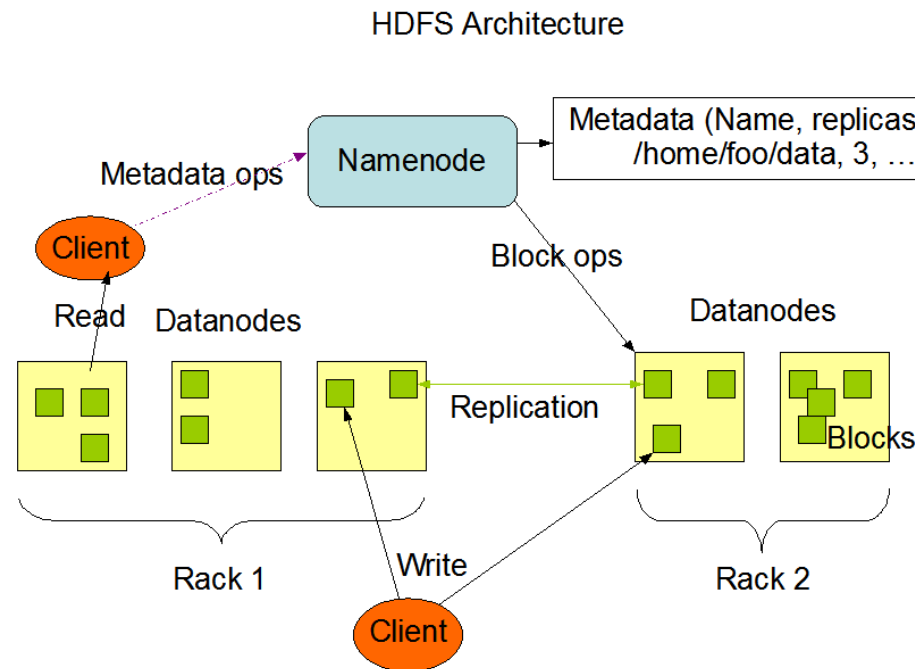
- **Want to control how keys get partitioned**
  - Inputs to map tasks are created by contiguous splits of input file
  - Reduce needs to ensure that records with the same intermediate key end up at the same worker
- **System uses a default partition function:**
  - **$\text{hash}(\text{key}) \bmod R$**
- **Sometimes useful to override the hash function:**
  - E.g.,  $\text{hash}(\text{hostname}(\text{URL})) \bmod R$  ensures URLs from a host end up in the same output file

# Hadoop

- Open source project managed by the Apache Software Foundation
- Current Framework includes:
  - Implementation of MapReduce
  - YARN
  - Hadoop Distributed File System (HDFS)
  - Hadoop Commons
- Users include Amazon, Facebook, and Ebay<sup>1</sup>

<sup>1</sup>: <https://wiki.apache.org/Hadoop/PoweredBy>

# Hadoop Distributed File System



[https://hadoop.apache.org/docs/r1.2.1/hdfs\\_design.html](https://hadoop.apache.org/docs/r1.2.1/hdfs_design.html)



# Pig

- High-level scripting platform
- Provides an abstraction from MapReduce

```
input_lines = LOAD '/tmp/word.txt' AS (line:chararray);
words = FOREACH input_lines GENERATE FLATTEN(TOKENIZE(line)) AS word;
filtered_words = FILTER words BY word MATCHES '\\w+';
word_groups = GROUP filtered_words BY word;
word_count = FOREACH word_groups GENERATE COUNT(filtered_words) AS
count, group AS word;
ordered_word_count = ORDER word_count BY count DESC;
STORE ordered_word_count INTO '/tmp/results.txt';
```

# Hive

- Data warehouse software
- SQL-like interface (Hive-QL)
- Specify what not how!

# Example: Join By Map-Reduce

- Compute the natural join  $R(A,B) \bowtie S(B,C)$
- $R$  and  $S$  are each stored in files
- Tuples are pairs  $(a,b)$  or  $(b,c)$

