

Large-Scale Data Processing with MapReduce



Slides courtesy of Jimmy Lin, University of Maryland - <http://www.umiacs.umd.edu/~jimmylin/>



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

Google[™]
processes 20 PB a day (2008)

eBay[®]
9 PB of user data +
>50 TB/day (11/2011)

36 PB of user data +
80-90 TB/day (6/2010)



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



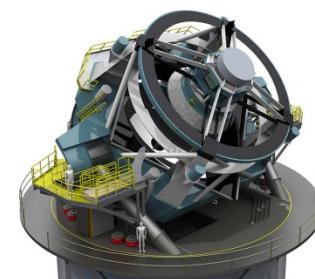
640K ought to be
enough for anybody.

JPMorganChase 150 PB on 50k+ servers
running 15k apps



Wayback Machine: 3 PB +
100 TB/month (3/2009)

LHC: ~15 PB a year
(at full capacity)



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

How much data?

cheap commodity clusters
+ simple, distributed programming models
= **data-intensive computing for the masses!**

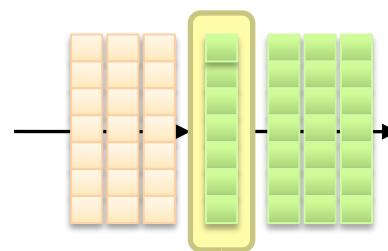
Parallel computing is hard!

Fundamental issues

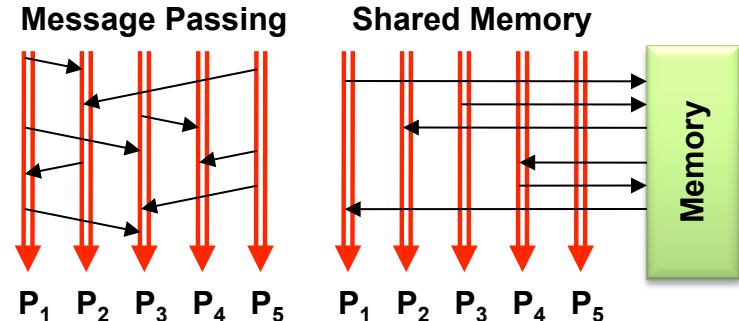
scheduling, data distribution, synchronization, inter-process communication, robustness, fault tolerance, ...

Architectural issues

Flynn's taxonomy (SIMD, MIMD, etc.), network typology, bisection bandwidth UMA vs. NUMA, cache coherence



Different programming models

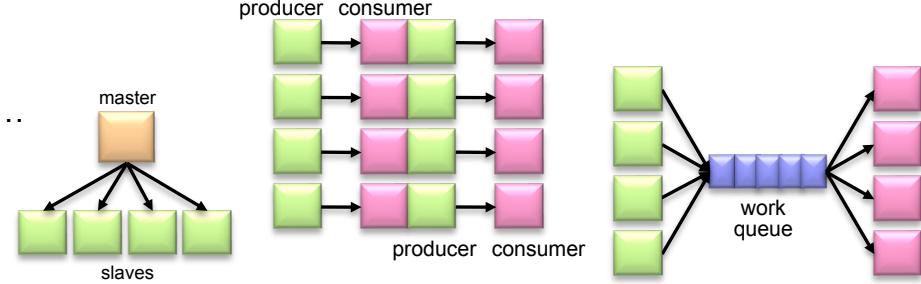


Common problems

livelock, deadlock, data starvation, priority inversion...
dining philosophers, sleeping barbers, cigarette smokers, ...

Different programming constructs

mutexes, conditional variables, barriers, ...
masters/slaves, producers/consumers, work queues, ...



The reality: programmer shoulders the burden of managing concurrency...

(We want to develop new algorithms, not debugging race conditions)

Where the rubber meets the road

- Concurrency is difficult to reason about
 - At the scale of datacenters (even across datacenters)
 - In the presence of failures
 - In terms of multiple interacting services
- 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: Ricardo Guimarães Herrmann



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



The datacenter *is* the computer!

What's the point?

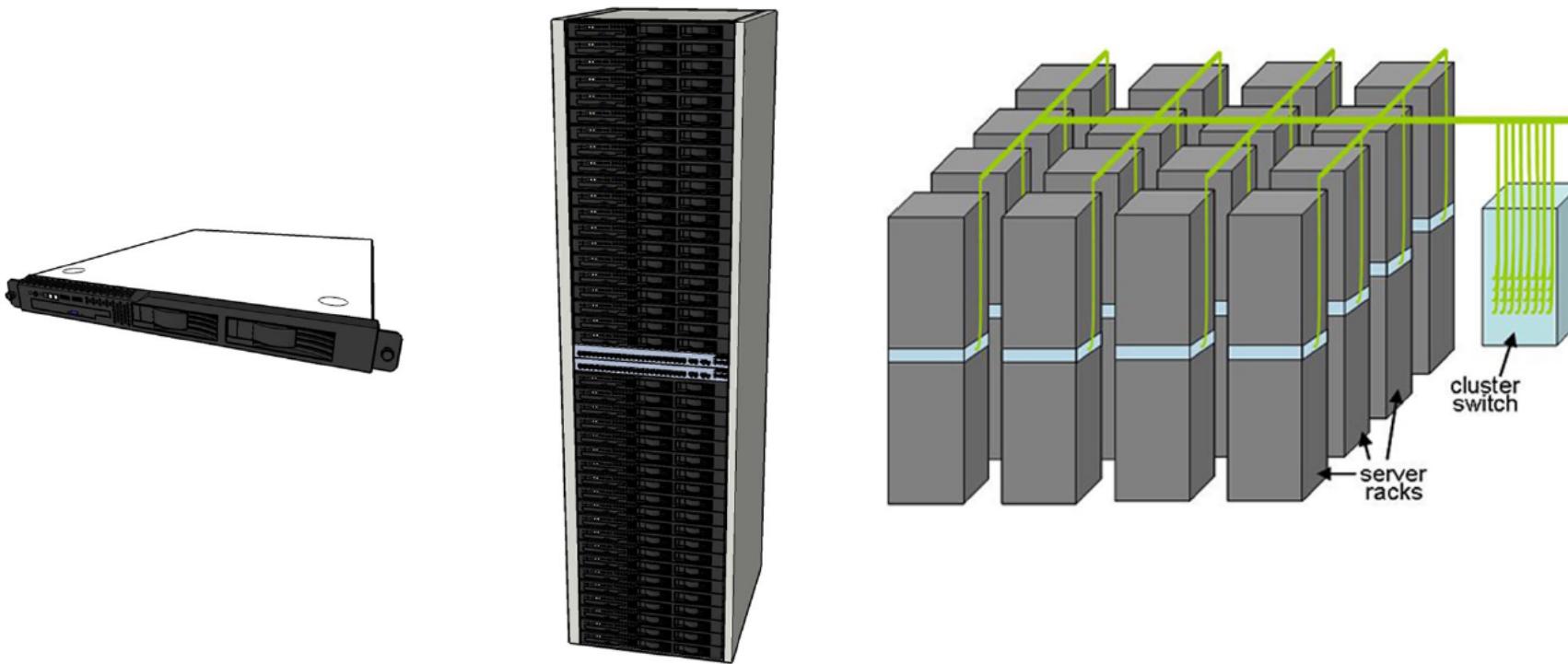
- It's all about the right level of abstraction
- Hide system-level details from the developers
 - No more race conditions, lock contention, etc.
- Separating the *what* from *how*
 - Developer specifies the computation that needs to be performed
 - Execution framework (“runtime”) handles actual execution

The datacenter *is* the computer!

“Big Ideas”

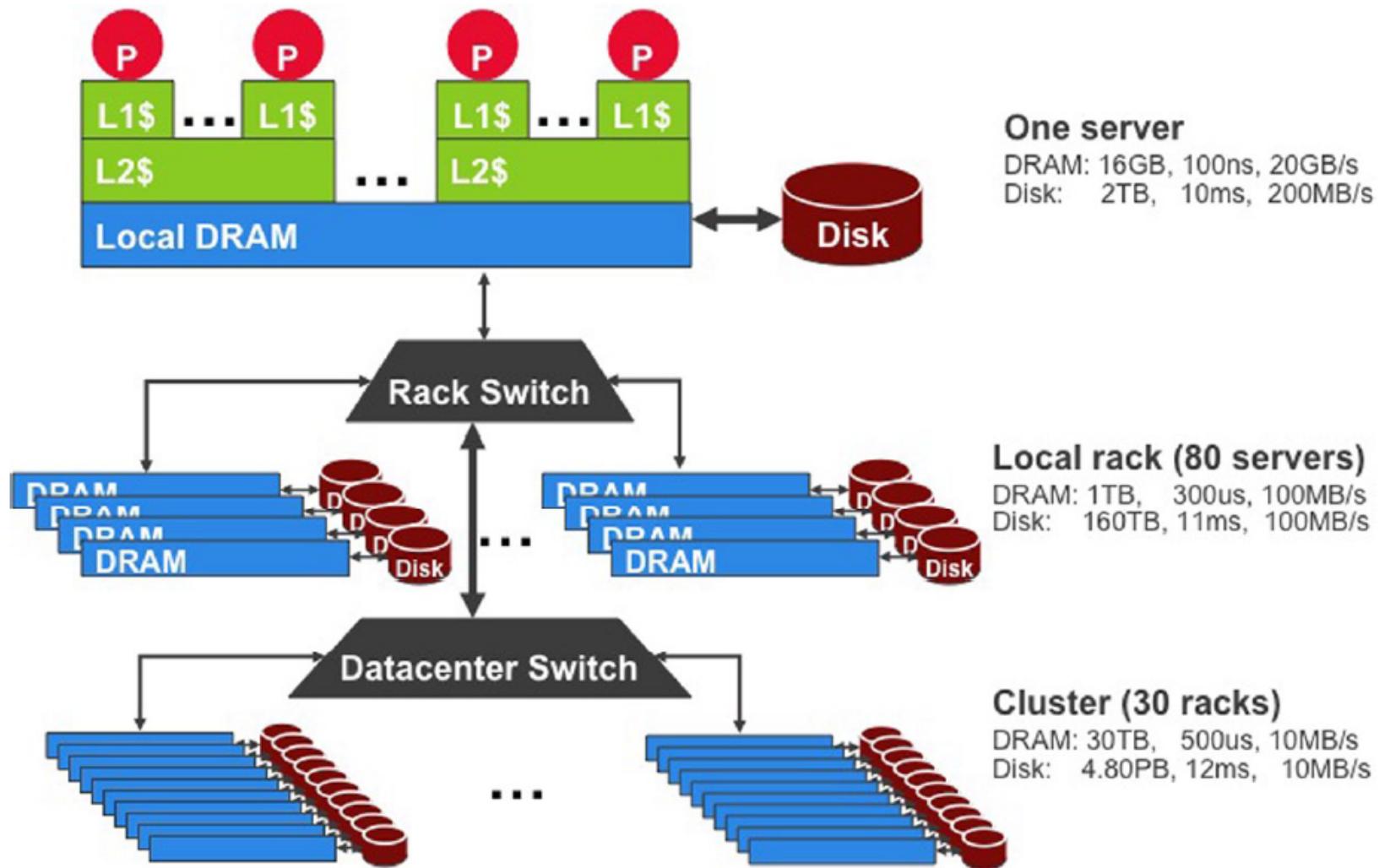
- Scale “out”, not “up”
 - Limits of SMP and large shared-memory machines
- Move processing to the data
 - Clusters have limited bandwidth
- Process data sequentially, avoid random access
 - Seek times are expensive, disk throughput is reasonable
- Seamless scalability
 - From the mythical man-month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

Building Blocks

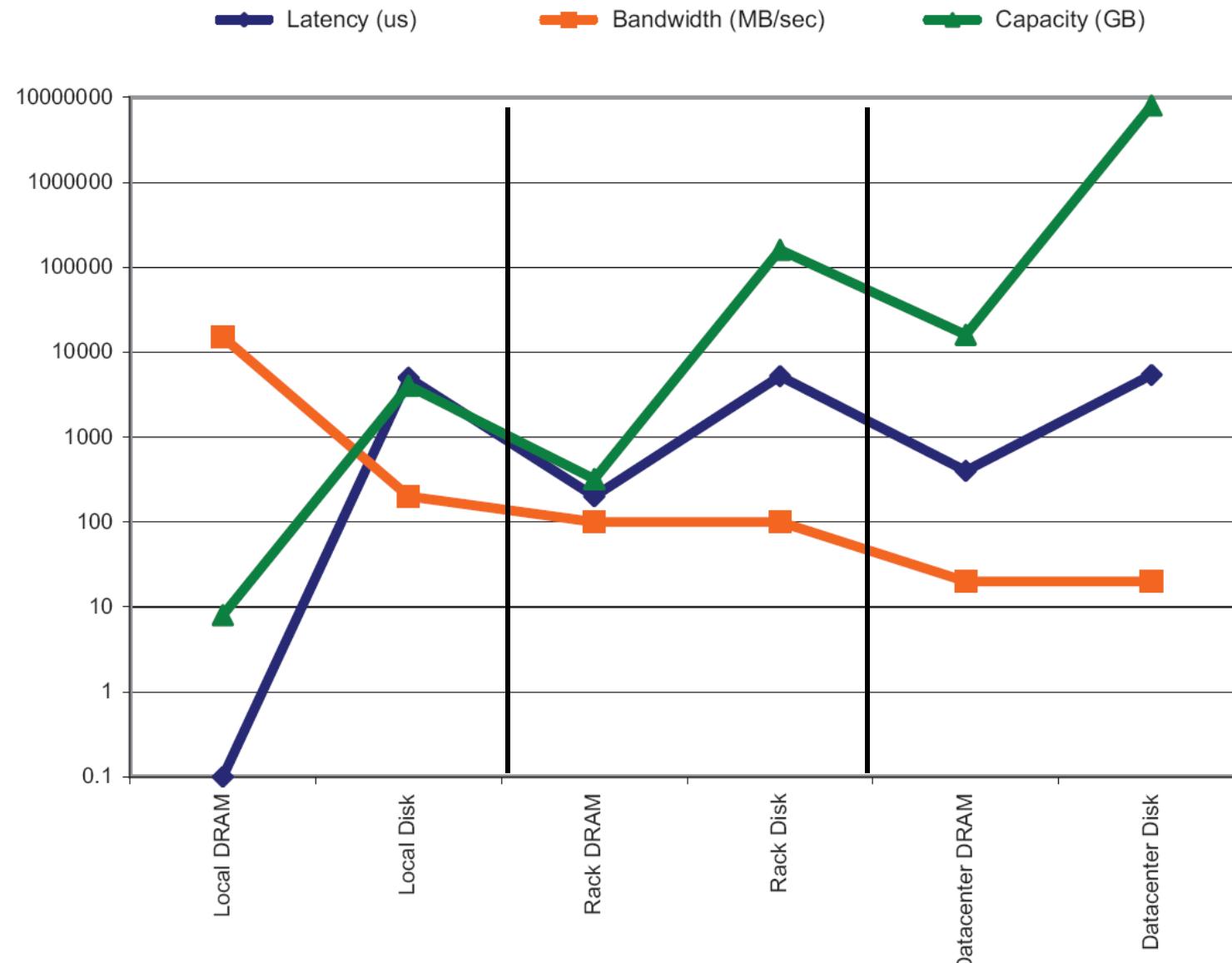


Source: Barroso and Urs Hözle (2009)

Storage Hierarchy



Storage Hierarchy



Source: Barroso and Urs Hözle (2009)

Introduction to MapReduce

Typical Large-Data Problem

- Iterate over a large number of records

Map Extract something of interest from each

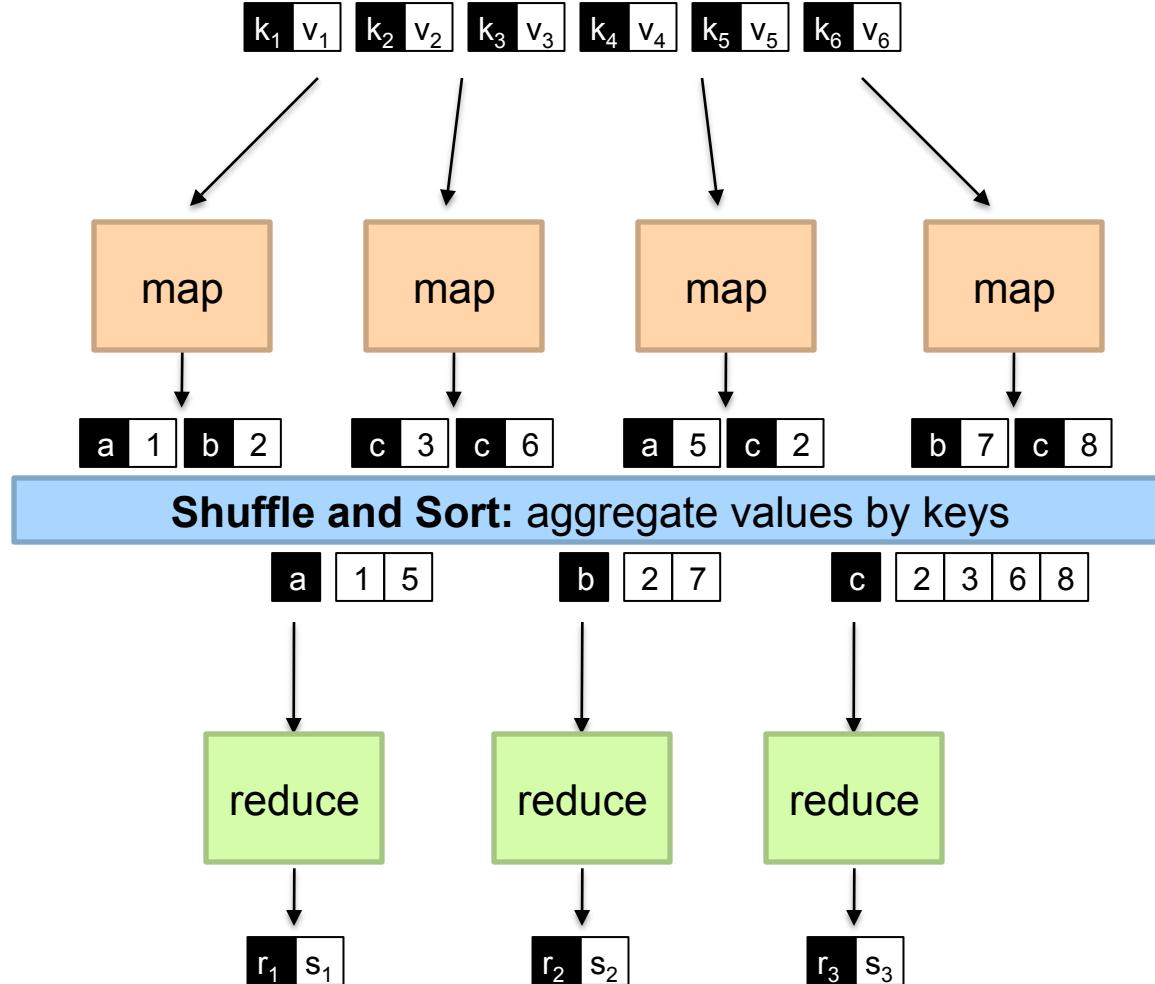
- Shuffle and sort intermediate results
- Aggregate intermediate results
- Generate final output

Reduce

Key idea: provide a functional abstraction for these two operations

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...

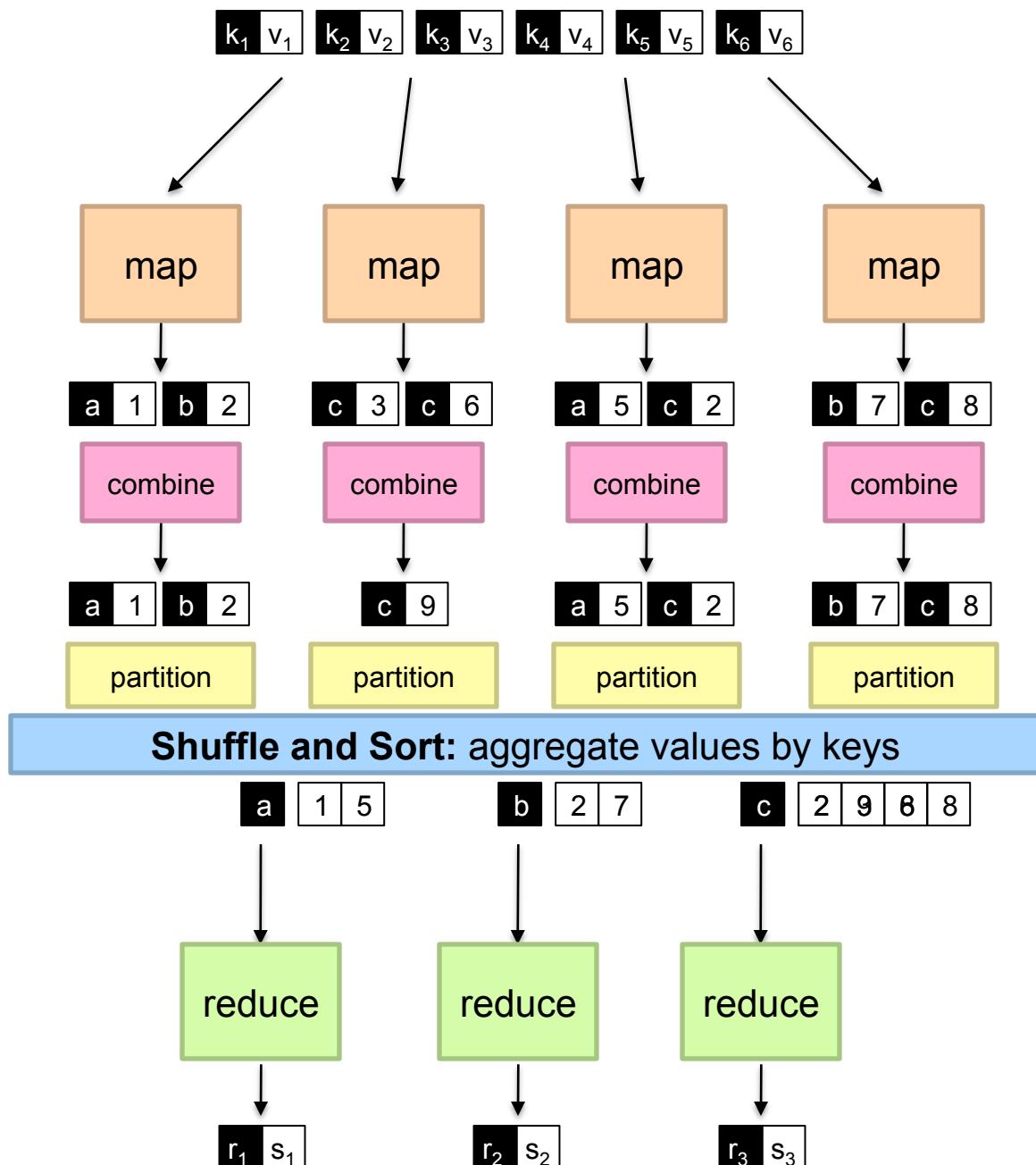
What's “everything else”?

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

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



“Hello World”: Word Count

```
1: class MAPPER
2:   method MAP(docid a, doc d)
3:     for all term t  $\in$  doc d do
4:       EMIT(term t, count 1)

1: class REDUCER
2:   method REDUCE(term t, counts [c1, c2, ...])
3:     sum  $\leftarrow$  0
4:     for all count c  $\in$  counts [c1, c2, ...] do
5:       sum  $\leftarrow$  sum + c
6:     EMIT(term t, count s)
```

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
 - Original development led by Yahoo
 - Now an Apache open source project
 - Emerging as the *de facto* big data stack
 - Rapidly expanding software ecosystem
- Lots of custom research implementations
 - For GPUs, cell processors, etc.
 - Includes variations of the basic programming model



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
- 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)

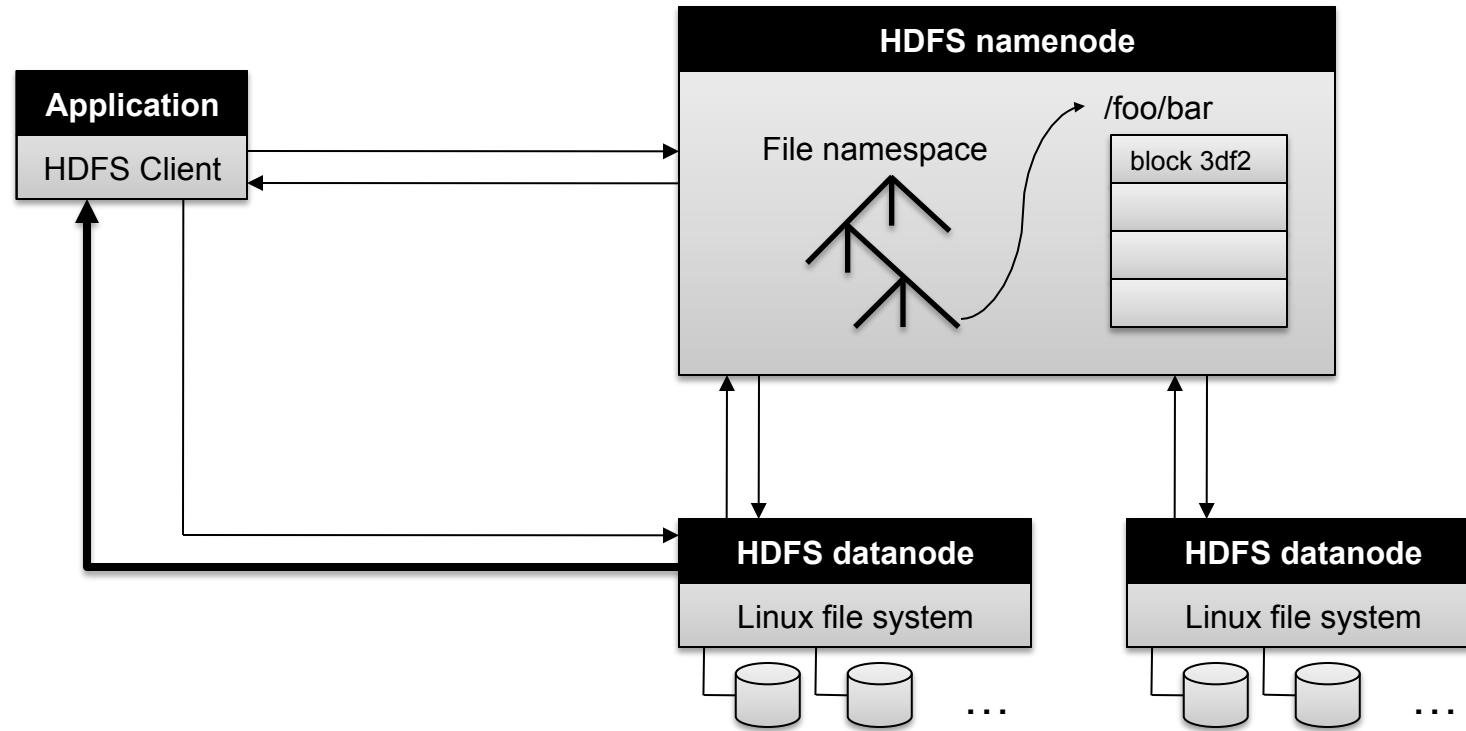
HDFS = GFS clone (same basic ideas)

From GFS to HDFS

- Terminology differences:
 - GFS master = Hadoop namenode
 - GFS chunkservers = Hadoop datanodes
- Functional differences:
 - File appends in HDFS is relatively new
 - HDFS performance is (likely) slower

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

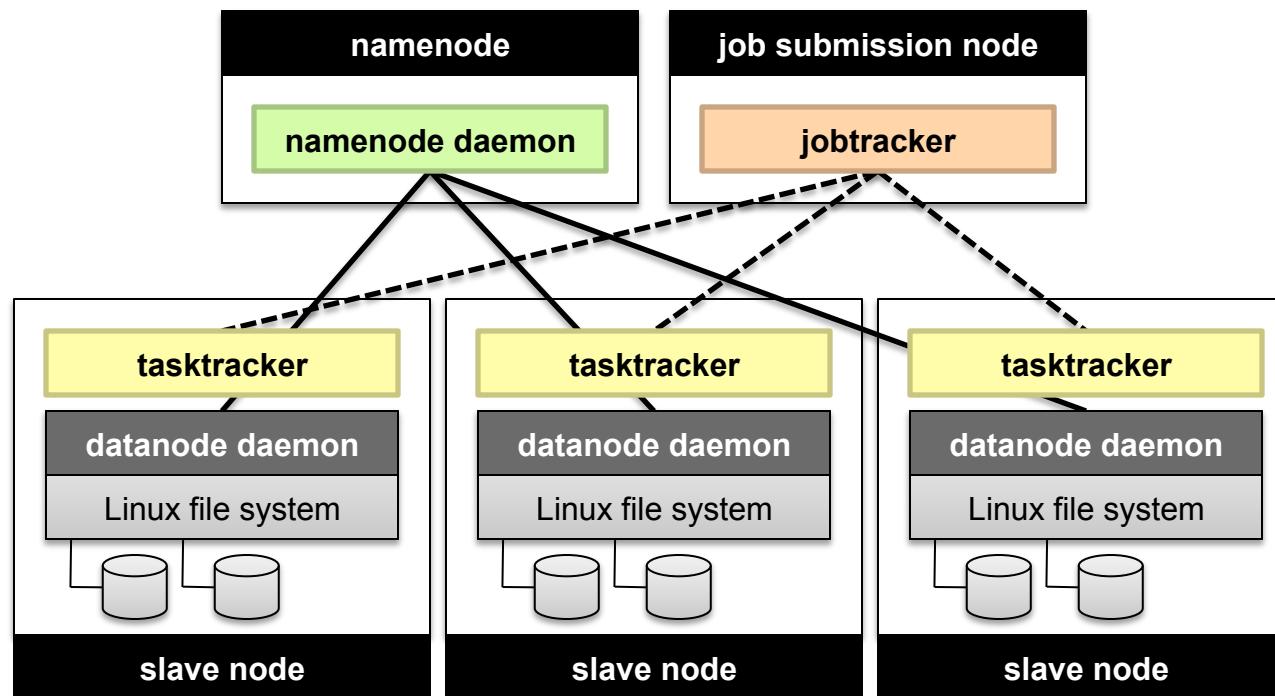
HDFS Architecture



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...



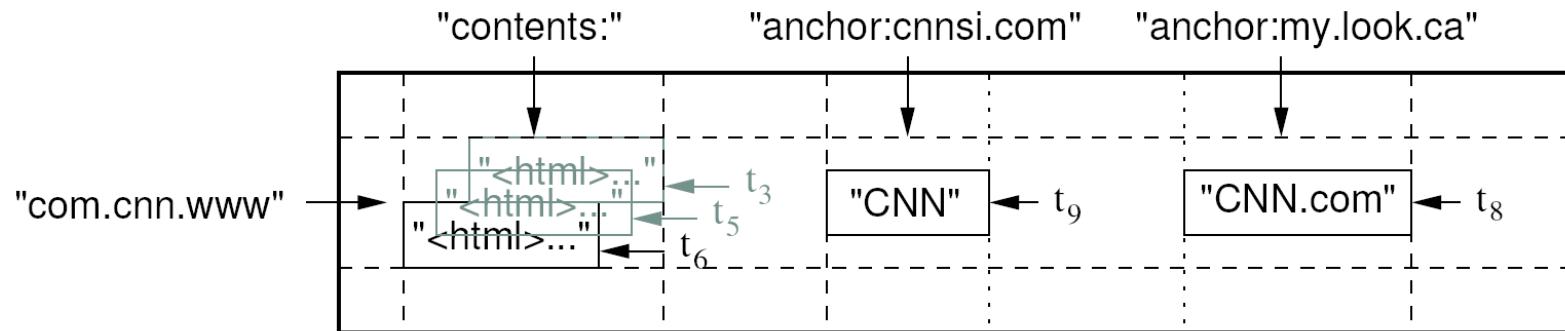
Hadoop Ecosystem Tour

From GFS to Bigtable

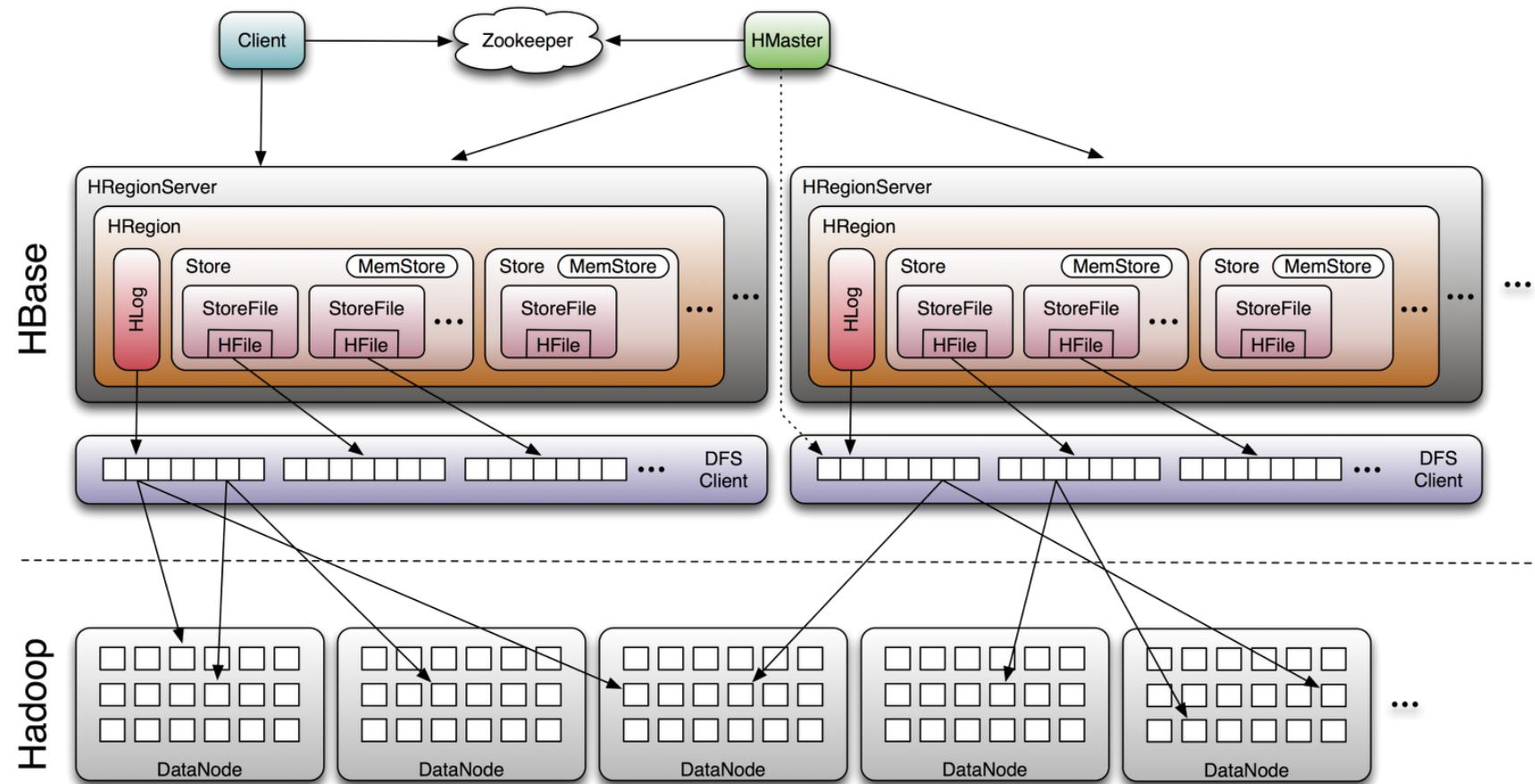
- Google's GFS is a distributed file system
- Bigtable is a storage system for structured data
 - Built on top of GFS
 - Solves many GFS issues: real-time access, short files, short reads
 - Serves as a source and a sink for MapReduce jobs

Bigtable: Data Model

- A table is a sparse, distributed, persistent multidimensional sorted map
- Map indexed by a row key, column key, and a timestamp
 - (row:string, column:string, time:int64) → uninterpreted byte array
- Supports lookups, inserts, deletes
 - Single row transactions only



HBase



Need for High-Level Languages

- Hadoop is great for large-data processing!
 - But writing Java programs for everything is verbose and slow
 - Analysts don't want to (or can't) write Java
- Solution: develop higher-level data processing languages
 - Hive: HQL is like SQL
 - Pig: Pig Latin is a dataflow language

Hive and Pig

- Hive: data warehousing application in Hadoop
 - Query language is HQL, variant of SQL
 - Tables stored on HDFS as flat files
 - Developed by Facebook, now open source
- Pig: large-scale data processing system
 - Scripts are written in Pig Latin, a dataflow language
 - Developed by Yahoo!, now open source
 - Roughly 1/3 of all Yahoo! internal jobs
- Common idea:
 - Provide higher-level language to facilitate large-data processing
 - Higher-level language “compiles down” to Hadoop jobs



Hive: Example

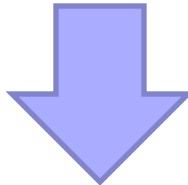
- Hive looks similar to an SQL database
- Relational join on two tables:
 - Table of word counts from Shakespeare collection
 - Table of word counts from the bible

```
SELECT s.word, s.freq, k.freq FROM shakespeare s
JOIN bible k ON (s.word = k.word) WHERE s.freq >= 1 AND k.freq >= 1
ORDER BY s.freq DESC LIMIT 10;
```

the	25848	62394
I	23031	8854
and	19671	38985
to	18038	13526
of	16700	34654
a	14170	8057
you	12702	2720
my	11297	4135
in	10797	12445
is	8882	6884

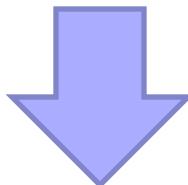
Hive: Behind the Scenes

```
SELECT s.word, s.freq, k.freq FROM shakespeare s  
JOIN bible k ON (s.word = k.word) WHERE s.freq >= 1 AND k.freq >= 1  
ORDER BY s.freq DESC LIMIT 10;
```



(Abstract Syntax Tree)

```
(TOK_QUERY (TOK_FROM (TOK_JOIN (TOK_TABREF shakespeare s) (TOK_TABREF bible k) (= (. (TOK_TABLE_OR_COL s)  
word) (. (TOK_TABLE_OR_COL k) word)))) (TOK_INSERT (TOK_DESTINATION (TOK_DIR TOK_TMP_FILE)) (TOK_SELECT  
(TOK_SELEXPR (. (TOK_TABLE_OR_COL s) word)) (TOK_SELEXPR (. (TOK_TABLE_OR_COL s) freq)) (TOK_SELEXPR (.  
(TOK_TABLE_OR_COL k) freq))) (TOK_WHERE (AND (>= (. (TOK_TABLE_OR_COL s) freq) 1) (>= (. (TOK_TABLE_OR_COL k)  
freq) 1))) (TOK_ORDERBY (TOK_TABSORTCOLNAMEDESC (. (TOK_TABLE_OR_COL s) freq))) (TOK_LIMIT 10)))
```



(one or more of MapReduce jobs)

Hive: Behind the Scenes

STAGE DEPENDENCIES:

Stage-1 is a root stage
Stage-2 depends on stages: Stage-1
Stage-0 is a root stage

STAGE PLANS:

Stage: Stage-1
Map Reduce
Alias -> Map Operator Tree:
s
 TableScan
 alias: s
 Filter Operator
 predicate:
 expr: (freq >= 1)
 type: boolean
 Reduce Output Operator
 key expressions:
 expr: word
 type: string
 sort order: +
 Map-reduce partition columns:
 expr: word
 type: string
 tag: 0
 value expressions:
 expr: freq
 type: int
 expr: word
 type: string
k
 TableScan
 alias: k
 Filter Operator
 predicate:
 expr: (freq >= 1)
 type: boolean
 Reduce Output Operator
 key expressions:
 expr: word
 type: string
 sort order: +
 Map-reduce partition columns:
 expr: word
 type: string
 tag: 1
 value expressions:
 expr: freq
 type: int

Stage: Stage-2
Map Reduce
Alias -> Map Operator Tree:
hdfs://localhost:8022/tmp/hive-training/364214370/10002
 Reduce Output Operator
 key expressions:
 expr: _col1
 type: int
 sort order: -
 tag: -1
 value expressions:
 expr: _col0
 type: string
 expr: _col1
 type: int
 expr: _col2
 type: int
 Reduce Operator Tree:
 Extract
 Limit
 File Output Operator
 compressed: false
 GlobalTableId: 0
 table:
 input format: org.apache.hadoop.mapred.TextInputFormat
 output format: org.apache.hadoop.hive.io.IgnoreKeyTextOutputFormat

Stage: Stage-0
Fetch Operator
limit: 10

Pig: Example

Task: Find the top 10 most visited pages in each category

Visits

User	Url	Time
Amy	cnn.com	8:00
Amy	bbc.com	10:00
Amy	flickr.com	10:05
Fred	cnn.com	12:00

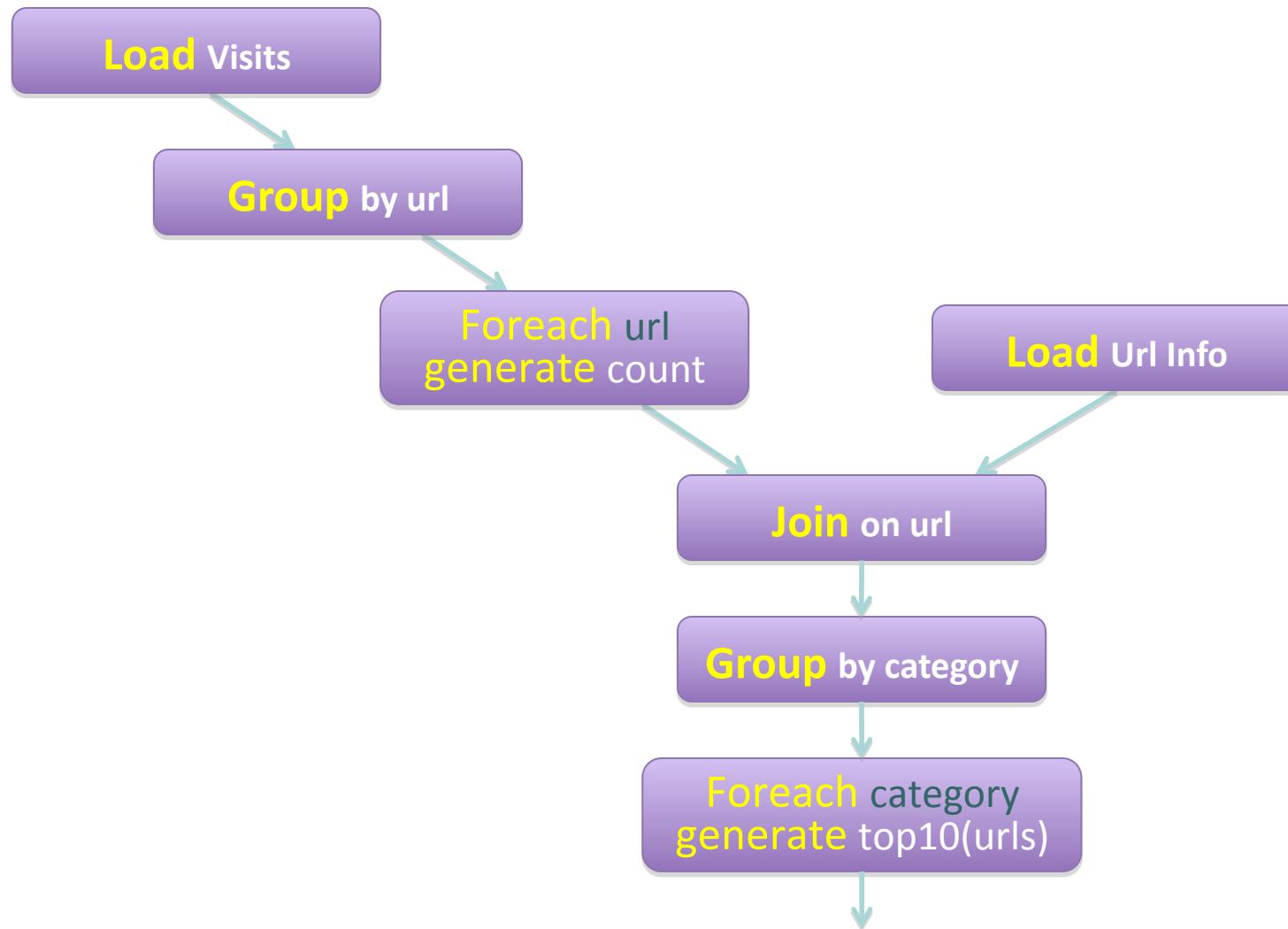


Url Info

Url	Category	PageRank
cnn.com	News	0.9
bbc.com	News	0.8
flickr.com	Photos	0.7
espn.com	Sports	0.9



Pig Query Plan



Pig Script

```
visits = load '/data/visits' as (user, url, time);  
gVisits = group visits by url;  
visitCounts = foreach gVisits generate url, count(visits);  
urlInfo = load '/data/urlInfo' as (url, category, pRank);  
visitCounts = join visitCounts by url, urlInfo by url;  
gCategories = group visitCounts by category;  
topUrls = foreach gCategories generate top(visitCounts,10);  
  
store topUrls into '/data/topUrls';
```

Pig Script in Hadoop

