

CSE 444: Database Internals

Lectures 21 MapReduce

Magda Balazinska - CSE 444, Spring 2013

1

References

- [MapReduce: Simplified Data Processing on Large Clusters](#). Jeffrey Dean and Sanjay Ghemawat. OSDI'04.

Magda Balazinska - CSE 444, Spring 2013

2

Outline

- Review high-level MR ideas from 344
- Discuss implementation in greater detail

Magda Balazinska - CSE 444, Spring 2013

3

Map Reduce Review

- Google: [Dean 2004]
- Open source implementation: Hadoop
- MapReduce = high-level programming model and implementation for large-scale parallel data processing

Magda Balazinska - CSE 444, Spring 2013

4

MapReduce Motivation

- Not designed to be a DBMS
- Designed to simplify task of writing parallel programs
 - A simple programming model that applies to many large-scale computing problems
- Hides messy details in MapReduce runtime library:
 - Automatic parallelization
 - Load balancing
 - Network and disk transfer optimizations
 - Handling of machine failures
 - Robustness
 - **Improvements to core library benefit all users of library!**

Magda Balazinska - CSE 444, Spring 2013

5

content in part from: Jeff Dean

Data Processing at Massive Scale

- Want to process petabytes of data and more
- Massive parallelism:
 - 100s, or 1000s, or 10000s servers (think data center)
 - Many hours
- Failure:
 - If medium-time-between-failure is 1 year
 - Then 10000 servers have one failure / hour

Magda Balazinska - CSE 444, Spring 2013

6

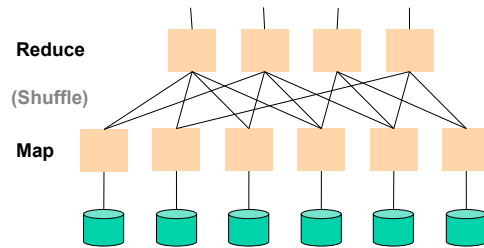
Data Storage: GFS/HDFS

- MapReduce job input is a file
- Common implementation is to store files in a highly scalable file system such as **GFS/HDFS**
 - GFS: Google File System
 - HDFS: Hadoop File System
- Each data file is split into M blocks (64MB or more)
- Blocks are stored on random machines & replicated
- Files are append only

Magda Balazinska - CSE 444, Spring 2013

7

Observation: Your favorite parallel algorithm...



Magda Balazinska - CSE 444, Spring 2013

8

Typical Problems Solved by MR

- Read a lot of data
- **Map**: extract something you care about from each record
- Shuffle and Sort
- **Reduce**: aggregate, summarize, filter, transform
- Write the results

Outline stays the same,
map and reduce change to fit
the problem

Magda Balazinska - CSE 444, Spring 2013

9

slide source: Jeff Dean

Data Model

Files !

A file = a bag of **(key, value)** pairs

A MapReduce program:

- Input: a bag of **(inputkey, value)** pairs
- Output: a bag of **(outputkey, value)** pairs

Magda Balazinska - CSE 444, Spring 2013

10

Step 1: the **MAP** Phase

User provides the **MAP**-function:

- Input: **(input key, value)**
- Output: **bag** of **(intermediate key, value)**

System applies map function in parallel to all **(input key, value)** pairs in the input file

Magda Balazinska - CSE 444, Spring 2013

11

Step 2: the **REDUCE** Phase

User provides the **REDUCE** function:

- Input:
(intermediate key, bag of values)
- Output (original MR paper): **bag** of output **(values)**
- Output (Hadoop): **bag** of **(output key, values)**

System groups all pairs with the same intermediate key, and passes the bag of values to the REDUCE function

Magda Balazinska - CSE 444, Spring 2013

12

Example

- Counting the number of occurrences of each word in a large collection of documents
- Each Document
 - The **key** = document id (**did**)
 - The **value** = set of words (**word**)

```
map(String key, String value):
// key: document name
// value: document contents
for each word w in value:
    EmitIntermediate(w, "1");
```

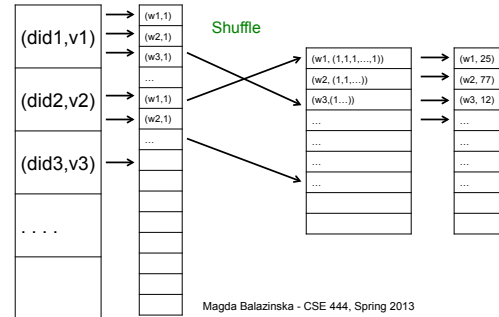
```
reduce(String key, Iterator values):
// key: a word
// values: a list of counts
int result = 0;
for each v in values:
    result += ParseInt(v);
Emit(AsString(result));
```

Magda Balazinska - CSE 444, Spring 2013

13

MAP

REDUCE



Magda Balazinska - CSE 444, Spring 2013

14

Jobs v.s. Tasks

- A **MapReduce Job**
 - One single "query", e.g. count the words in all docs
 - More complex queries may consist of multiple jobs
- A **Map Task**, or a **Reduce Task**
 - A group of instantiations of the map-, or reduce-function, which are scheduled on a single worker

Magda Balazinska - CSE 444, Spring 2013

15

Workers

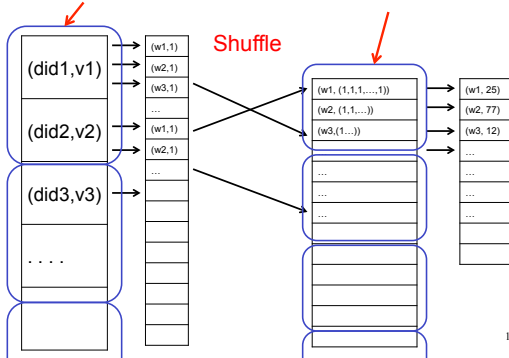
- A **worker** is a process that executes one task at a time
- Typically there is one worker per processor, hence 4 or 8 per node
- Often talk about "slots"
 - E.g., Each server has 2 map slots and 2 reduce slots

Magda Balazinska - CSE 444, Spring 2013

16

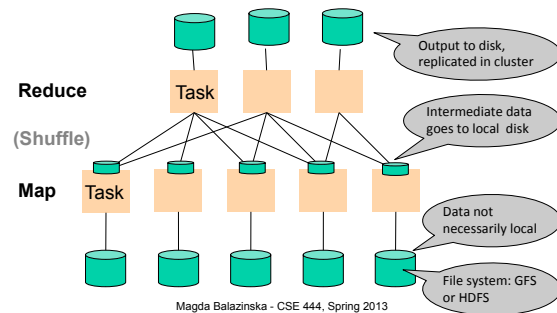
MAP Tasks

REDUCE Tasks



17

Parallel MapReduce Details



Magda Balazinska - CSE 444, Spring 2013

MapReduce Implementation

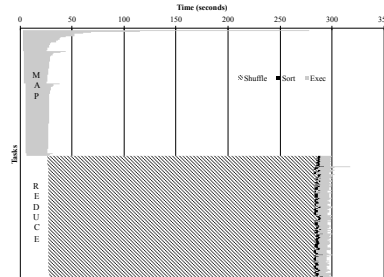
- There is one master node
- Input file gets partitioned further into M' splits
 - Each split is a contiguous piece of the input file
- Master assigns *workers* (=servers) to the M' map tasks, keeps track of their progress
- Workers write their output to local disk
- Output of each map task is partitioned into R regions
- Master assigns workers to the R reduce tasks
- Reduce workers read regions from the map workers' local disks

Magda Balazinska - CSE 444, Spring 2013

19

Example MapReduce Execution

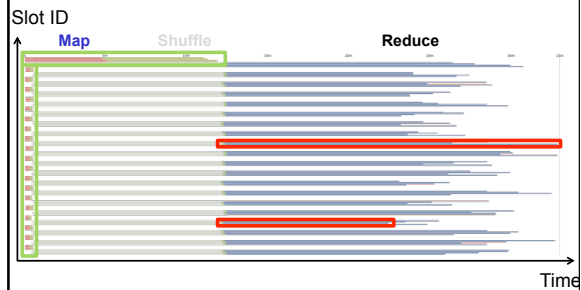
PageRank Application



Magda Balazinska - CSE 444, Spring 2013

20

Example: CloudBurst

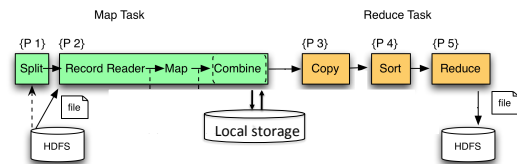


CloudBurst. Lake Washington Dataset (1.1GB). 80 Mappers 80 Reducers.

Magda Balazinska - CSE 444, Spring 2013

21

MapReduce Phases



Magda Balazinska - CSE 444, Spring 2013

22

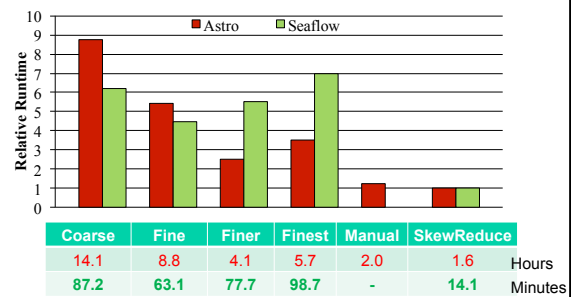
Interesting Implementation Details

- Worker failure:
 - Master pings workers periodically,
 - If down then reassigns its task to *another* worker
 - (≠ a parallel DBMS restarts whole query)
- How many map and reduce tasks:
 - Larger is better for load balancing
 - But more tasks also add overheads
 - (≠ parallel DBMS spreads ops across all nodes)

Magda Balazinska - CSE 444, Spring 2013

23

MapReduce Granularity Illustration



Magda Balazinska - CSE 444, Spring 2013

24

Interesting Implementation Details

Backup tasks:

- **Straggler** = a machine that takes unusually long time to complete one of the last tasks. Eg:
 - Bad disk forces frequent correctable errors (30MB/s → 1MB/s)
 - The cluster scheduler has scheduled other tasks on that machine
- Stragglers are a main reason for slowdown
- Solution: *pre-emptive backup execution of the last few remaining in-progress tasks*

Magda Balazinska - CSE 444, Spring 2013

25

Parallel DBMS vs MapReduce

- Parallel DBMS
 - Relational data model and schema
 - Declarative query language: SQL
 - Many pre-defined operators: relational algebra
 - Can easily combine operators into complex queries
 - Query optimization, indexing, and physical tuning
 - Streams data from one operator to the next without blocking
 - Can do more than just run queries: Data management
 - Updates and transactions, constraints, security, etc.

Magda Balazinska - CSE 444, Spring 2013

26

Parallel DBMS vs MapReduce

- MapReduce
 - Data model is a file with key-value pairs!
 - No need to “load data” before processing it
 - Easy to write user-defined operators
 - Can easily add nodes to the cluster (no need to even restart)
 - Uses less memory since processes one key-group at a time
 - Intra-query fault-tolerance thanks to results on disk
 - Intermediate results on disk also facilitate scheduling
 - Handles adverse conditions: e.g., stragglers
 - Arguably more scalable... but also needs more nodes!

Magda Balazinska - CSE 444, Spring 2013

27

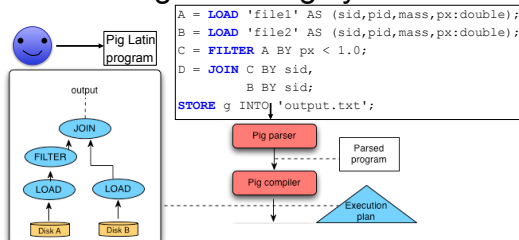
Declarative Languages on MR

- PIG Latin (Yahoo!)
 - New language, like Relational Algebra
 - Open source
- HiveQL (Facebook)
 - SQL-like language
 - Open source
- SQL / Tenzing (Google)
 - SQL on MR
 - Proprietary

Magda Balazinska - CSE 444, Spring 2013

28

Background: Pig system



Magda Balazinska - CSE 444, Spring 2013

29

MapReduce State

- Lots of extensions to address limitations
 - Capabilities to write DAGs of MapReduce jobs
 - Declarative languages (see 344)
 - Ability to read from structured storage (e.g., indexes)
 - Etc.
- Most companies use both types of engines
- Increased integration of both engines

Magda Balazinska - CSE 444, Spring 2013

30