

Introduction to Database Systems CSE 414

Lecture 20: Map-Reduce and Spark

CSE 414 - Autumn 2018

1

Announcements

- HW6 is two parts
 - Running your Spark code locally
 - Running your Spark code on AWS
- Do all the local coding first, then run on AWS last.
- Useful:
<http://spark.apache.org/docs/latest/rdd-programming-guide.html>

CSE 414 - Autumn 2018

2

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

Paradigm stays the same,
change map and reduce
functions for different problems

CSE 414 - Autumn 2018

5
slide source: Jeff Dean

Data Model

Files!

A file = a bag of (**key**, **value**) pairs
Sounds familiar after HW5?

A MapReduce program:

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

CSE 414 - Autumn 2018

6

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");
```

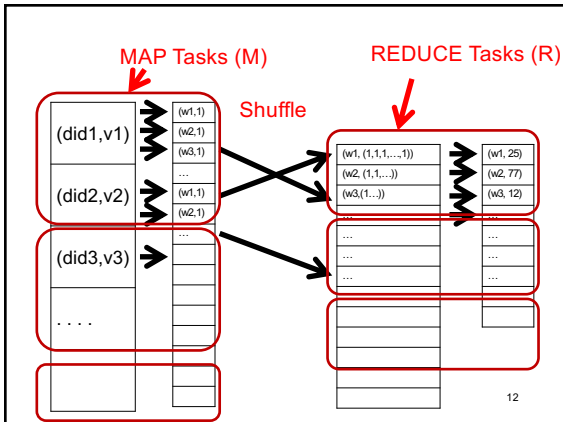
```
reduce(String key, Iterator values):  
  // key: a word  
  // values: a list of counts  
  int result = 0;  
  for each v in values:  
    result += ParseInt(v);  
  emit(AsString(result));
```

CSE 414 - Autumn 2018

11

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



Fault Tolerance

- If one server fails once every year...
... then a job with 10,000 servers will fail in less than one hour
- MapReduce handles fault tolerance by writing intermediate files to disk:
 - Mappers write file to local disk
 - Reducers read the files (=reshuffling); if the server fails, the reduce task is restarted on another server

Implementation

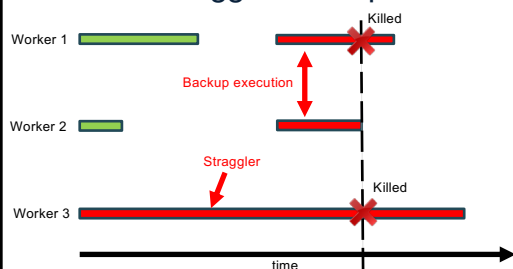
- There is one master node
- Master partitions input file into *M splits*, by key
- Master assigns *workers* (=servers) to the *M map tasks*, keeps track of their progress
- Workers write their output to local disk, partition into *R regions*
- Master assigns workers to the *R reduce tasks*
- Reduce workers read regions from the map workers' local disks

Interesting Implementation Details

Backup tasks:

- *Straggler* = a machine that takes unusually long time to complete one of the last tasks. E.g.:
 - 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*

Straggler Example



Using MapReduce in Practice: Implementing RA Operators in MR

Relational Operators in MapReduce

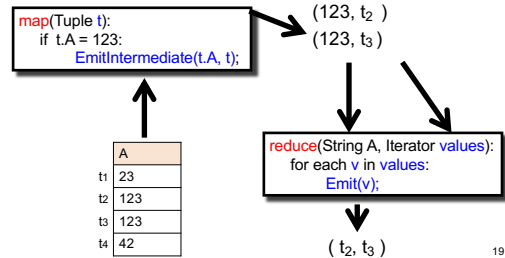
Given relations $R(A,B)$ and $S(B,C)$ compute:

- Selection: $\sigma_{A=123}(R)$
- Group-by: $\gamma_{A,\text{sum}(B)}(R)$
- Join: $R \bowtie S$ (Saved for later)

CSE 414 - Autumn 2018

18

Selection $\sigma_{A=123}(R)$



19

Selection $\sigma_{A=123}(R)$

```

map(Tuple t):
  if t.A = 123:
    EmitIntermediate(t.A, t);
  
```

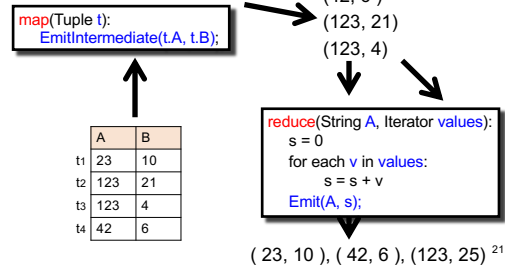
```

reduce(String A, Iterator values):
  for each v in values:
    Emit(v);
  
```

No need for reduce.
But need system hacking in Hadoop
to remove reduce from MapReduce

20

Group By $\gamma_{A,\text{sum}(B)}(R)$



(23, 10), (42, 6), (123, 25)²¹

Conclusions

- MapReduce offers a simple abstraction, and handles distribution + fault tolerance
- Speedup/scaleup achieved by allocating dynamically map tasks and reduce tasks to available server. However, skew is possible (e.g., one huge reduce task)
- Writing intermediate results to disk is necessary for fault tolerance, but very slow.
- Spark replaces this with "Resilient Distributed Datasets" = main memory + lineage

CSE 414 - Autumn 2018

27

Spark

A Case Study of the MapReduce Programming Paradigm

CSE 414 - Autumn 2018

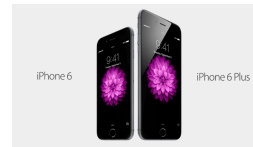
28

HW6

- HW6 will ask you to write SQL queries and MapReduce tasks using Spark
- You will get to “implement” SQL using MapReduce tasks
 - Can you beat Spark’s implementation?

CSE 414 - Autumn 2018

29



Parallel Data Processing @ 2010



CSE 414 - Autumn 2018

30

Spark

- Open source system from UC Berkeley
- Distributed processing over HDFS
- Differences from MapReduce (CSE 322):
 - Multiple steps, including iterations
 - Stores intermediate results in main memory
 - Closer to relational algebra (familiar to you)
- Details:
<http://spark.apache.org/examples.html>

CSE 414 - Autumn 2018

31

Spark

- Spark supports interfaces in Java, Scala, and Python
 - Scala: extension of Java with functions/closures
- We will illustrate use the Spark Java interface in this class
- Spark also supports a SQL interface (SparkSQL), and compiles SQL to its native Java interface

CSE 414 - Autumn 2018

32

Programming in Spark

- A Spark program consists of:
 - Transformations (map, reduce, join...). **Lazy**
 - Actions (count, reduce, save...). **Eager**
- **Eager**: operators are executed immediately
- **Lazy**: operators are not executed immediately
 - A *operator tree* is constructed in memory instead
 - Similar to a relational algebra tree

CSE 414 - Autumn 2018

33

Collections in Spark

- $RDD<T>$ = an RDD collection of type T
 - Distributed on many servers, not nested
 - Operations are done in parallel
 - Recoverable via lineage; more later
 - We use JavaRDD in HW 6
- $Seq<T>$ = a sequence
 - Local to one server, may be nested
 - Operations are done sequentially

CSE 414 - Autumn 2018

34

Transformations:	
<code>map(f : T -> U):</code>	<code>RDD<T> -> RDD<U></code>
<code>mapToPair(f : T -> K, V):</code>	<code>RDD<T> -> RDD<K, V></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 <small>35</small>

Transformations:	
<code>map(f : T -> U):</code>	<code>RDD<T> -> RDD<U></code>
<code>mapToPair(f : T -> K, V):</code>	<code>RDD<T> -> RDD<K, V></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 <small>36</small>

Example

Given a large log file `hdfs://logfile.log` retrieve all 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();
```

CSE 414 - Autumn 2018 37

Example

Recall: anonymous functions (lambda expressions) starting in Java 8

```
errors = lines.filter(l -> l.startsWith("ERROR"));
```

is the same as:

```
class FilterFn implements Function<Row, Boolean>{
  Boolean call (Row l)
  { return l.startsWith("ERROR"); }
}
errors = lines.filter(new FilterFn());
```

CSE 414 - Autumn 2018 38

Example

Recall: anonymous functions (lambda expressions) starting in Java 8

```
errors = lines.filter(l -> l.startsWith("ERROR"));
```

is the same as:

```
class FilterFn implements Function<Row, Boolean>{
  Boolean call (Row l)
  { return l.startsWith("ERROR"); }
}
errors = lines.filter(new FilterFn());
```

CSE 414 - Autumn 2018 39

Example

Given a large log file `hdfs://logfile.log` retrieve all 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();
```

CSE 414 - Autumn 2018 40

Example

Given a large log file `hdfs://logfile.log`
retrieve all lines that:

- Start with "ERROR"
- Contain the string "sqlite"

```
s = SparkSession.builder().getOrCreate();
lines = s.read().textFile("hdfs://logfile.log");
errors = lines.filter(1 -> 1.startsWith("ERROR"));
sqlerrors = errors.filter(1 -> 1.contains("sqlite"));
sqlerrors.collect();
```

Transformation:

Not executed yet...

Action:

triggers execution of entire program

CSE 414 - Autumn 2018

41

Example

Given a large log file `hdfs://logfile.log`
retrieve all lines that:

- Start with "ERROR"
- Contain the string "sqlite"

```
s = SparkSession.builder().getOrCreate();
sqlerrors = s.read().textFile("hdfs://logfile.log")
    .filter(1 -> 1.startsWith("ERROR"))
    .filter(1 -> 1.contains("sqlite"))
    .collect();
```

"Call chaining" style

42

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(1 -> 1.startsWith("ERROR"))
    .filter(1 -> 1.contains("sqlite"))
    .collect();
```

CSE 414 - Autumn 2018

43

Example

The RDD s:

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

Parallel step 1

```
s = SparkSession.builder().getOrCreate();
sqlerrors = s.read().textFile("hdfs://logfile.log")
    .filter(1 -> 1.startsWith("ERROR"))
    .filter(1 -> 1.contains("sqlite"))
    .collect();
```

CSE 414 - Autumn 2018

44

Example

The RDD s:

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

Parallel step 1

```
s = SparkSession.builder().getOrCreate();
sqlerrors = s.read().textFile("hdfs://logfile.log")
    .filter(1 -> 1.startsWith("ERROR"))
    .filter(1 -> 1.contains("sqlite"))
    .collect();
```

CSE 414 - Autumn 2018

45

Example

The RDD s:

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

Parallel step 1

```
s = SparkSession.builder().getOrCreate();
sqlerrors = s.read().textFile("hdfs://logfile.log")
    .filter(1 -> 1.startsWith("ERROR"))
    .filter(1 -> 1.contains("sqlite"))
    .collect();
```

Parallel step 2

CSE 414 - Autumn 2018

46

Fault Tolerance

- When a job is executed on x100 or x1000 servers, the probability of a failure is high
- Example: if a server fails once/year, then a job with 10000 servers fails once/hour
- Different solutions:
 - Parallel database systems: restart. Expensive.
 - MapReduce: write everything to disk, redo. Slow.
 - Spark: redo only what is needed. Efficient.

47

Resilient Distributed Datasets

- RDD = Resilient Distributed Dataset
 - Distributed, immutable and records its *lineage*
 - Lineage = expression that says how that relation was computed = a relational algebra plan
- Spark stores intermediate results as RDD
- If a server crashes, its RDD in main memory is lost. However, the driver (=master node) knows the lineage, and will simply recompute the lost partition of the RDD

CSE 414 - Autumn 2018

48

Persistence

```
lines = s.read().textFile("hdfs://logfile.log");
errors = lines.filter(1->l.startsWith("ERROR"));
sqlerrors = errors.filter(1->l.contains("sqlite"));
sqlerrors.collect();
```

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

CSE 414 - Autumn 2018

49

Persistence

RDD:

hdfs://logfile.log

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

result

```
lines = s.read().textFile("hdfs://logfile.log");
errors = lines.filter(1->l.startsWith("ERROR"));
sqlerrors = errors.filter(1->l.contains("sqlite"));
sqlerrors.collect();
```

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

CSE 414 - Autumn 2018

50

Persistence

RDD:

hdfs://logfile.log

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

result

```
lines = s.read().textFile("hdfs://logfile.log");
errors = lines.filter(1->l.startsWith("ERROR"));
sqlerrors = errors.filter(1->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(1->l.startsWith("ERROR"));
errors.persist();
sqlerrors = errors.filter(1->l.contains("sqlite"));
sqlerrors.collect();
```

New RDD

Spark can recompute the result from errors

CSE 414 - Autumn 2018

51

Persistence

RDD:

hdfs://logfile.log

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

result

```
lines = s.read().textFile("hdfs://logfile.log");
errors = lines.filter(1->l.startsWith("ERROR"));
sqlerrors = errors.filter(1->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(1->l.startsWith("ERROR"));
errors.persist();
sqlerrors = errors.filter(1->l.contains("sqlite"));
sqlerrors.collect();
```

New RDD

Spark can recompute the result from errors

CSE 414 - Autumn 2018

52

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
persisting on disk

CSE 414 - Autumn 2018 53

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

CSE 414 - Autumn 2018 54

Recap: Programming in Spark

- A Spark/Scala program consists of:
 - Transformations (map, reduce, join...). **Lazy**
 - Actions (count, reduce, save...). **Eager**
- RDD<T> = an RDD collection of type T
 - Partitioned, recoverable (through lineage), not nested
- Seq<T> = a sequence
 - Local to a server, may be nested

CSE 414 - Autumn 2018 55

Transformations:	
map(f : T -> U):	RDD<T> -> RDD<U>
mapToPair(f : T -> K, V):	RDD<T> -> RDD<K, V>
flatMap(f: T -> Seq(U)):	RDD<T> -> RDD<U>
filter(f:T->Bool):	RDD<T> -> RDD<T>
groupByKey():	RDD<(K,V)> -> RDD<(K,Seq[V])>
reduceByKey(F:(V,V)-> V):	RDD<(K,V)> -> RDD<(K,V)>
union():	(RDD<T>, RDD<T>) -> RDD<T>
join():	(RDD<(K,V)>, RDD<(K,W)>) -> RDD<(K, (V,W))>
cogroup():	(RDD<(K,V)>, RDD<(K,W)>) -> RDD<(K, (Seq<V>, Seq<W>))>
crossProduct():	(RDD<T>, RDD<U>) -> RDD<(T,U)>
Actions:	
count():	RDD<T> -> Long
collect():	RDD<T> -> Seq<T>
reduce(f:(T,T)->T):	RDD<T> -> T
save(path:String):	Outputs RDD to a storage system e.g., HDFS

56

Transformations:	
map(f : T -> U):	RDD<T> -> RDD<U>
mapToPair(f : T -> K, V):	RDD<T> -> RDD<K, V>
flatMap(f: T -> Seq(U)):	RDD<T> -> RDD<U>
filter(f:T->Bool):	RDD<T> -> RDD<T>
groupByKey():	RDD<(K,V)> -> RDD<(K,Seq[V])>
reduceByKey(F:(V,V)-> V):	RDD<(K,V)> -> RDD<(K,V)>
union():	(RDD<T>, RDD<T>) -> RDD<T>
join():	(RDD<(K,V)>, RDD<(K,W)>) -> RDD<(K, (V,W))>
cogroup():	(RDD<(K,V)>, RDD<(K,W)>) -> RDD<(K, (Seq<V>, Seq<W>))>
crossProduct():	(RDD<T>, RDD<U>) -> RDD<(T,U)>
Actions:	
count():	RDD<T> -> Long
collect():	RDD<T> -> Seq<T>
reduce(f:(T,T)->T):	RDD<T> -> T
save(path:String):	Outputs RDD to a storage system e.g., HDFS

57

Spark 2.0

The DataFrame and Dataset Interfaces

CSE 414 - Autumn 2018 58

DataFrames

- Like RDD, also an immutable distributed collection of data
- Organized into *named columns* rather than individual objects
 - Just like a relation
 - Elements are untyped objects called Row's
- Similar API as RDDs with additional methods
 - `people = spark.read().textFile(...);`
 - `ageCol = people.col("age");`
 - `ageCol.plus(10); // creates a new DataFrame`

CSE 414 - Autumn 2018

59

Datasets

- Similar to DataFrames, except that elements must be typed objects
- E.g.: `Dataset<People>` rather than `Dataset<Row>`
- Can detect errors during compilation time
- DataFrames are aliased as `Dataset<Row>` (as of Spark 2.0)
- You will use both Datasets and RDD APIs in HW6

CSE 414 - Autumn 2018

60

Datasets API: Sample Methods

- Functional API
 - `agg(Column expr, Column... exprs)`
Aggregates on the entire Dataset without groups.
 - `groupBy(String col1, String... cols)`
Groups the Dataset using the specified columns, so that we can run aggregation on them.
 - `join(Dataset<?> right)`
Join with another DataFrame.
 - `orderBy(Column... sortExprs)`
Returns a new Dataset sorted by the given expressions.
 - `select(Column... cols)`
Selects a set of column based expressions.
- "SQL" API
 - `SparkSession.sql("select * from R");`
- Look familiar?

CSE 414 - Autumn 2018

61

Conclusions

- Parallel databases
 - Predefined relational operators
 - Optimization
 - Transactions
- MapReduce
 - User-defined map and reduce functions
 - Must implement/optimize manually relational ops
 - No updates/transactions
- Spark
 - Predefined relational operators
 - Must optimize manually
 - No updates/transactions

CSE 414 - Autumn 2018

62