

CSCI-620

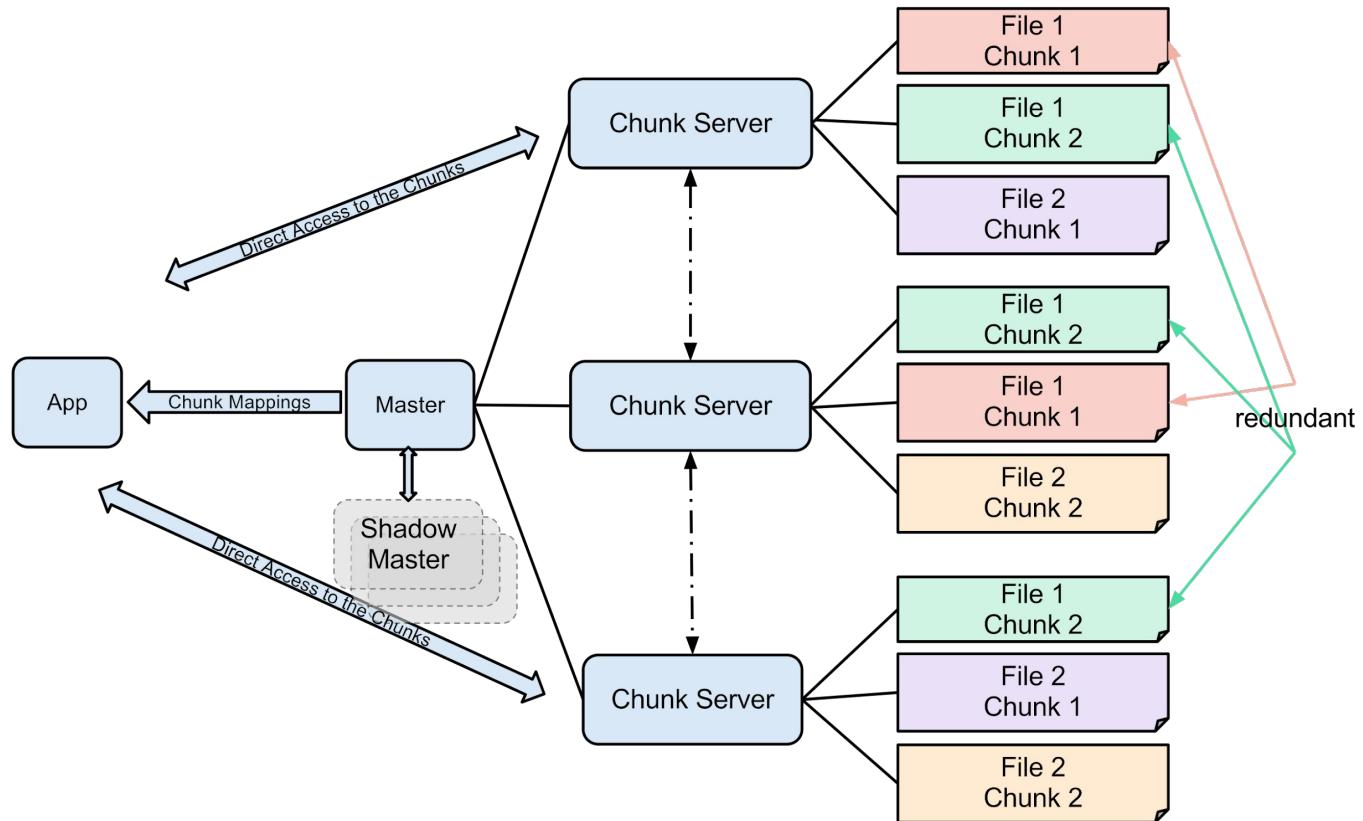
Distributed Data Processing

Distributed File Systems

- ▶ Many datasets are too large to fit on a single machine
- ▶ Unstructured data may not be easy to insert into a database
- ▶ Distributed file systems store data across a large number of servers

Google File System

- ▶ A distributed file system used by Google in the early 2000s
- ▶ Designed to run on a large number of cheap servers
- ▶ Somewhat slower, but tolerant of failures



Google File System

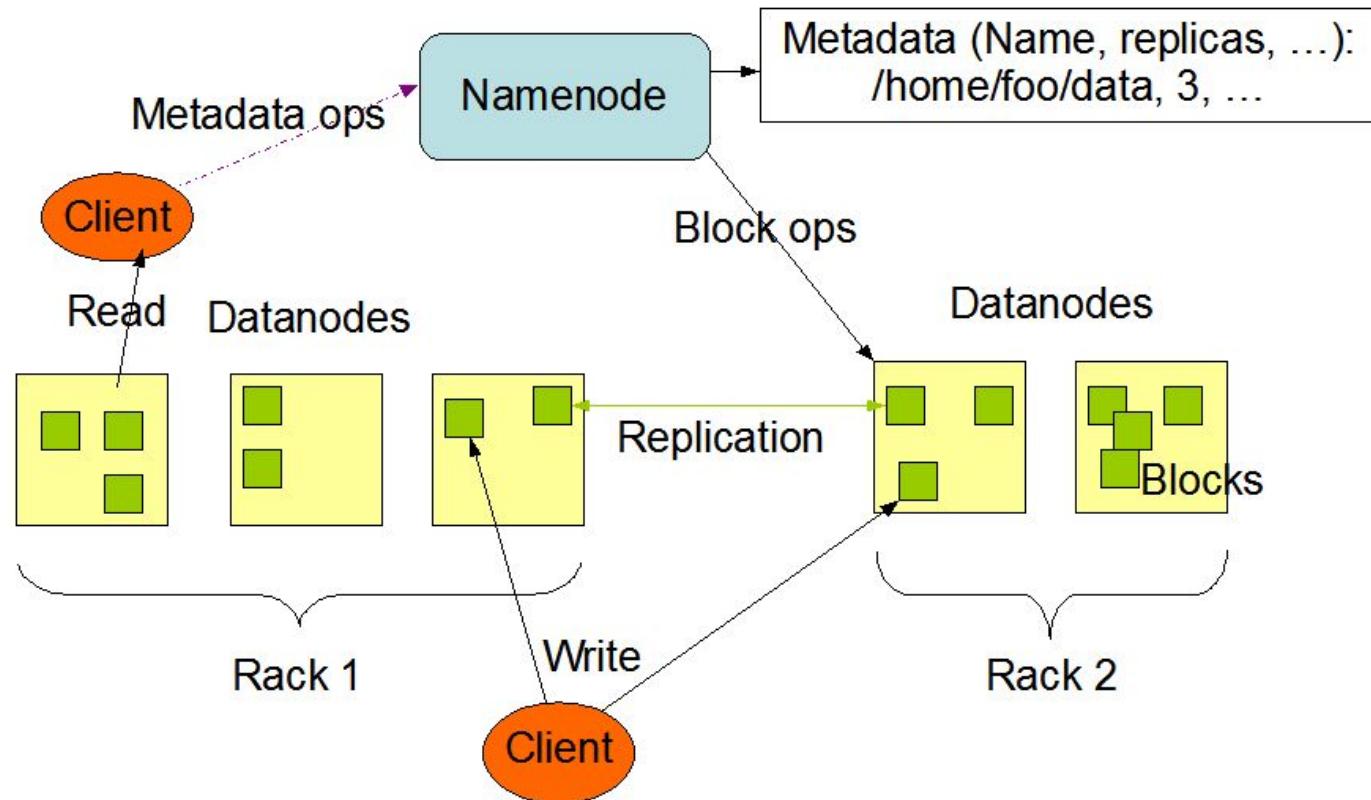


History

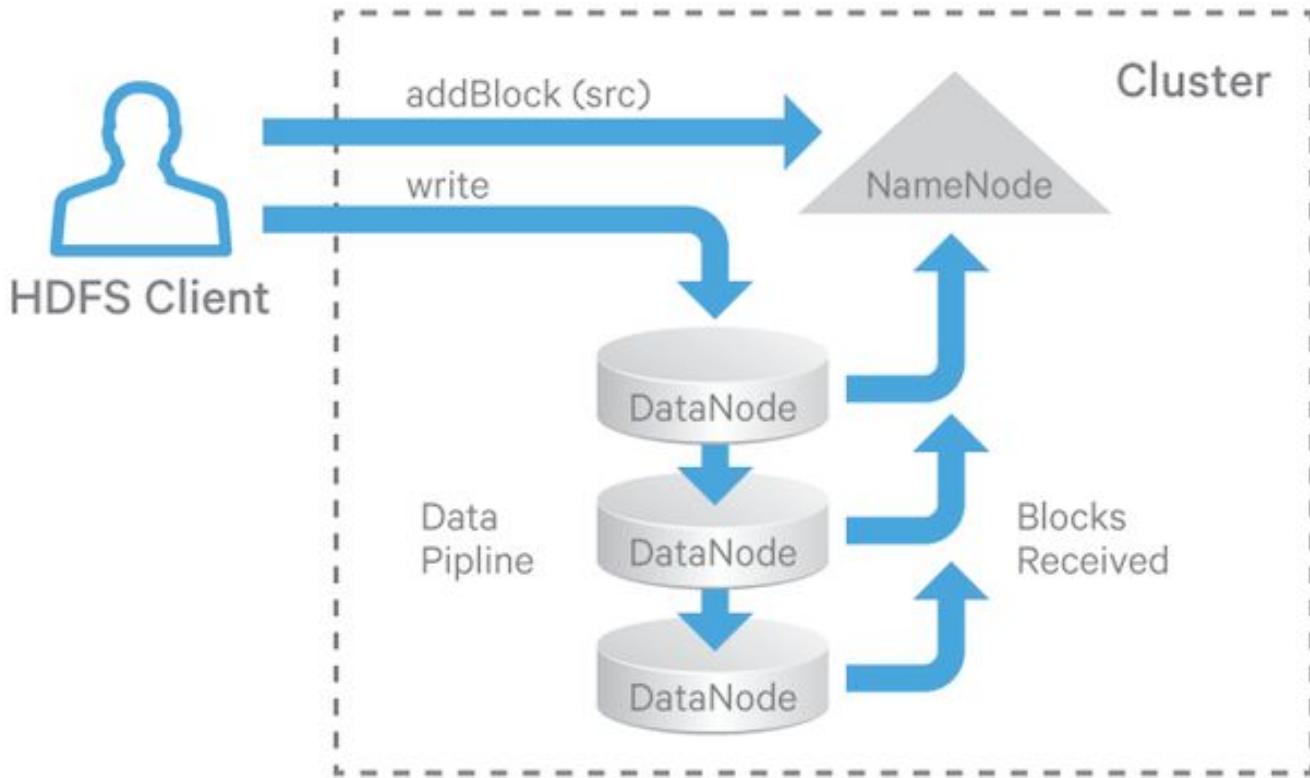
- ▶ 2003 Google File System
- ▶ 2004 MapReduce
- ▶ 2006 Hadoop
- ▶ 2007 HBase
- ▶ 2008 Hadoop wins TeraSort contest
- ▶ 2009 Spark
- ▶ 2009 Hadoop sorts one petabyte
- ▶ 2010 Hive

Hadoop File System

- ▶ Open source version of GFS
- ▶ Very similar architecture
- ▶ Deployments of 1,000s of nodes of HDFS exist



HDFS Architecture



HDFS Replication

Replication

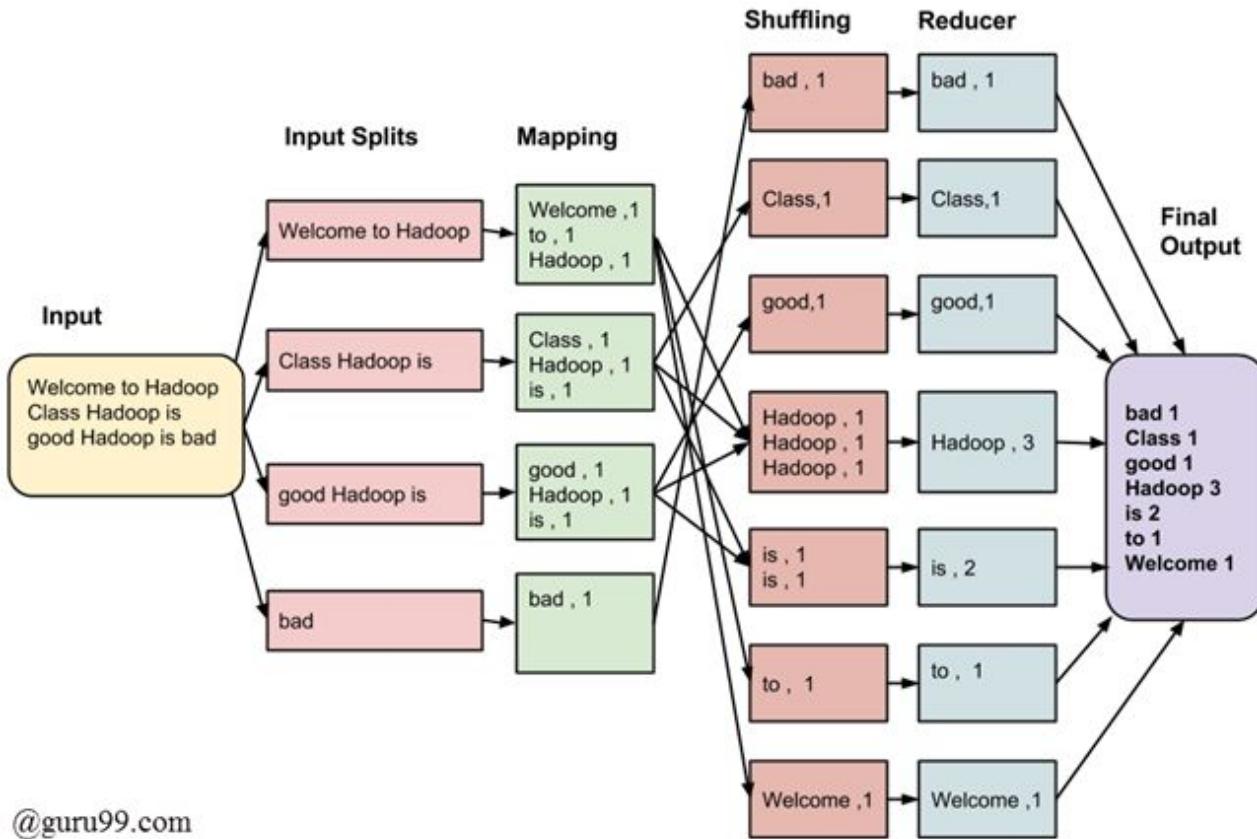
- ▶ Multiple copies are kept in case of failure so data isn't lost
- ▶ Can also makes reads faster since we have multiple copies to choose from
- ▶ Data is automatically split across multiple racks in a data center

MapReduce

- ▶ A programming model which consists of writing *map* and *reduce* functions
- ▶ Map accepts key/value pairs and produce a sequence of K/V pairs
- ▶ Data is shuffled to group keys together
- ▶ Reduce accepts values with the same key and produce a new K/V pair

MapReduce Execution

- ▶ Map tasks are assigned to machines based on the input data
- ▶ Map tasks produce their output
- ▶ Mapper output is shuffled and sorted
- ▶ Reduce tasks are scheduled and run
- ▶ Reduce output is stored to disk





```
import sys

# input comes from STDIN (standard input)
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()
    # split the line into words
    words = line.split()
    # increase counters
    for word in words:
        # write the results to STDOUT (standard output);
        # what we output here will be the input for the
        # Reduce step, i.e. the input for reducer.py
        #
        # tab-delimited; the trivial word count is 1
        print '%s\t%s' % (word, 1)
```

Map function



```
import sys

current_word = None
current_count = 0
word = None

# input comes from STDIN
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()

    # parse the input we got from mapper.py
    word, count = line.split('\t', 1)

    # (continued on next slide)
```

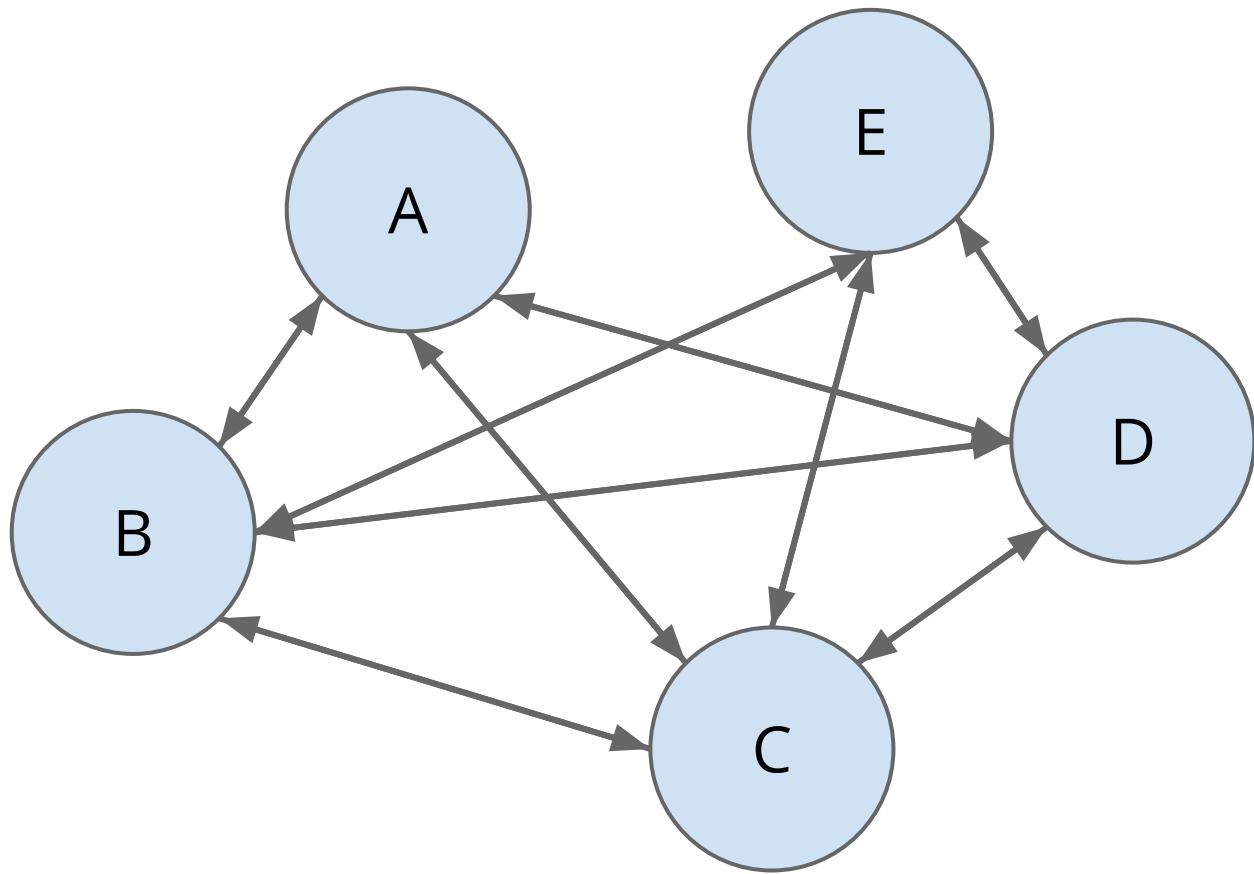
Reduce function (1/2)



```
# (continued within loop)

if current_word == word:
    current_count += int(count)
else:
    if current_word:
        print('%s\t%s' % (current_word, current_count))
    current_word = word
    current_count = count
```

Reduce function (2/2)



$A \rightarrow B C D$
 $B \rightarrow A C D E$
 $C \rightarrow A B D E$
 $D \rightarrow A B C E$
 $E \rightarrow B C D$

Mutual friends



$A \rightarrow B C D$

AB	BCD
AC	BCD
AD	BCD

$C \rightarrow A B D E$

CA	ABDE
CB	ABDE
CD	ABDE
CE	ABDE

$E \rightarrow B C D$

EB	BCD
EC	BCD
ED	BCD

$B \rightarrow A C D E$

BA	ACDE
BC	ACDE
BD	ACDE
BE	ACDE

$D \rightarrow A B C E$

DA	ABCE
DB	ABCE
DC	ABCE
DE	ABCE



$A \rightarrow B C D$

$C \rightarrow A B D E$

$E \rightarrow B C D$

AB BCD

AC ABDE

BE BCD

AC BCD

BC ABDE

CE BCD

AD BCD

CD ABDE

DE BCD

CE ABDE

$B \rightarrow A C D E$

$D \rightarrow A B C E$

AB ACDE

AD ABCE

BC ACDE

BD ABCE

BD ACDE

CD ABCE

BE ACDE

DE ABCE

Map



AB	ACDE
AB	BCD

BD	ABCE
BD	ACDE

DE	ABCE
DE	BCD

AC	ABDE
AC	BCD

BE	ACDE
BE	BCD

AD	ABCE
AD	BCD

CD	ABDE
CD	ABCE

BC	ABDE
BC	ACDE

CE	ABDE
CE	BCD

Sorted Map Output

A B C D

B D A C E

D E B C

A C B D

B E C D

A D B C

C D A B E

B C A D E

C E B D

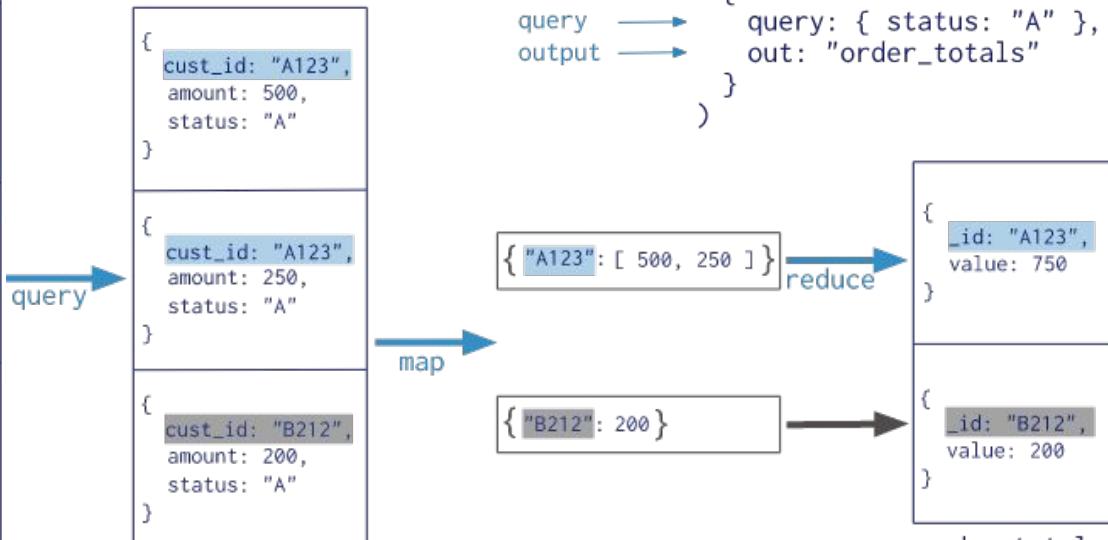
Reduce



```
{  
  cust_id: "A123",  
  amount: 500,  
  status: "A"  
}  
  
{  
  cust_id: "A123",  
  amount: 250,  
  status: "A"  
}  
  
{  
  cust_id: "B212",  
  amount: 200,  
  status: "A"  
}  
  
{  
  cust_id: "A123",  
  amount: 300,  
  status: "D"  
}
```

orders

Collection
db.orders.mapReduce(
 map → function() { emit(this.cust_id, this.amount); },
 reduce → function(key, values) { return Array.sum(values) },
 {
 query → { status: "A" },
 output → { out: "order_totals" }
 }
)



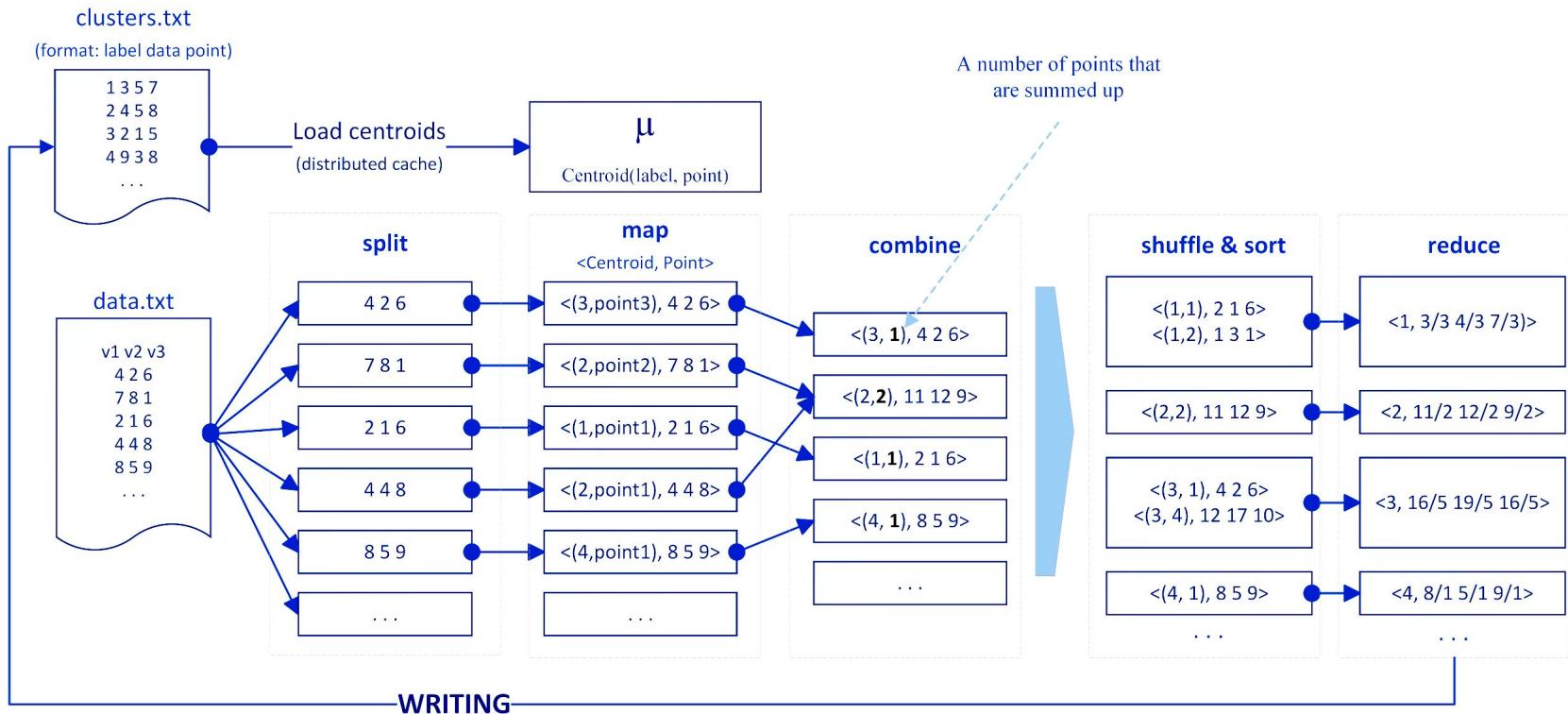
MapReduce in MongoDB

MapReduce vs Aggregation

- ▶ Aggregation pipelines are often easier to read than MapReduce programs
- ▶ MapReduce in MongoDB is generally less efficient and can't be optimized
- ▶ Collections in MongoDB can be “sharded” or split so MapReduce can happen in parallel as in Hadoop

Iterative MapReduce

- ▶ Sometimes an algorithm cannot be executed in a single MapReduce job
- ▶ MapReduce jobs can be chained together to provide a solution
- ▶ Often happens with algorithms which iterate until convergence (e.g. k-means, PageRank, etc.)



MapReduce k-Means

Iterative MapReduce

- ▶ Reduce output must be read from disk again with each new job
- ▶ In some cases, the input data is read from disk many times
- ▶ There is no easy way to share the work done between iterations

Apache Spark

- ▶ Computation model is much richer than just MapReduce
- ▶ Transformations on input data can be written lazily and batched together
- ▶ Intermediate results can be cached and reused in future calculations

Spark Model

- ▶ Series of lazy **transformations** which are followed by **actions** that force evaluation of all transformations
- ▶ Each step produces a resilient distributed dataset (RDD)
- ▶ Intermediate results can be cached on memory or disk, optionally serialized

RDD Lineage

- ▶ For each RDD, we keep a *lineage*, the operations which created it
- ▶ Spark can then recompute any data which is lost without storing to disk
- ▶ We can still decide to keep a copy in memory or on disk for performance

Spark Model

Example

```
val in = spark.read.format("csv")
    .option("header", "true").csv("data.csv")
val over = in.filter($"y" > $"x").cache

val dist = over.select(hypot($"x", $"y"))
dist.write.csv("dist.csv")

val angle = over.select(atan($"y" / $"x"))
angle.write.csv("angle.csv")
```

Spark Model

Transformations

```
val in = spark.read.format("csv")
    .option("header", "true").csv("data.csv")
val over = in.filter($"y" > $"x").cache
```

```
val dist = over.select(hypot($"x", $"y"))
dist.write.csv("dist.csv")
```

```
val angle = over.select(atan($"y" / $"x"))
angle.write.csv("angle.csv")
```

Spark Model

Actions

```
val in = spark.read.format("csv")
    .option("header", "true").csv("data.csv")
val over = in.filter($"y" > $"x").cache

val dist = over.select(hypot($"x", $"y"))
dist.write.csv("dist.csv")

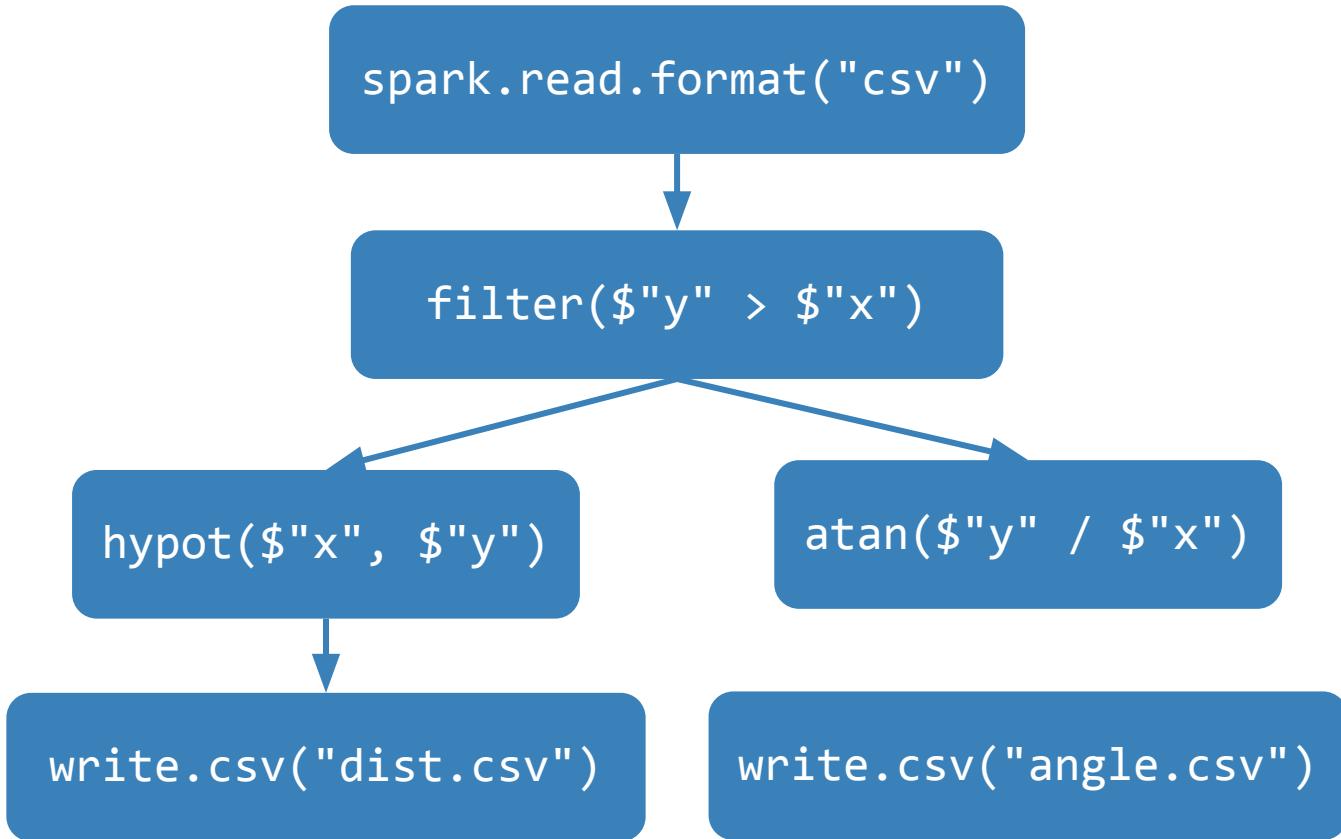
val angle = over.select(atan($"y" / $"x"))
angle.write.csv("angle.csv")
```

Spark Transforms

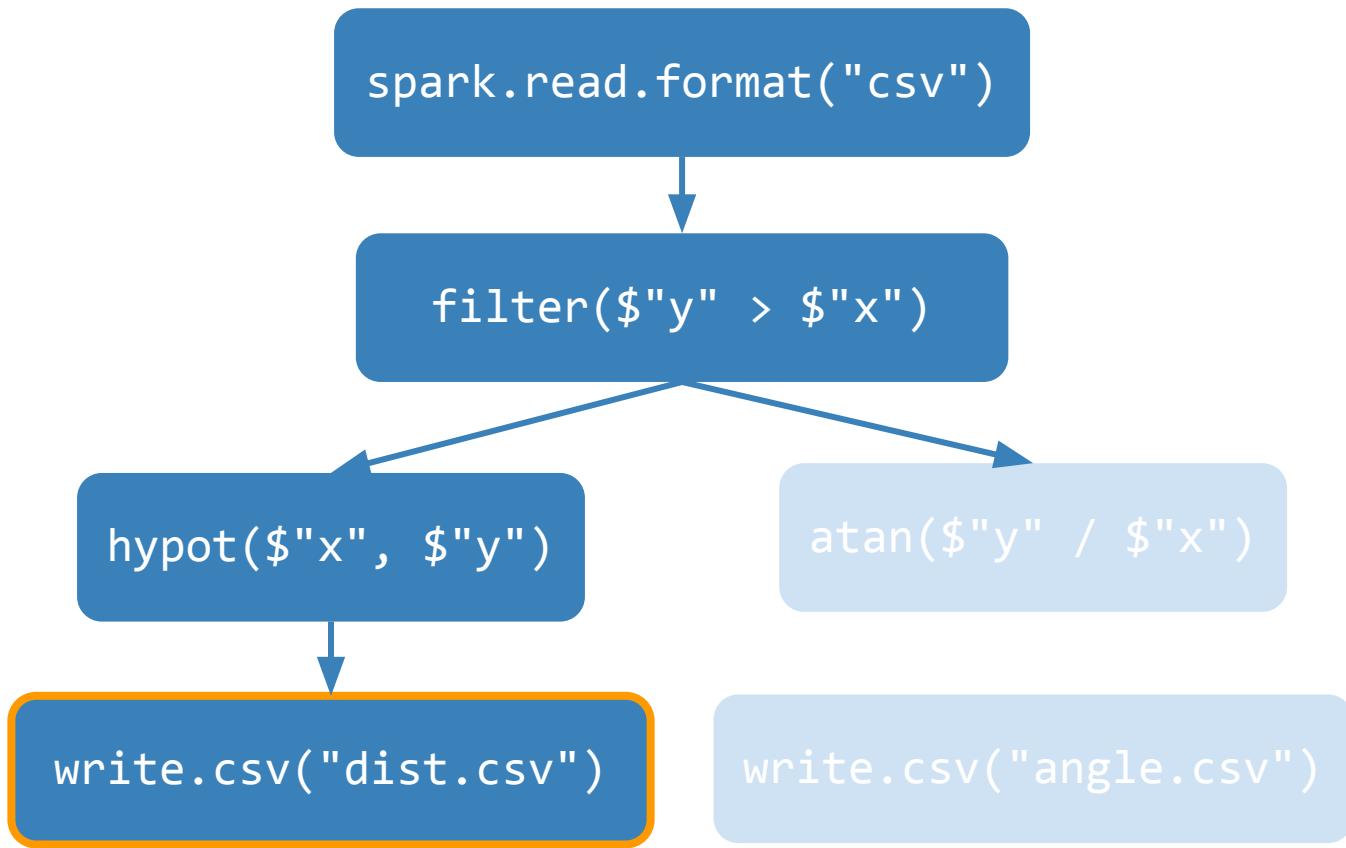
- ▶ `distinct()`
- ▶ `filter(fn)`
- ▶ `intersection(other)`
- ▶ `join(other)`
- ▶ `map(fn)`
- ▶ `union(other)`

Spark Actions

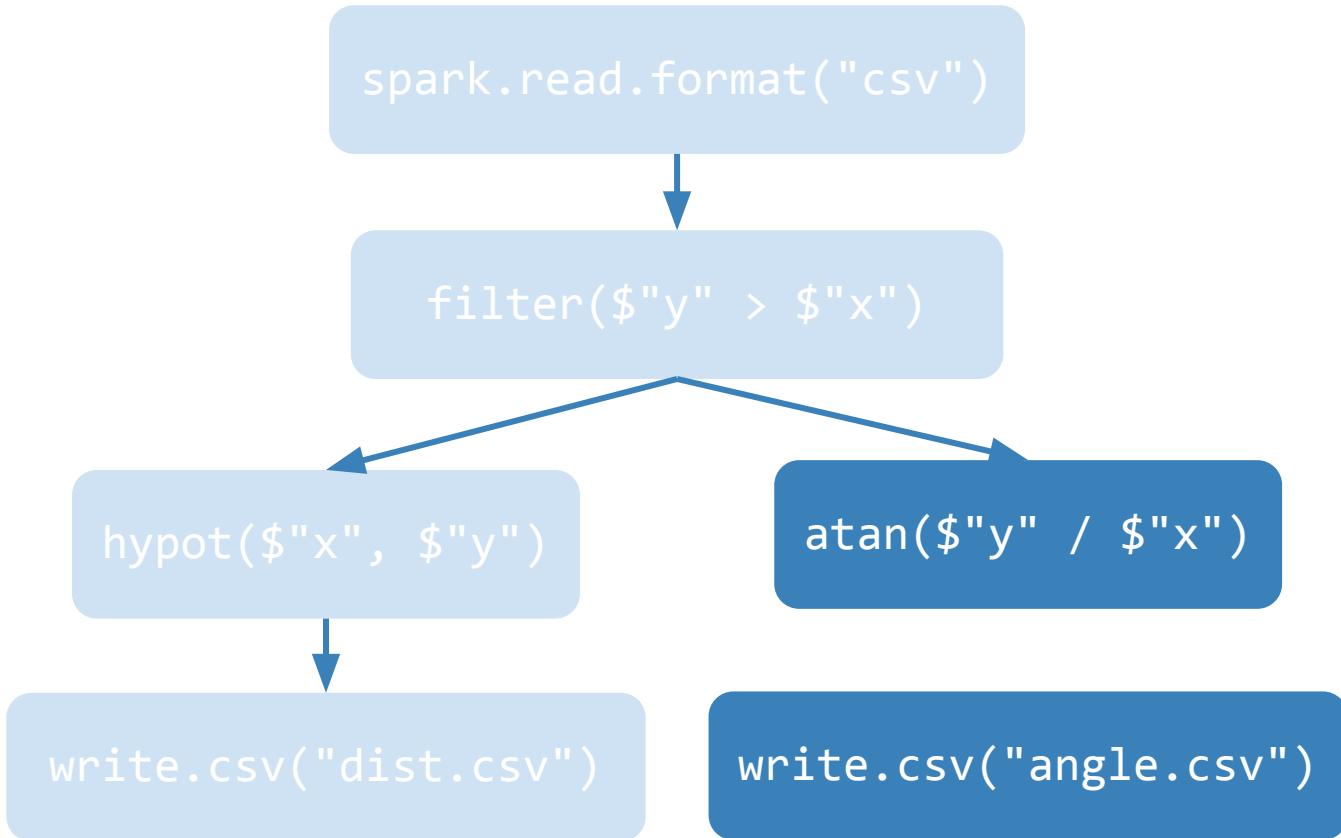
- ▶ collect()
- ▶ count()
- ▶ first()
- ▶ take(n)
- ▶ reduce(fn)
- ▶ foreach(fn)
- ▶ saveAsTextFile()



Task DAG



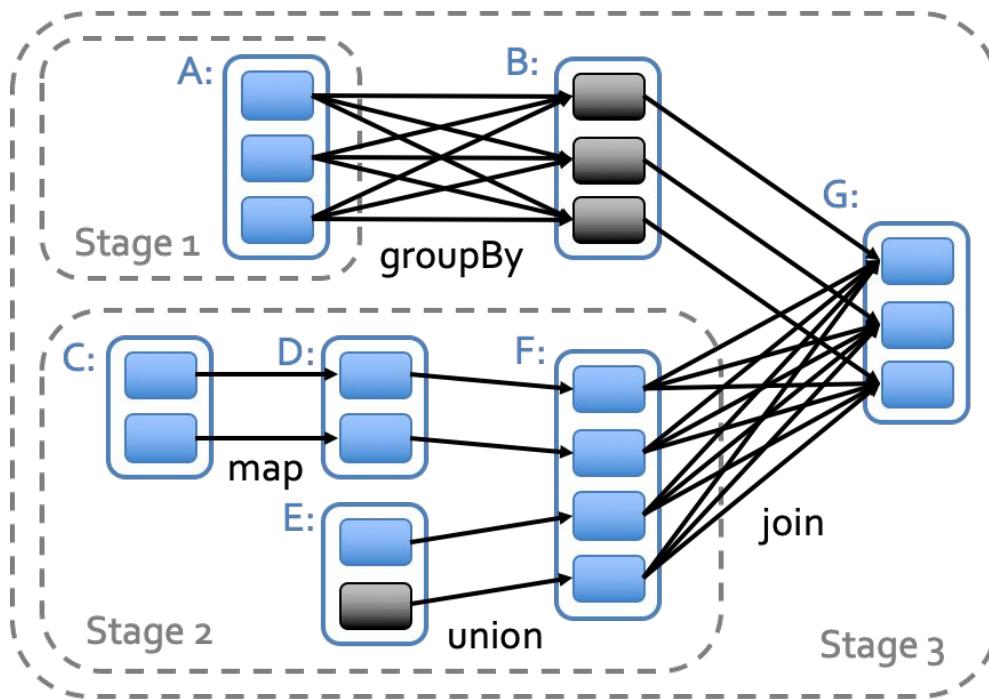
Task DAG



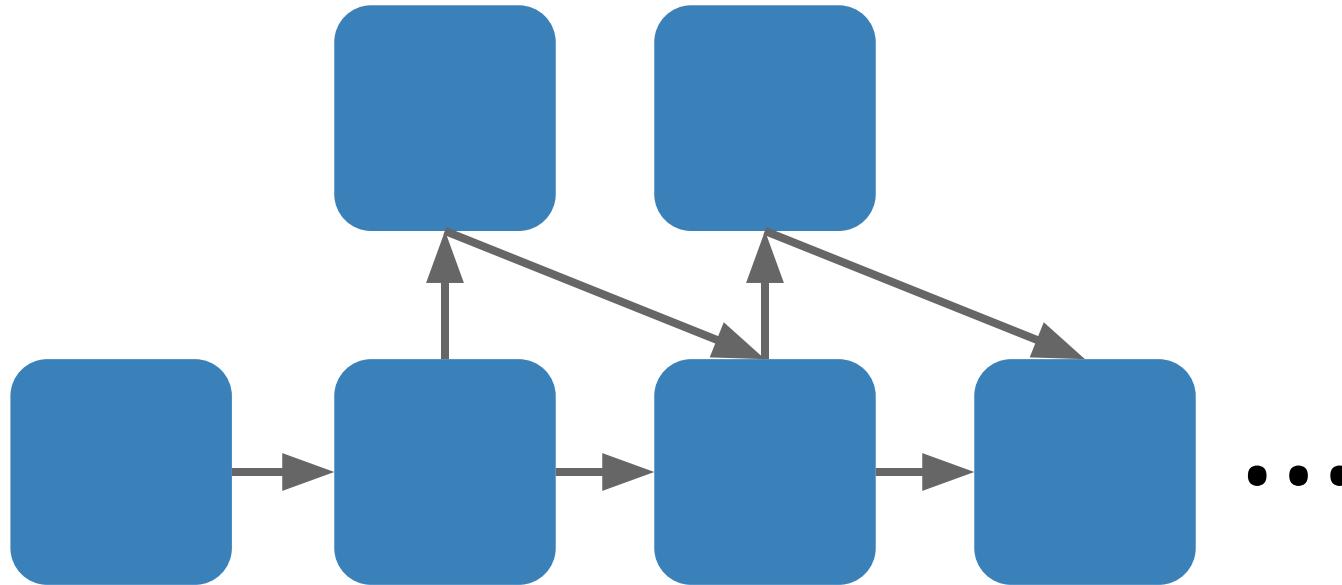
Task DAG

Spark Stages

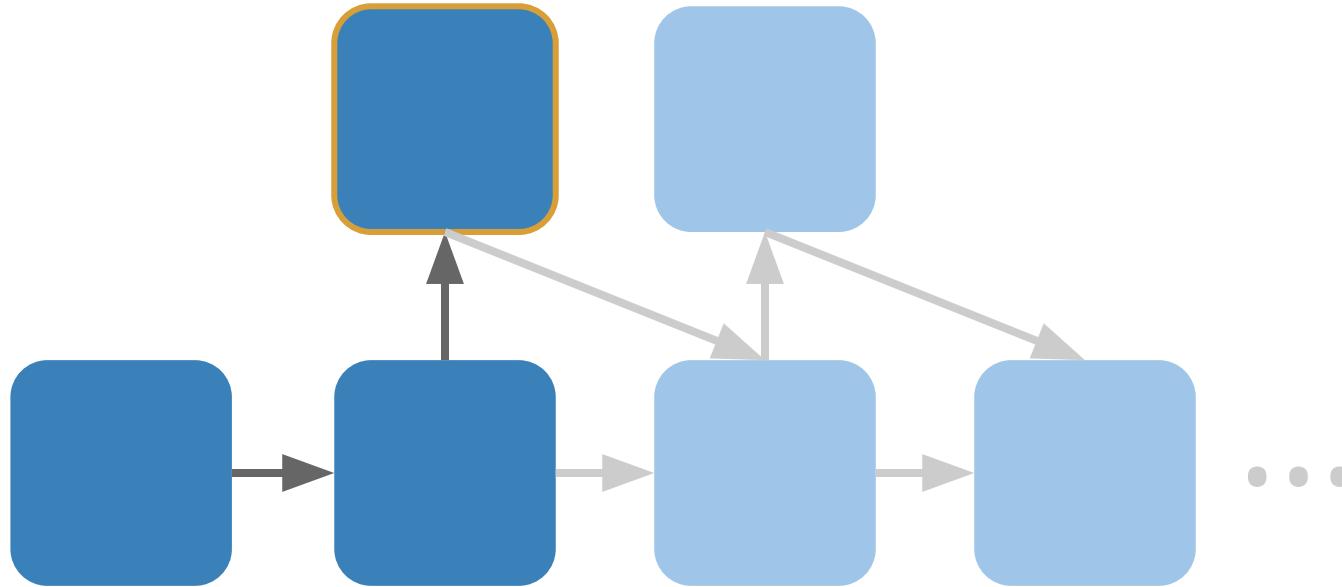
- ▶ Some operations such as join, sort and grouping, require shuffling
- ▶ This is similar to the sorting before the reduce phase in MapReduce
- ▶ Execution is grouped into *stages* depending on where shuffles occur
- ▶ Tasks execute on partitions in parallel



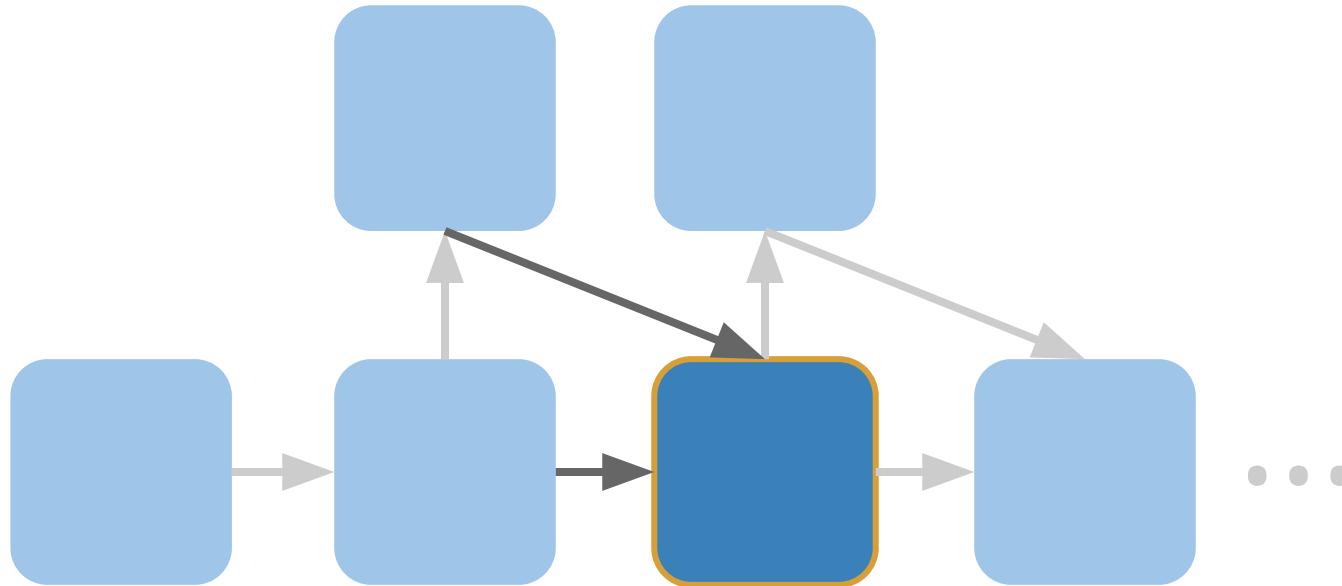
Spark Scheduler



Spark Iteration



Spark Iteration



Spark Iteration

Spark k-means

1. Load data and “persist”
2. Randomly sample for initial clusters
3. Broadcast the initial clusters
4. Loop over data and find new centers
5. Repeat from 3 until converged

Spark k-means

- ▶ Unlike MapReduce, we can keep the input in memory and load them once
- ▶ Broadcasting means we can quickly send the centers to all machines
- ▶ All cluster assignments do not need to be written to disk every time



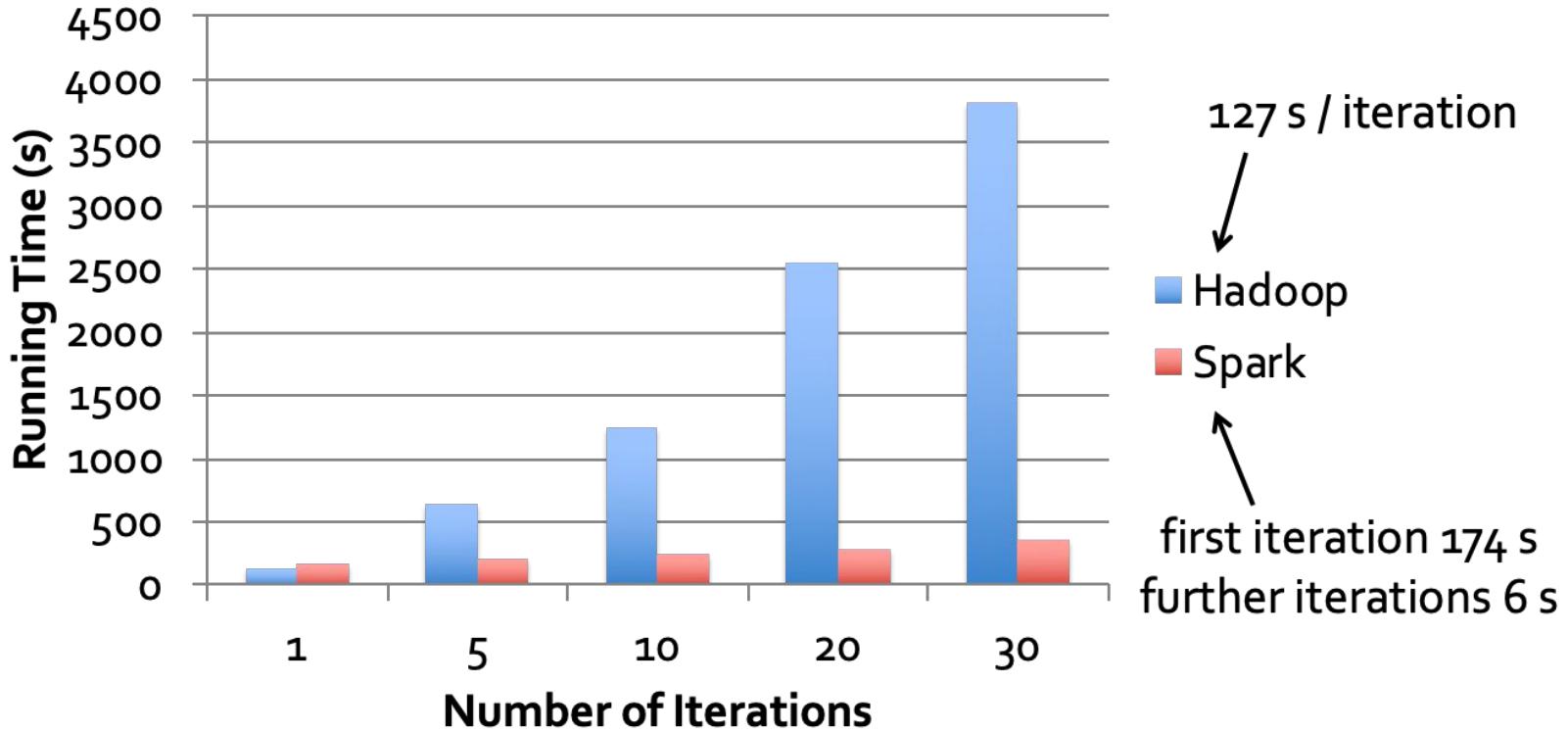
```
val data = spark.textFile(...).map(readPoint).cache()

var w = Vector.random(D)

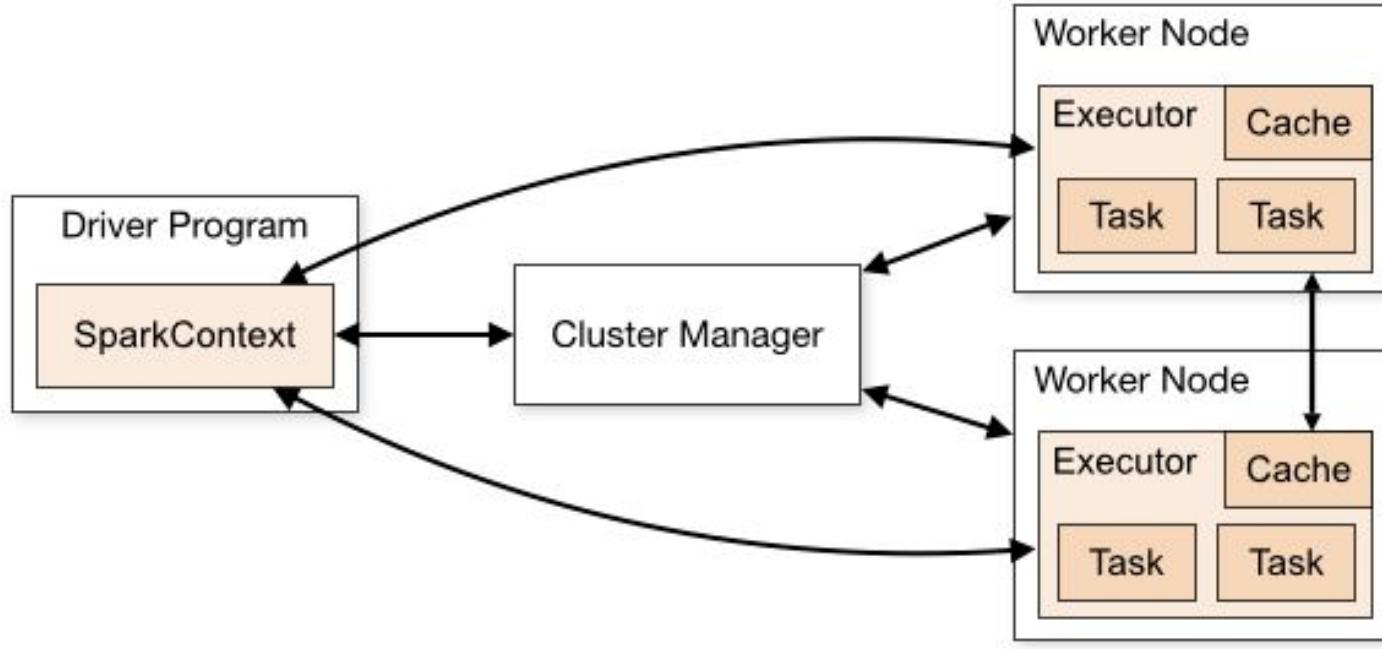
for (i <- 1 to ITERATIONS) {
    val gradient = data.map(p =>
        (1 / (1 + exp(-p.y*(w dot p.x))) - 1) * p.y * p.x
    ).reduce(_ + _)
    w -= gradient
}

println("Final w: " + w)
```

Logistic Regression



Logistic Regression



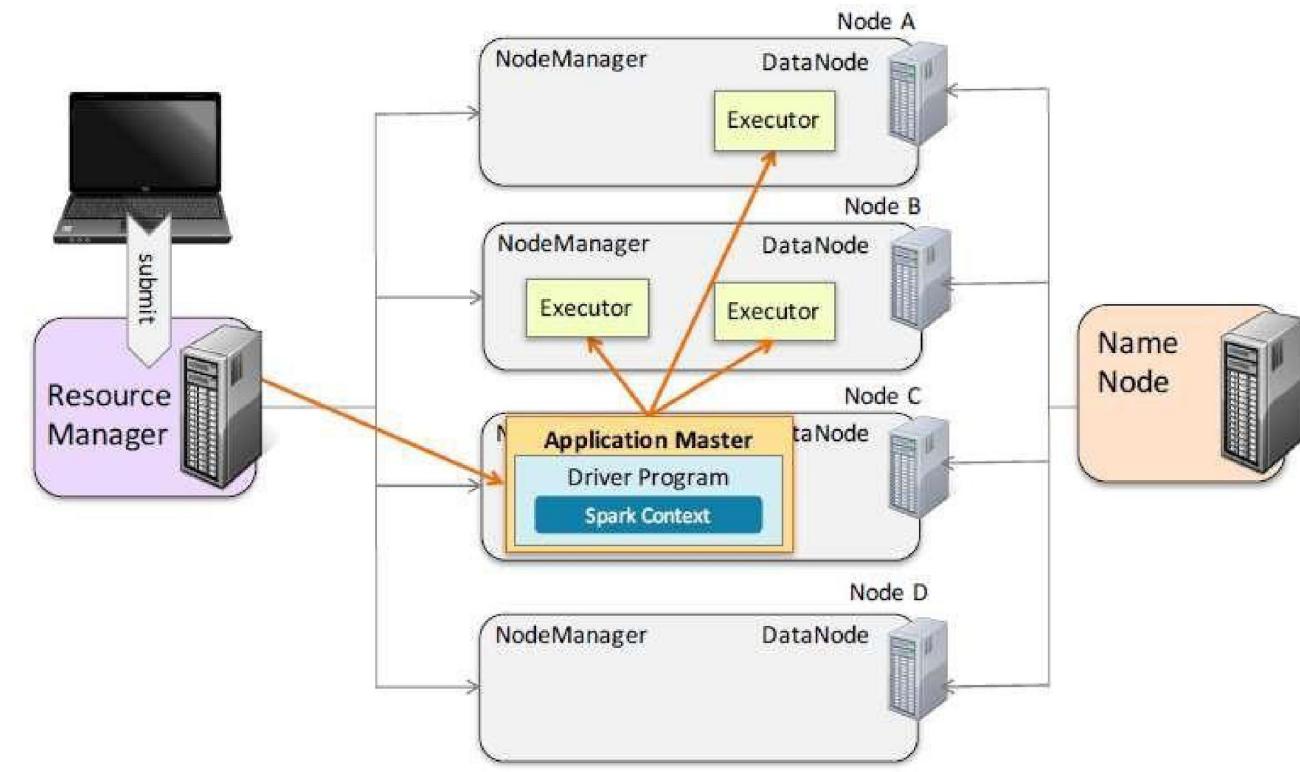
Spark Architecture

Spark Driver

- ▶ Runs the application to define RDDs and submit jobs to run on the cluster
- ▶ Creates a Spark *context* which represents the execution environment
- ▶ Tracks the progress of all jobs and displays a web interface

Spark Cluster Manager

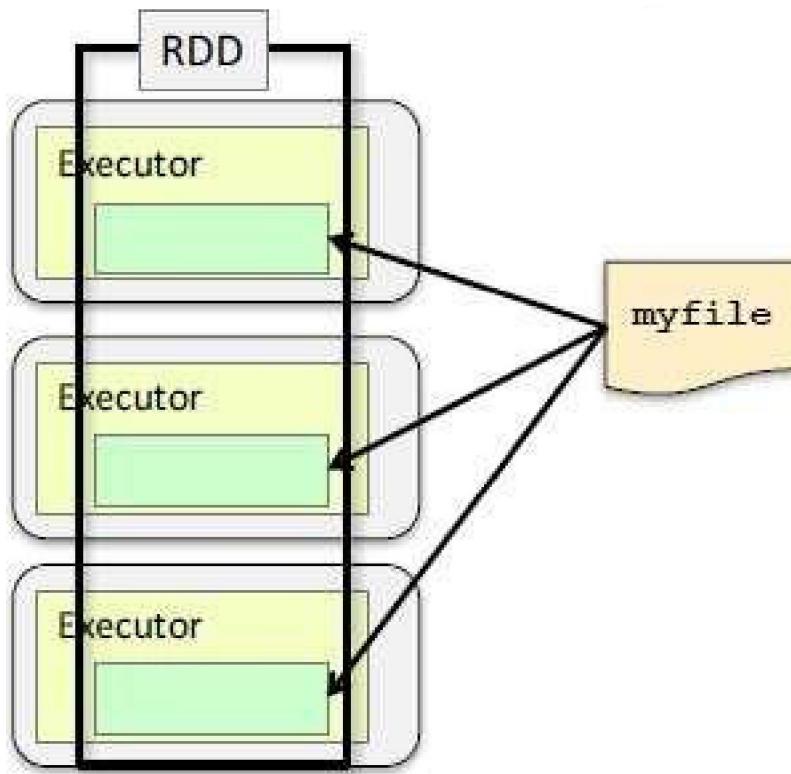
- ▶ Multiple cluster managers can be used
 - ▷ Standalone - single machine
 - ▷ Apache Mesos
 - ▷ Hadoop YARN
 - ▷ Kubernetes
 - ▷ (and others)
- ▶ Handles creating workers to run each distributed Spark job



Spark on YARN

Spark Executors

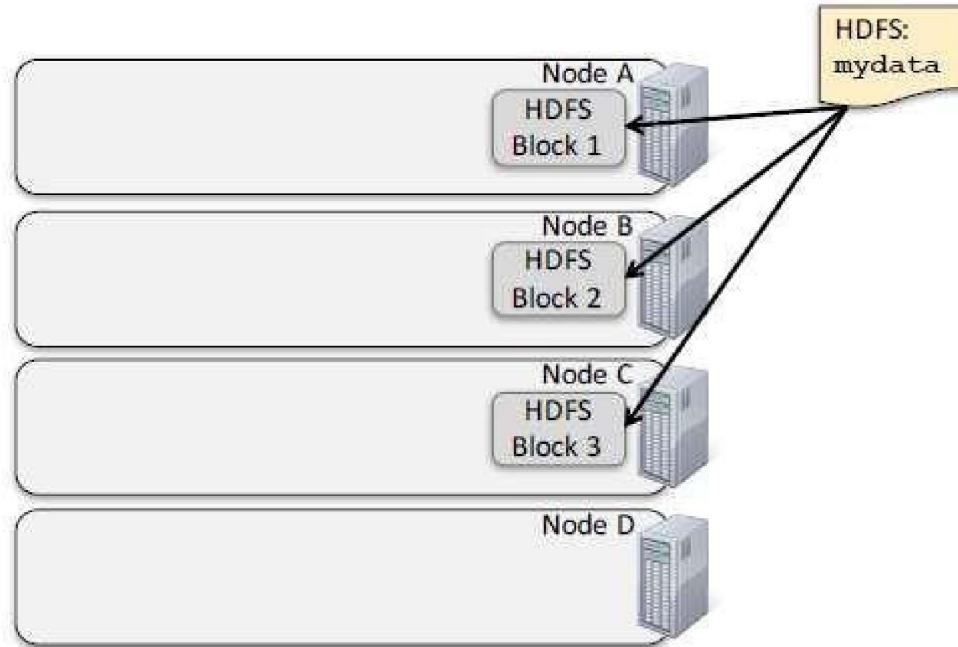
- ▶ Executors run across multiple machines in a Spark cluster
- ▶ Receives a series of transformations to execute on RDDs from the drivers
- ▶ If an executor fails, the driver can ask the cluster manager to start another



Parallelism in Spark



```
$ hdfs dfs -put mydata
```

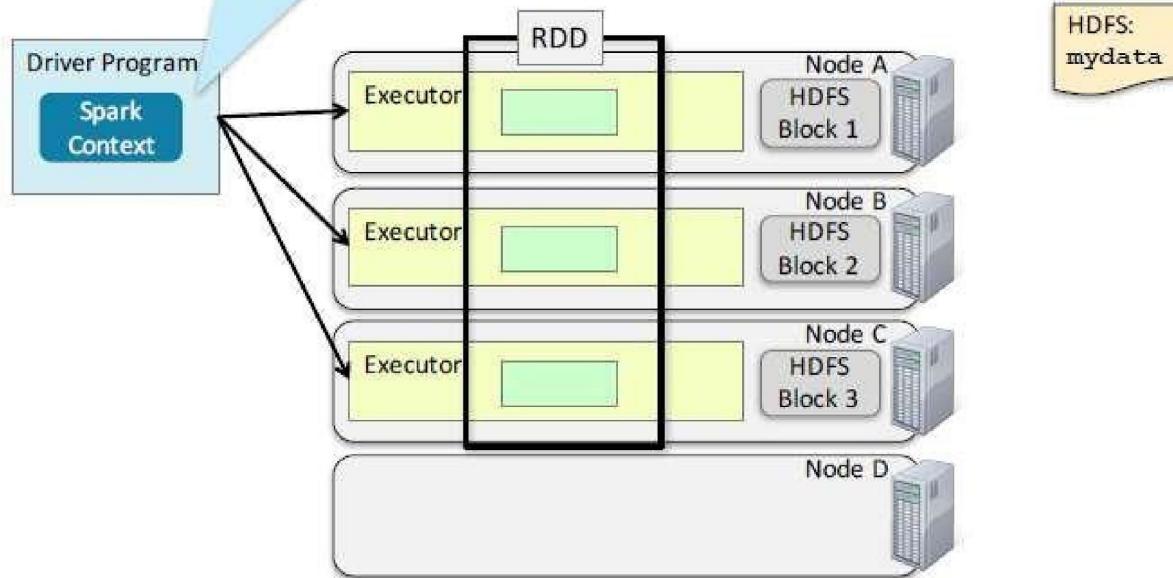


Parallelism in Spark



```
sc.textFile("hdfs://...mydata").collect()
```

By default, Spark partitions file-based RDDs by block. Each block loads into a single partition.



Parallelism in Spark



File: purplecow.txt

I've never seen a purple cow.
I never hope to see one;
But I can tell you, anyhow,
I'd rather see than be one.

Language: Scala

View
Edit
Run
Help

Spark Execution



```
> val mydata = sc.textFile("purplecow.txt")
```

Language: Scala

File: purplecow.txt

I've never seen a purple cow,
I never hope to see one;
But I can tell you, anyhow,
I'd rather see than be one.

RDD: mydata

Spark Execution



File: purplecow.txt

I've never seen a purple cow.
I never hope to see one;
But I can tell you, anyhow,
I'd rather see than be one.

Language: Scala

```
> val mydata = sc.textFile("purplecow.txt")
> val mydata_uc = mydata.map(line =>
  line.toUpperCase())
```

RDD: mydata

RDD: mydata_uc

Spark Execution



File: purplecow.txt

I've never seen a purple cow.
I never hope to see one;
But I can tell you, anyhow,
I'd rather see than be one.

Language: Scala

```
> val mydata = sc.textFile("purplecow.txt")
> val mydata_uc = mydata.map(line =>
line.toUpperCase())
> val mydata_filt = mydata_uc.filter(line
=> line.startsWith("I"))
```

RDD: mydata

RDD: mydata_uc

RDD: mydata_filt

Spark Execution



File: purplecow.txt

I've never seen a purple cow.
I never hope to see one;
But I can tell you, anyhow,
I'd rather see than be one.

Language: Scala

```
> val mydata = sc.textFile("purplecow.txt")
> val mydata_uc = mydata.map(line =>
  line.toUpperCase())
> val mydata_filt = mydata_uc.filter(line
  => line.startsWith("I"))
> mydata_filt.count()
3
```

RDD: mydata

I've never seen a purple cow.
I never hope to see one;
But I can tell you, anyhow,
I'd rather see than be one.

RDD: mydata_uc

I'VE NEVER SEEN A PURPLE COW.
I NEVER HOPE TO SEE ONE;
BUT I CAN TELL YOU, ANYHOW,
I'D RATHER SEE THAN BE ONE.

RDD: mydata_filt

I'VE NEVER SEEN A PURPLE COW.
I NEVER HOPE TO SEE ONE;
I'D RATHER SEE THAN BE ONE.

Spark Execution



rdd1	rdd2
Chicago	San Francisco
Boston	Boston
Paris	Amsterdam
San Francisco	Mumbai
Tokyo	McMurdo Station

`rdd1.subtract(rdd2)`



Tokyo
Paris
Chicago

`rdd1.zip(rdd2)`



(Chicago, San Francisco)
(Boston, Boston)
(Paris, Amsterdam)
(San Francisco, Mumbai)
(Tokyo, McMurdo Station)



rdd1
Chicago
Boston
Paris
San Francisco
Tokyo

rdd2
San Francisco
Boston
Amsterdam
Mumbai
McMurdo Station

`rdd1.union(rdd2)`



Chicago
Boston
Paris
San Francisco
Tokyo
San Francisco
Boston
Amsterdam
Mumbai
McMurdo Station

`rdd1.intersection(rdd2)`



Boston
San Francisco

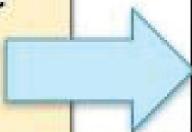
Multi-RDD Operations

(key1 , value1)
(key2 , value2)
(key3 , value3)
...

Pair RDDs

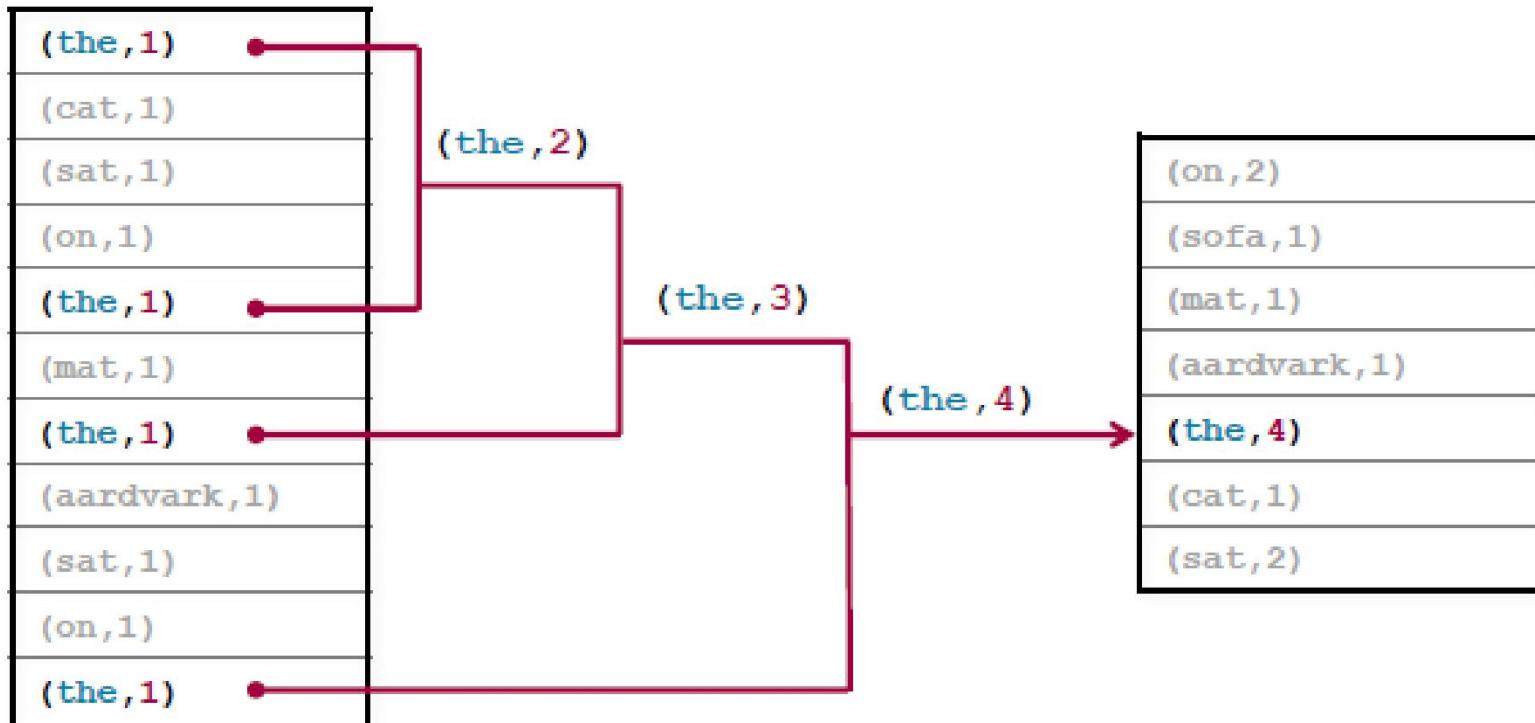
```
> val users = sc.textFile(file).  
  map(line => line.split('\t')).  
  map(fields => (fields(0), fields(1)))
```

user001\tFred Flintstone
user090\tBugs Bunny
user111\tHarry Potter
...

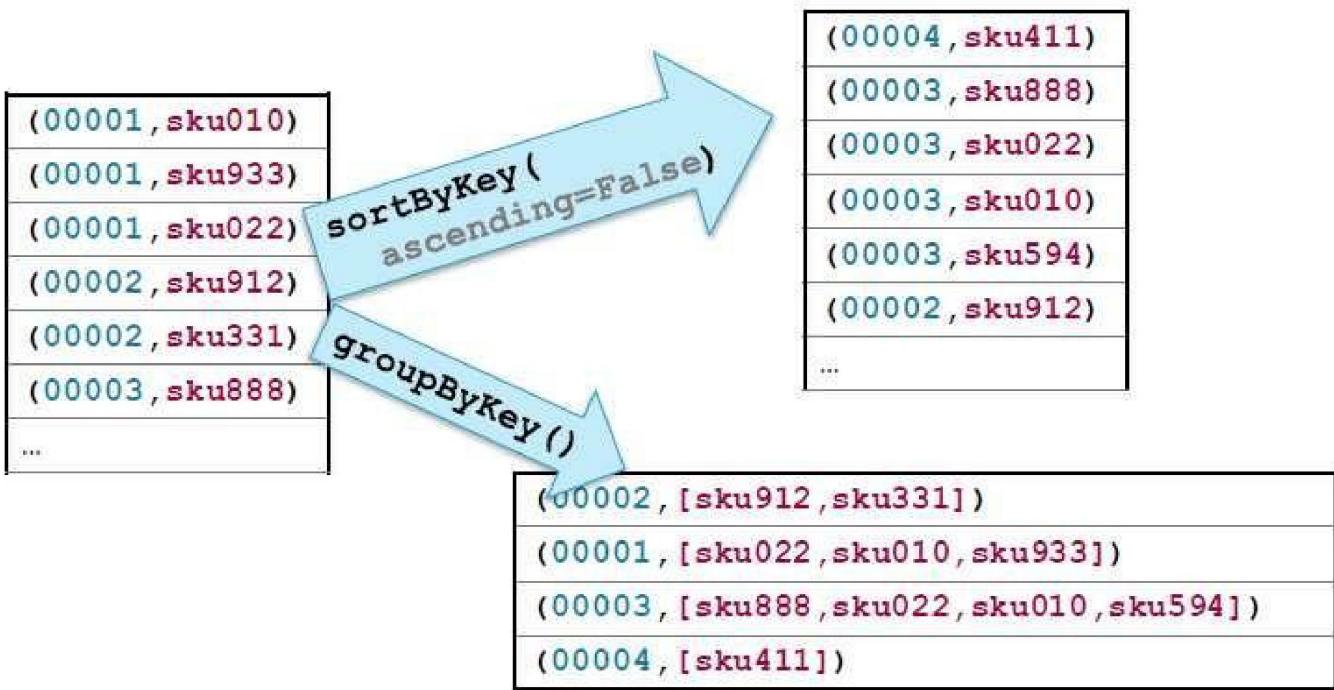


(user001, Fred Flintstone)
(user090, Bugs Bunny)
(user111, Harry Potter)
...

Pair RDDs

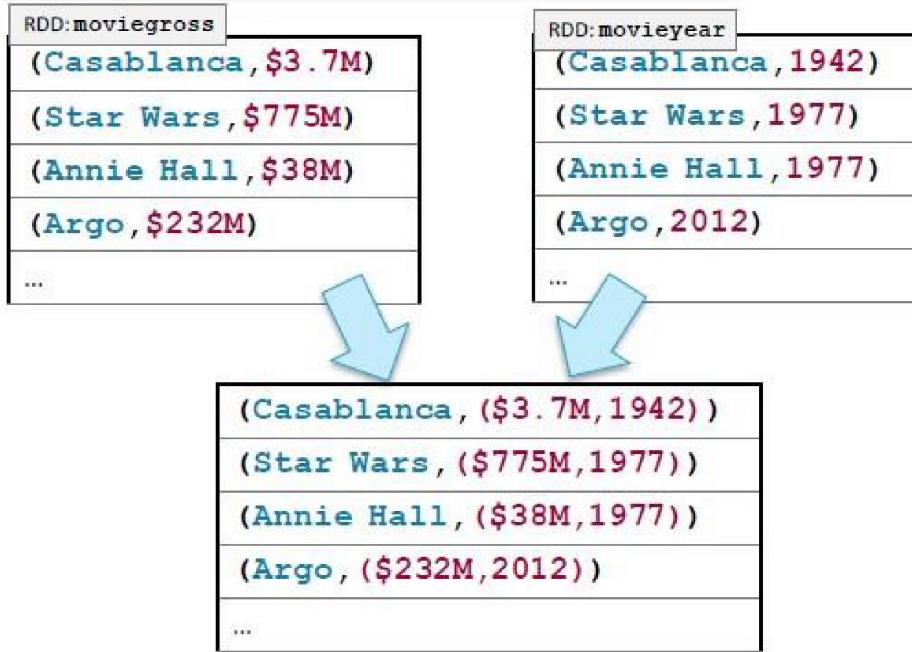


reduceByKey



Sorting and grouping

```
> movies = moviegross.join(movieyear)
```



Joining PairRDDs