

CSE544

Data Management

Lectures 14

Parallel Query Processing

Outline

- MapReduce
- Spark
- Snowflake

References

- Jeffrey Dean and Sanjay Ghemawat, [MapReduce: Simplified Data Processing on Large Clusters](#). OSDI'04
- Matei Zaharia, Mosharaf Chowdhury, Michael J. Franklin, Scott Shenker, Ion Stoica: [Spark: Cluster Computing with Working Sets](#). HotCloud 2010
- Dageville et al.: [The Snowflake Elastic Data Warehouse](#), SIGMOD'2016

MapReduce

Distributed File System (DFS)

- For very large files: TBs, PBs
- Each file partitioned into *chunks* (64MB)
- Each chunk replicated (≥ 3 times) – why?
- Implementations:
 - Google's DFS: **GFS**, proprietary
 - Hadoop's DFS: **HDFS**, open source

MapReduce

- Google:
 - Started around 2000
 - Paper published 2004
 - Discontinued September 2019
- Free variant: Hadoop
- MapReduce = high-level programming model and implementation for large-scale parallel data processing

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

Step 1: the **MAP** Phase

User provides the **MAP**-function:

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

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

Step 2: the **REDUCE** Phase

User provides the **REDUCE** function:

- Input: **(intermediate key, bag of values)**
- Output: bag of output **(values)**

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

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));
```

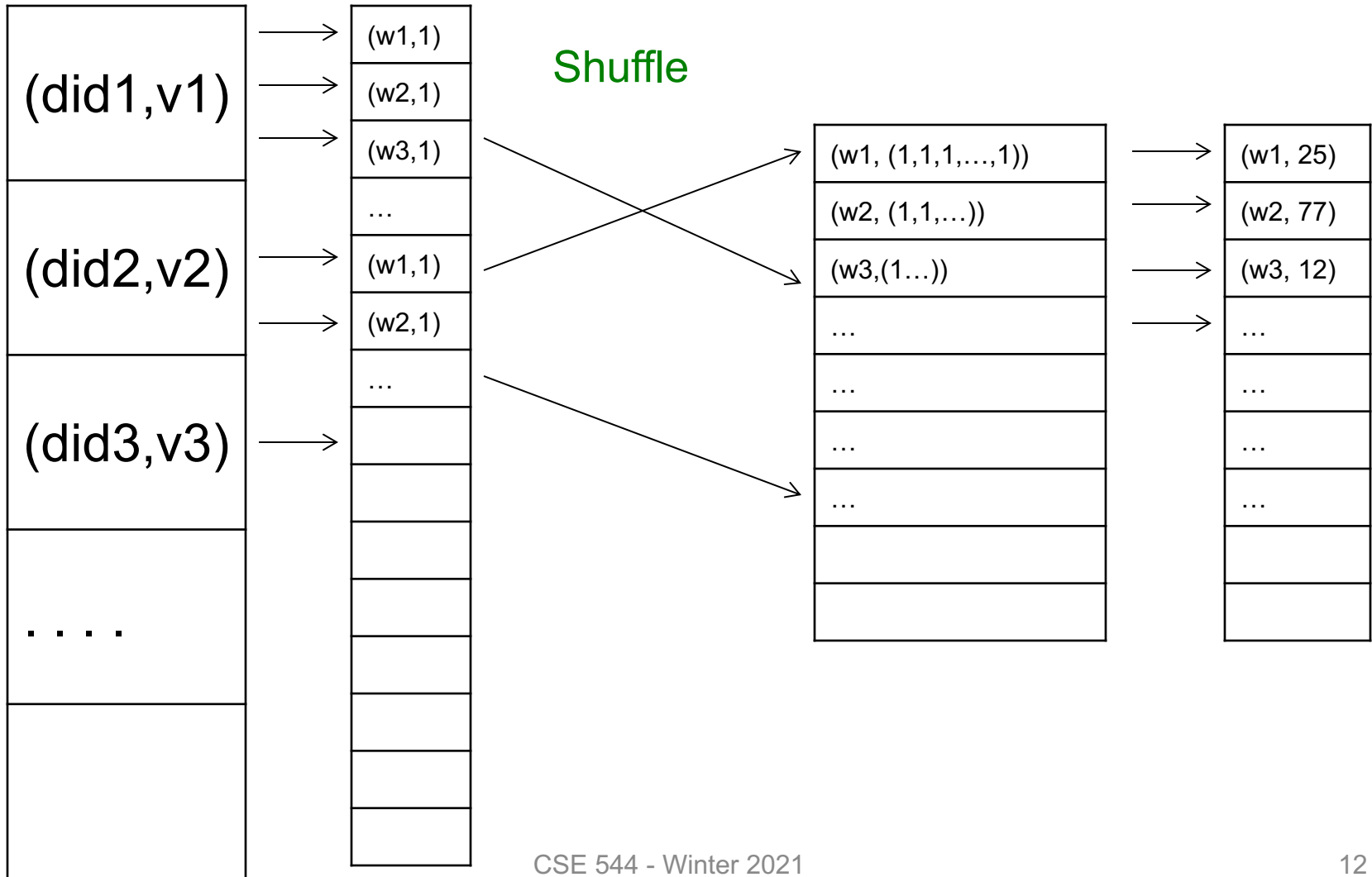
MapReduce = GroupBy-Aggregate

Occurrence(docID, word)

```
select  word, count(*)  
from    Occurrence  
group by word
```

MAP

REDUCE



Jobs v.s. Tasks

- A **MapReduce Job**
 - One simple “query”, e.g. count words in docs
 - Complex queries may require many jobs
- A **Map Task**, or a **Reduce Task**
 - A group of instantiations of the map-, or reduce-function, to be scheduled on a single worker

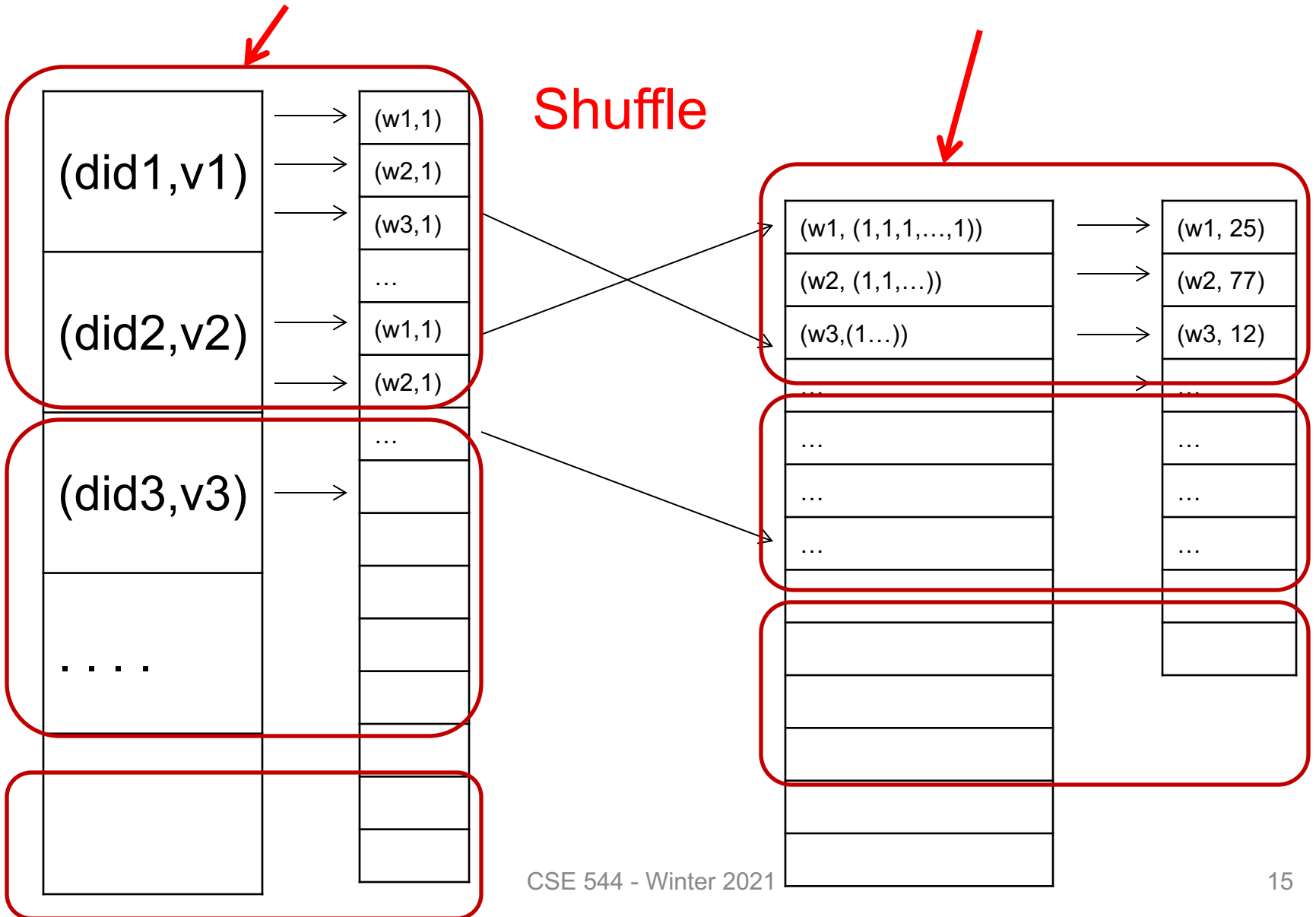
Workers

- 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

MAP Tasks

REDUCE Tasks

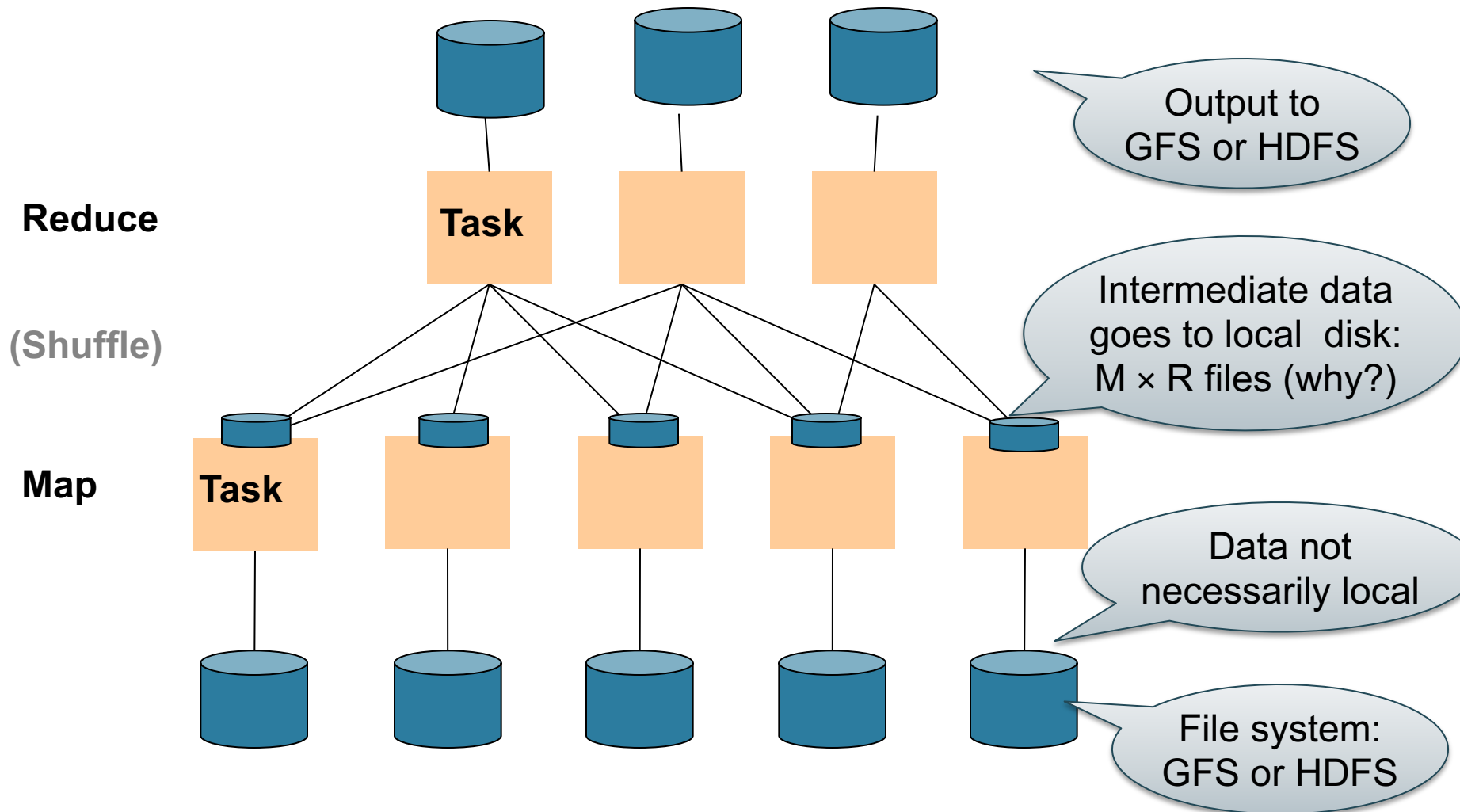
Shuffle



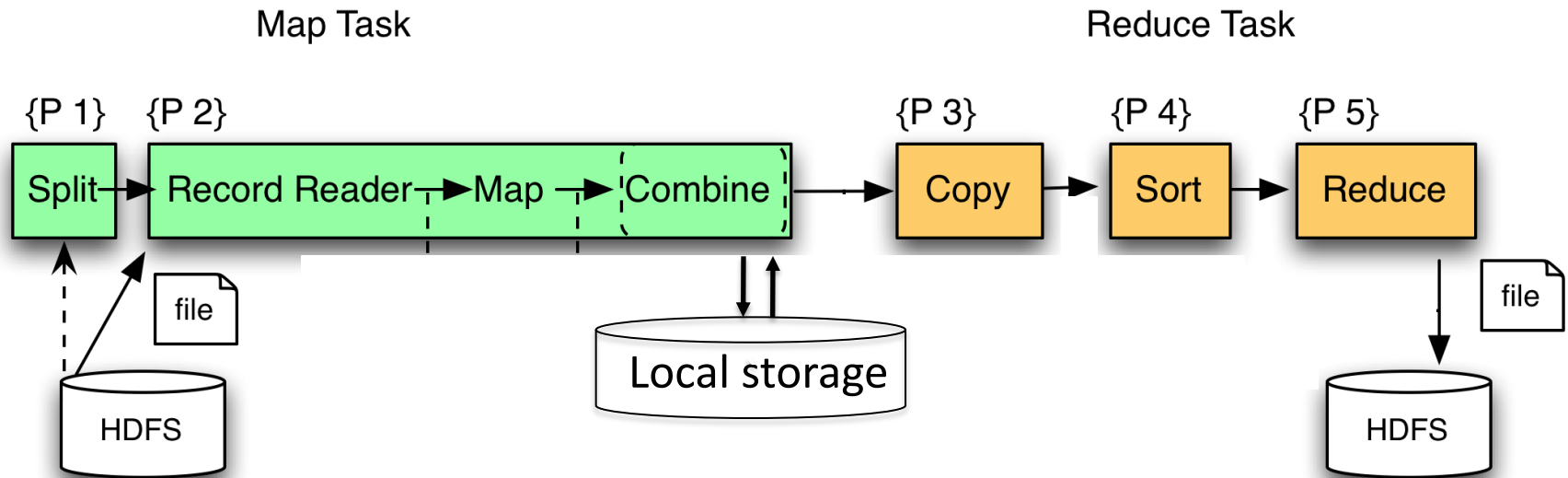
Fault Tolerance

- If one server fails once every year...
... then a job on 10,000 servers fails in 1h
- MapReduce handles fault tolerance by writing intermediate files to disk:
 - Mappers write file to local disk
 - Reducers read the local files (=reshuffling);
 - If reducer fails, new work re-reads local files

MapReduce Execution Details



MapReduce Phases



MapReduce v.s. Databases

Blog* by DeWitt and Stonebraker

- Schemas are good
- Indexes
- Skew (MR mitigates it somewhat – how?)
- The M*R problem – what is it?
- Parallel databases uses push (to sockets) instead of pull – what's the point?

*Original blog deleted, cached version still [here](#); slightly longer paper [here](#).

Spark

Spark

- Distributed processing over HDFS
- Multiple steps, including iterations
- Stores intermediate results in main memory
- Closer to relational algebra

Collections in Spark

- `RDD<T>` = an RDD collection of type T
 - Distributed
 - Not nested
 - Recoverable via lineage
- `Seq<T>` = a sequence
 - Local
 - May be nested

Programming in Spark

- Transformations (map, join...). **Lazy**
- Actions (count, reduce, save...). **Eager**

Example

Find in file `hdfs://logfile.log` the lines that:

- Start with “ERROR”
- Contain the string “sqlite”

```
s = SparkSession.builder()...getOrCreate();  
lines = s.read().textFile("hdfs://logfile.log");  
errors = lines.filter(l -> l.startsWith("ERROR"));  
sqlerrors = errors.filter(l -> l.contains("sqlite"));  
sqlerrors.collect();
```


Example

Find in file `hdfs://logfile.log` the lines that:

- Start with “ERROR”
- Contain the string “sqlite”

```
s = SparkSession.builder()...getOrCreate();  
  
sqlerrors = s.read().textFile("hdfs://logfile.log")  
    .filter(l -> l.startsWith("ERROR"))  
    .filter(l -> l.contains("sqlite"))  
    .collect();
```

“Call chaining” style

Example

The RDD s:

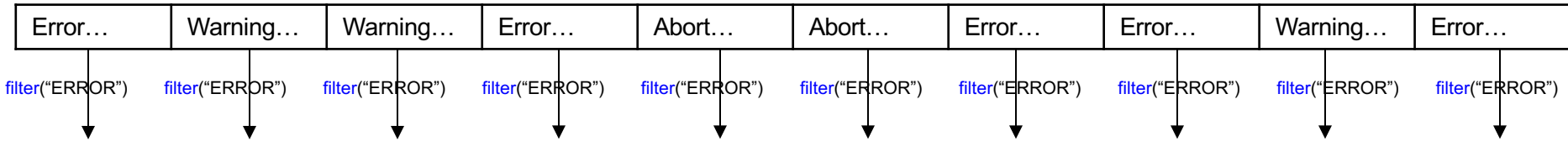
Error...	Warning...	Warning...	Error...	Abort...	Abort...	Error...	Error...	Warning...	Error...
----------	------------	------------	----------	----------	----------	----------	----------	------------	----------

```
s = SparkSession.builder()...getOrCreate();  
  
sqlerrors = s.read().textFile("hdfs://logfile.log")  
    .filter(l -> l.startsWith("ERROR"))  
    .filter(l -> l.contains("sqlite"))  
    .collect();
```

Example

The RDD s:

Parallel step 1

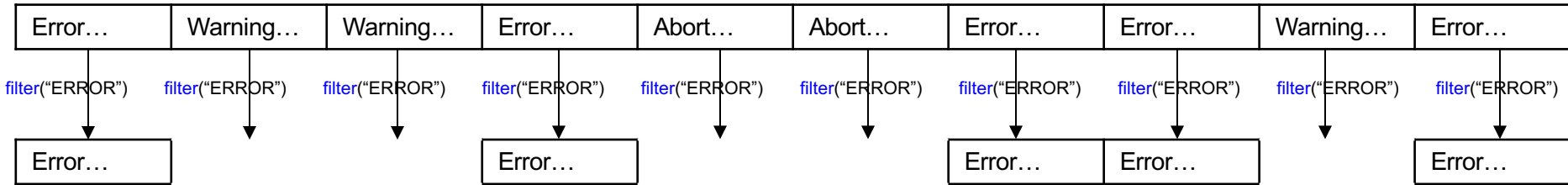


```
s = SparkSession.builder().getOrCreate();  
  
sqlerrors = s.read().textFile("hdfs://logfile.log")  
    .filter(l -> l.startsWith("ERROR"))  
    .filter(l -> l.contains("sqlite"))  
    .collect();
```

Example

The RDD s:

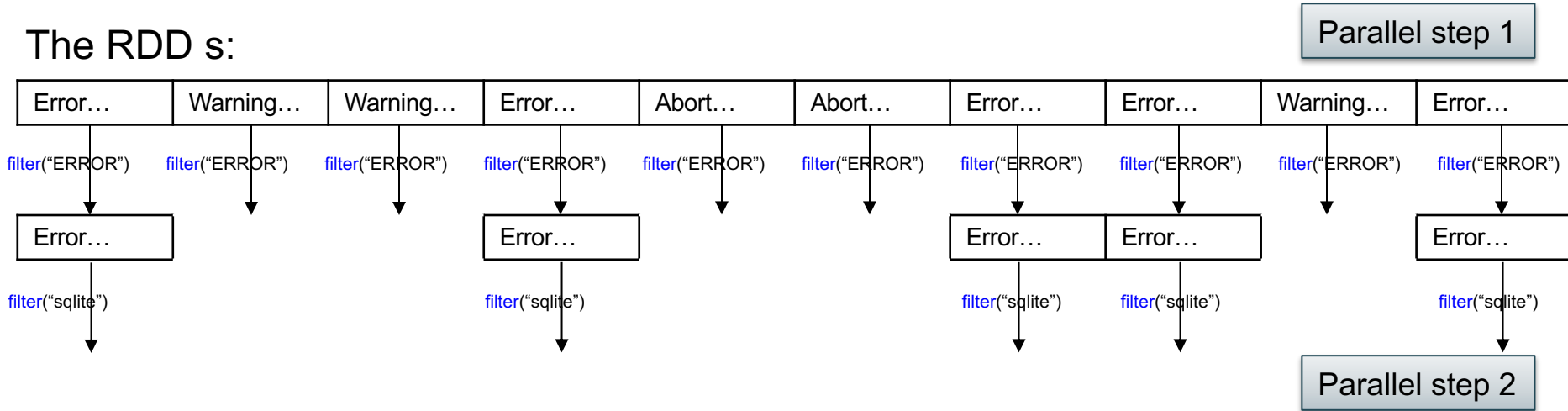
Parallel step 1



```
s = SparkSession.builder()...getOrCreate();  
  
sqlerrors = s.read().textFile("hdfs://logfile.log")  
    .filter(l -> l.startsWith("ERROR"))  
    .filter(l -> l.contains("sqlite"))  
    .collect();
```

Example

The RDD s:



```
s = SparkSession.builder()...getOrCreate();  
  
sqlerrors = s.read().textFile("hdfs://logfile.log")  
    .filter(l -> l.startsWith("ERROR"))  
    .filter(l -> l.contains("sqlite"))  
    .collect();
```

Fault Tolerance

- Parallel database systems: restart.
Expensive.
- MapReduce: write everything to disk,
redo. Slow.
- Spark: redo only what is needed.
Efficient.

Resilient Distributed Datasets

RDD = Resilient Distributed Dataset

- Distributed, immutable *lineage*
- Lineage = a relational algebra plan
- If a server crashes, recompute the lost partition of the RDD using the lineage

Persistence

```
lines = s.read().textFile("hdfs://logfile.log");  
errors = lines.filter(l->l.startsWith("ERROR"));  
sqlerrors = errors.filter(l->l.contains("sqlite"));  
sqlerrors.collect();
```

If any server fails before the end, then Spark must restart

Persistence

RDD:

hdfs://logfile.log

filter(...startsWith("ERROR"))
filter(...contains("sqlite"))

result

```
lines = s.read().textFile("hdfs://logfile.log");  
errors = lines.filter(l->l.startsWith("ERROR"));  
sqlerrors = errors.filter(l->l.contains("sqlite"));  
sqlerrors.collect();
```

If any server fails before the end, then Spark must restart

Persistence

RDD:

hdfs://logfile.log

filter(...startsWith("ERROR"))
filter(...contains("sqlite"))

result

```
lines = s.read().textFile("hdfs://logfile.log");  
errors = lines.filter(l->l.startsWith("ERROR"));  
sqlerrors = errors.filter(l->l.contains("sqlite"));  
sqlerrors.collect();
```

If any server fails before the end, then Spark must restart

```
lines = s.read().textFile("hdfs://logfile.log");  
errors = lines.filter(l->l.startsWith("ERROR"));  
errors.persist();  
sqlerrors = errors.filter(l->l.contains("sqlite"));  
sqlerrors.collect();
```

New RDD

Spark can recompute the result from errors

Persistence

RDD:

hdfs://logfile.log

filter(...startsWith("ERROR"))
filter(...contains("sqlite"))

result

```
lines = s.read().textFile("hdfs://logfile.log");  
errors = lines.filter(l->l.startsWith("ERROR"));  
sqlerrors = errors.filter(l->l.contains("sqlite"));  
sqlerrors.collect();
```

If any server fails before the end, then Spark must restart

hdfs://logfile.log

filter(..startsWith("ERROR"))

errors

filter(...contains("sqlite"))

result

```
lines = s.read().textFile("hdfs://logfile.log");  
errors = lines.filter(l->l.startsWith("ERROR"));  
errors.persist();  
sqlerrors = errors.filter(l->l.contains("sqlite"));  
sqlerrors.collect();
```

New RDD

Spark can recompute the result from errors

R(A,B)
S(A,C)

```
SELECT count(*) FROM R, S  
WHERE R.B > 200 and S.C < 100 and R.A = S.A
```

Example

```
R = strm.read().textFile("R.csv").map(parseRecord).persist();  
S = strm.read().textFile("S.csv").map(parseRecord).persist();
```

Parses each line into an object

persist

R(A,B)
S(A,C)

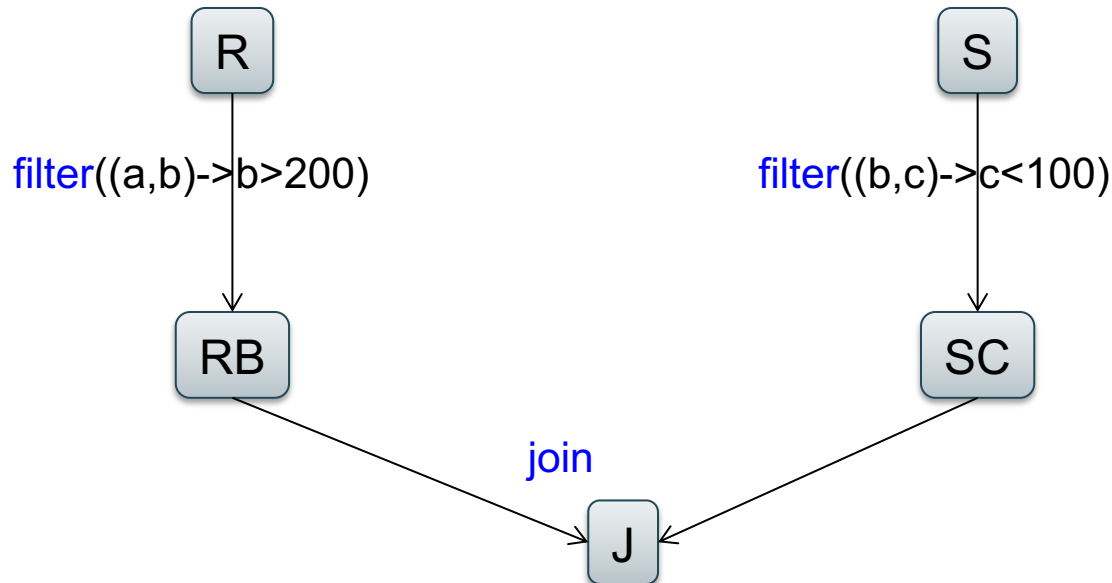
```
SELECT count(*) FROM R, S  
WHERE R.B > 200 and S.C < 100 and R.A = S.A
```

Example

```
R = strm.read().textFile("R.csv").map(parseRecord).persist();  
S = strm.read().textFile("S.csv").map(parseRecord).persist();  
RB = R.filter(t -> t.b > 200).persist();  
SC = S.filter(t -> t.c < 100).persist();  
J = RB.join(SC).persist();  
J.count();
```

transformations

action



Transformations:

<code>map(f : T -> U):</code>	<code>RDD<T> -> RDD<U></code>
<code>flatMap(f: T -> Seq(U)):</code>	<code>RDD<T> -> RDD<U></code>
<code>filter(f:T->Bool):</code>	<code>RDD<T> -> RDD<T></code>
<code>groupByKey():</code>	<code>RDD<(K,V)> -> RDD<(K,Seq[V])></code>
<code>reduceByKey(F:(V,V)-> V):</code>	<code>RDD<(K,V)> -> RDD<(K,V)></code>
<code>union():</code>	<code>(RDD<T>,RDD<T>) -> RDD<T></code>
<code>join():</code>	<code>(RDD<(K,V)>,RDD<(K,W)>) -> RDD<(K,(V,W))></code>
<code>cogroup():</code>	<code>(RDD<(K,V)>,RDD<(K,W)>)-> RDD<(K,(Seq<V>,Seq<W>))></code>
<code>crossProduct():</code>	<code>(RDD<T>,RDD<U>) -> RDD<(T,U)></code>

Actions:

<code>count():</code>	<code>RDD<T> -> Long</code>
<code>collect():</code>	<code>RDD<T> -> Seq<T></code>
<code>reduce(f:(T,T)->T):</code>	<code>RDD<T> -> T</code>
<code>save(path:String):</code>	Outputs RDD to a storage system e.g., HDFS

Spark 2.0

- DataFrames
 - Records, dynamically typed
- Datasets
 - Records, statically typed

Snowflake

Snowflake

- It is an SaaS – what is this? Give other examples of types of cloud services...

Snowflake

- It is an SaaS – what is this? Give other examples of types of cloud services...
- SaaS = software as a service
- Other examples:
 - Platform as a service (PaaS): e.g. Amazon's EC
 - Infrastructure as a service (virtual machines)
 - Software as a Service
 - Function as a Service: Amazon's Lambda

Snowflake

- Describe Snowflake's Data Storage

Snowflake

- Describe Snowflake's Data Storage

In class:

- S3:PUT/GET/DELETE
- Table → horizontal partition in files
- Blobs+PAX
- Temp storage → S3

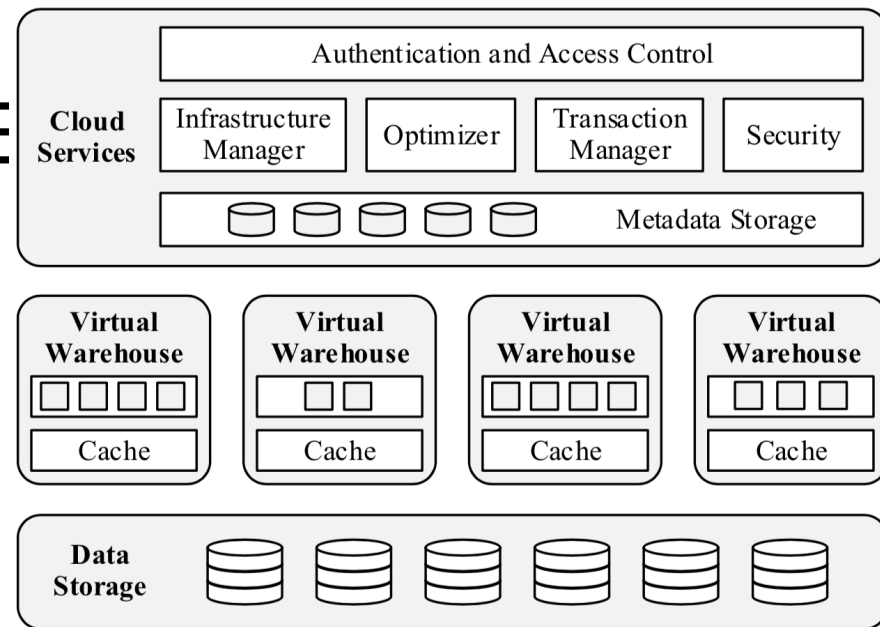


Figure 1: Multi-Cluster, Shared Data Architecture

Snowflake

- Describe Elasticity in Snowflake
- Describe failure handling in Snowflake

Snowflake

- Describe Elasticity in Snowflake
 - Virtual Warehouse (VW) serves one user
 - T-Shirt sizes: X-Small ... XX-Large
 - Small query may run on subset of VW
- Describe failure handling in Snowflake

Snowflake

- Describe Elasticity in Snowflake
 - Virtual Warehouse (VW) serves one user
 - T-Shirt sizes: X-Small ... XX-Large
 - Small query may run on subset of VW
- Describe failure handling in Snowflake
 - Restart the query
 - No partial retries (like MapReduce or Spark)

Snowflake

- Describe its execution engine

Snowflake

- Describe its execution engine
- Column-oriented (in class)
- Vectorized (“tuple batches” – in class)
- Push-based (in class)

Snowflake

- What does Snowflake use instead of indexes?

Snowflake

- What does Snowflake use instead of indexes?
- “Pruning”: for each file (recall: this is a horizontal partition of a table) and each attribute, it stores the min/max values in that column in that file; may skip files when not needed.

Conclusion

- Distributed data processing:
 - Spread the data to fit in main memory
 - Take advantage of parallelism
- “SQL is embarrassingly parallel”
 - Relational algebra: easy to parallelize
 - Hash-based algorithm suffer from skew