

# Databases and Big Data

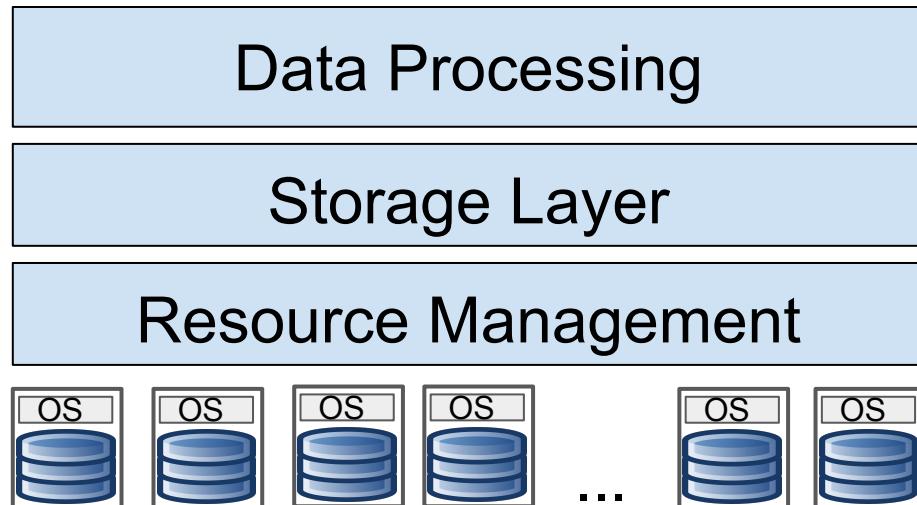
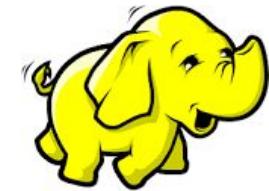
Spark

# Agenda

- Why Spark?
- Spark APIs
- Examples

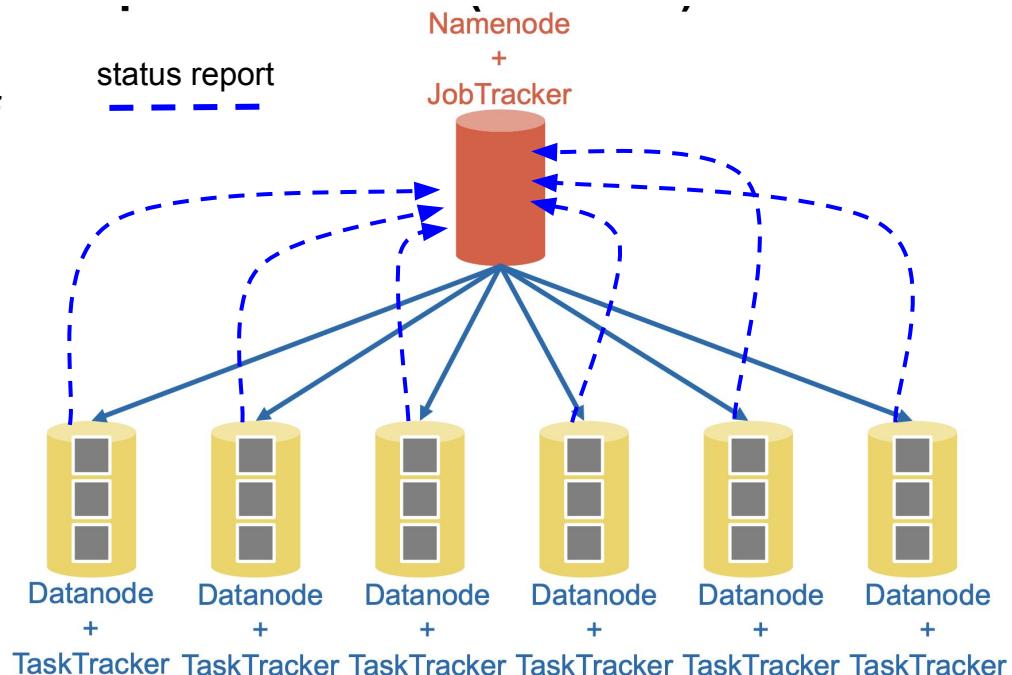
# Recap

- Hadoop stack
  - HDFS: distributed file system
  - MapReduce: data processing
- We haven't seen
  - How MapReduce handles failures



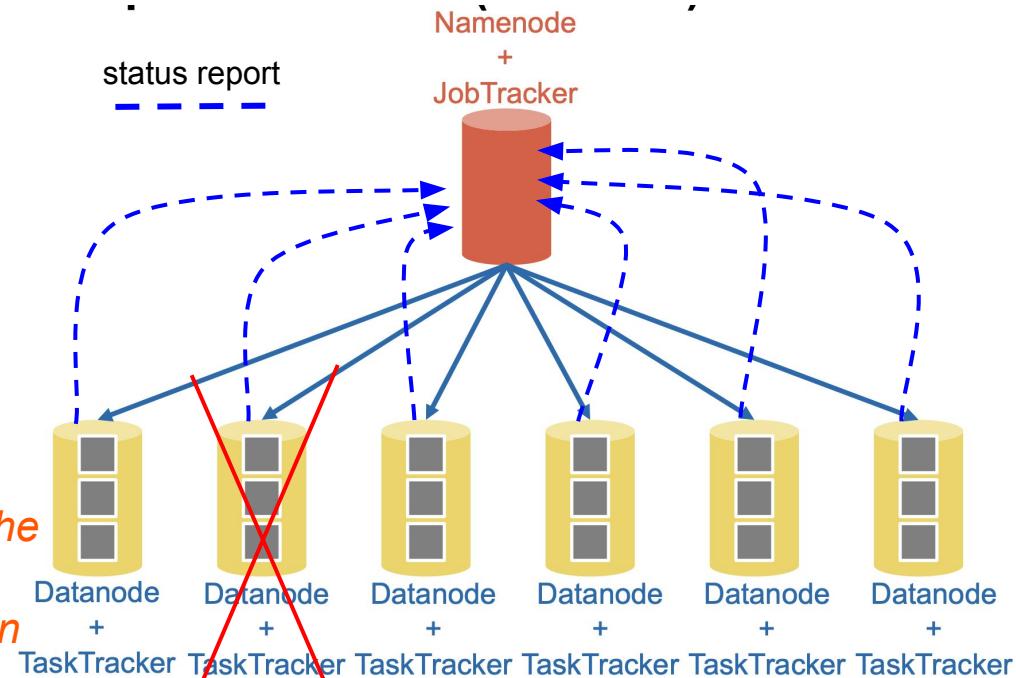
# Recap

- Architecture
  - JobTracker tracks progress of map/reduce tasks
  - TaskTracker reports progress



# MapReduce Fault Tolerance

- When a node fails:
  - Both DataNode and TaskTracker fail
- How did HDFS handle failures?
- JobTracker starts another tasks
  - Why possible?



*When reboot, Jobtracker needs to recover all the hdfs blocks. Namenode keep tracks of all the metadata on the node and is the master, so can coordinate the recovering.*

*If some node fail, just start a new machine and run*

# Problems with MapReduce

Not this!

- Probably just sour grapes
  - Remember how DB people reacted to MongoDB?

[The Database Column](#)

A multi-author blog on database technology and innovation.

## MapReduce: A major step backwards

By David DeWitt on January 17, 2008 4:20 PM | [Permalink](#) | [Comments \(44\)](#) | [TrackBacks \(1\)](#)

*[Note: Although the system attributes this post to a single author, it was written by David J. DeWitt and Michael Stonebraker.]*

On January 8, a Database Column reader asked for our views on new distributed database research efforts, and we'll begin here with our views on [MapReduce](#). revolution of so-called "cloud computing." This paradigm entails harnessing large numbers of (low-end) processors working in parallel to solve a computing problem "rather than utilizing a much smaller number of high-end servers.

# MapReduce

- To fulfill a real need at end of 1990s:
  - Search is the killer app
  - Need to create index on huge data (Web)
    - What index?
  - Data process is embarrassingly parallel
    - ETL

But then, needs changed

*Completely independent of each other.  
(suitable for map task, cause you can just throw more machines in and the task can run faster in parallel)*



# What Changed?

- Application trends in late 00s
  - **Iterative computation:**
    - Machine learning: SGD has multiple passes
  - **Interactive computation:**
    - Data science: exploratory analysis
      - Don't always know what you want in advance

**Both -> not supported by MapReduce framework**

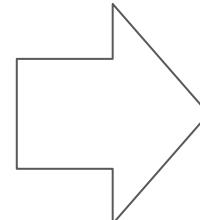
MapReduce designed for **batch computation**



# What Changed?

## Batch computation

- Do it once, or infrequently
- Throughput
- Examples:
  - Build index
  - Sort
  - ETL



## Iterative computation

- Do it many times
- Throughput
- Examples:
  - PageRank
  - kMeans
  - Graphs
  - SGD

## Interactive computation

- Find out what “it” is
- Ad-hoc queries
- Human in the loop
- Latency
- Example:
  - Ad-hoc queries
  - Exploratory data science

*TP: your system can handle that huge volume of data, main requirement*

*Latency: Fast queries (1-10s)*

# What Changed?

- Most applications **fit in memory**

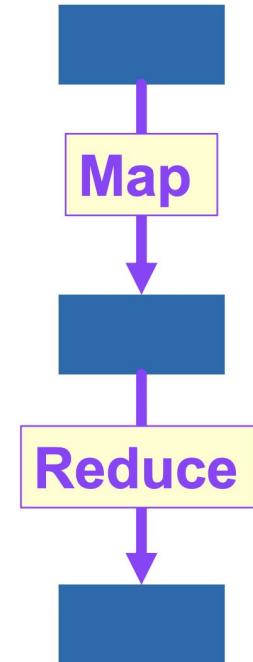


Memory (GB)	Facebook (% jobs)	Microsoft (% jobs)	Yahoo! (% jobs)
8	69	38	66
16	74	51	81
32	96	82	97.5
64	97	98	99.5
128	98.8	99.4	99.8
192	99.5	100	100
256	99.6	100	100

\*G Ananthanarayanan, A. Ghodsi, S. Shenker, I. Stoica, "Disk-Locality in Datacenter Computing Considered Irrelevant", HotO 2011

# MapReduce

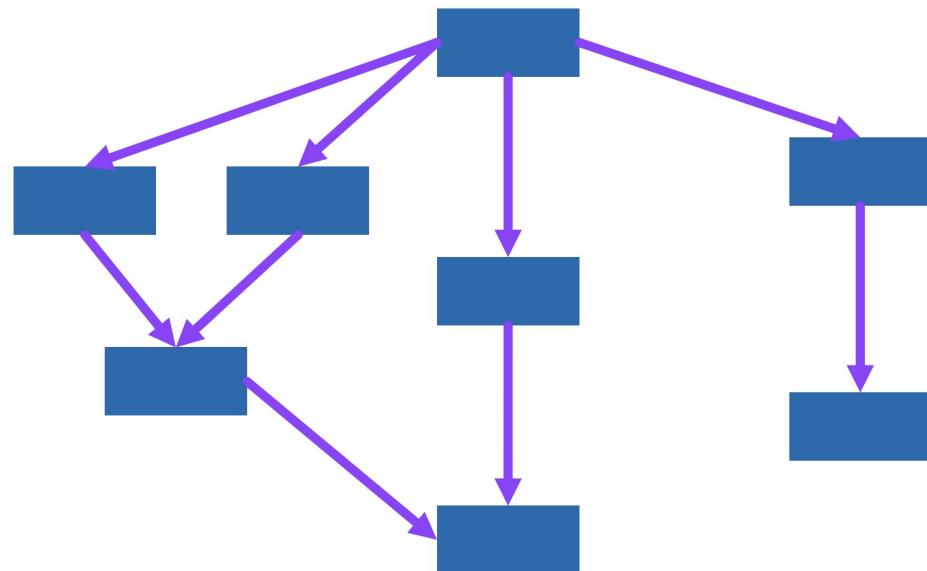
- Applications changed
- Also, MapReduce model too rigid
  - Just 2 steps per MapReduce job
  - You can chain multiple jobs together
    - **What's the catch?**



# MapReduce

- What if I want this type of computation?

- Direct Acyclic Graph
- More general



# Summary

- MapReduce couldn't keep up:
  - Applications that need more than 1 MapReduce jobs
  - Iterative computations that reuse a lot of data
  - Interactive computations that need fast response
- MapReduce didn't take advantage of memory

*Even if whole workload can fit into memory*

enters  
Spark

The diagram consists of two words: "enters" in blue at the top left and "Spark" in large orange letters below it. Five purple arrows originate from the letters "e", "n", "t", "e", and "r" of the word "enters". These arrows point to the corresponding letters "S", "p", "a", "r", and "k" in the word "Spark".

# Spark History

- Started in 2009
  - As a faster alternative for MapReduce



## **Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing**

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma,  
Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica

*University of California, Berkeley*



# Spark History

- Now a multi-billions company



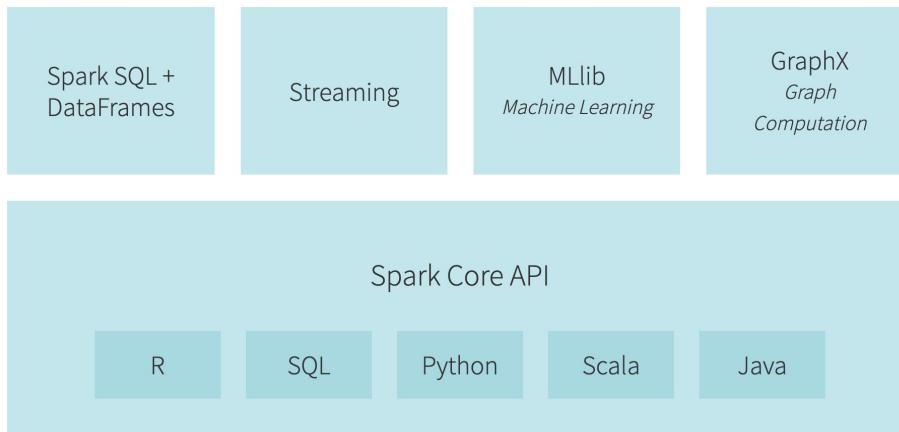
**Databricks announces \$400M round on  
\$6.2B valuation as analytics platform  
continues to grow**

Ron Miller @ron\_miller 2:14 am +08 • October 23, 2019

 Comment

# Spark History

- With a vibrant ecosystem

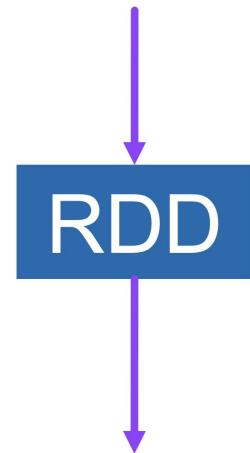


# Spark

- RDD Abstraction
  - RDDs in Spark are like objects in Java
  - Meaning:
    - Your Spark programs manipulate RDDs

Spark's first-class citizen

Resilient  
Distributed  
Dataset



# Spark RDD

it's just a

Big  
collection

Spark will try to fit partitions in  
memory



... and it is partitioned

# Life of a RDD

## RDD lifecycle



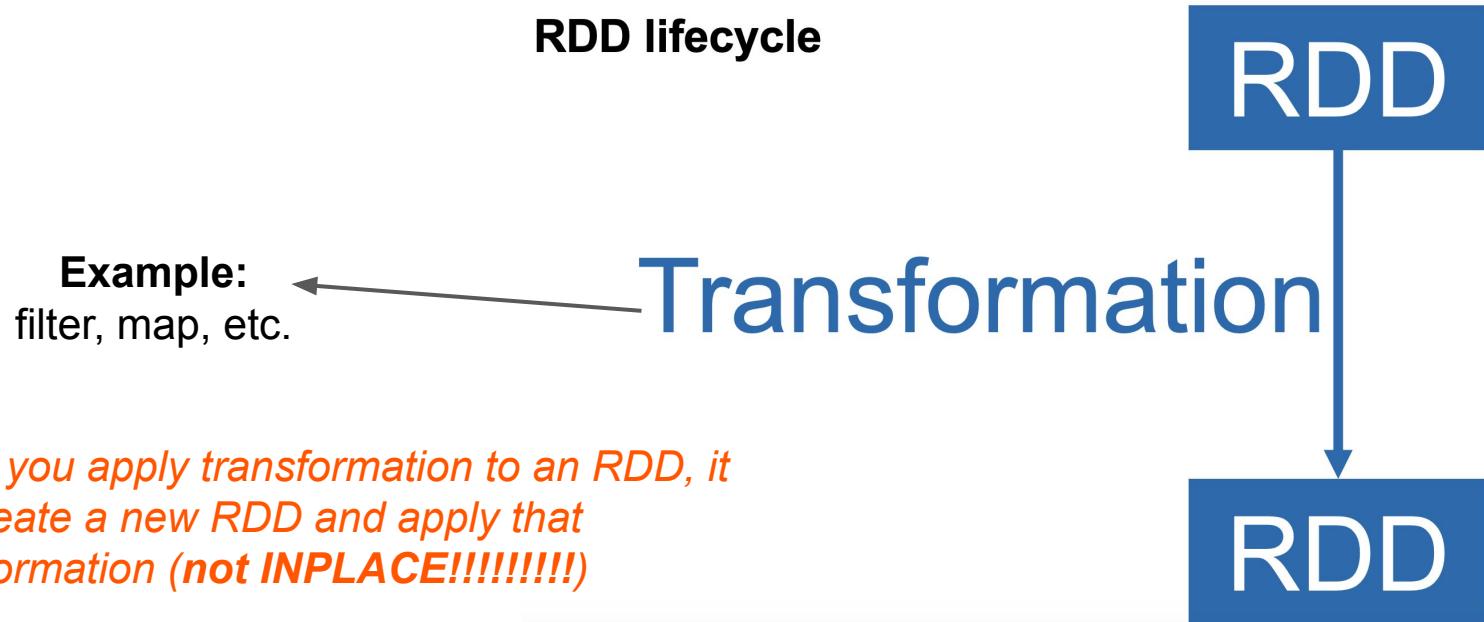
Creation

RDD

RDD

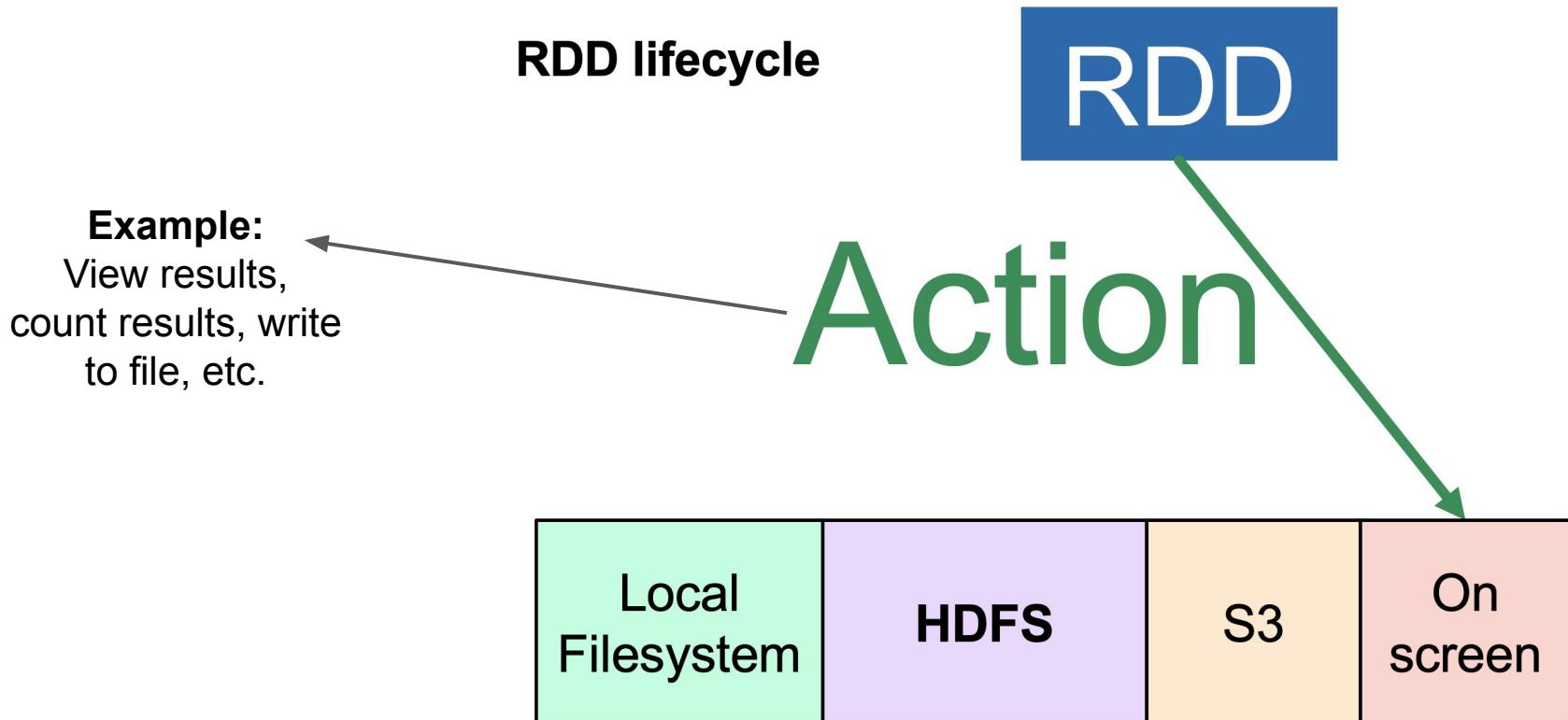


# Life of a RDD



**IMMUTABLE!!!!!!**

# Life of a RDD



# Life of a RDD

- Example

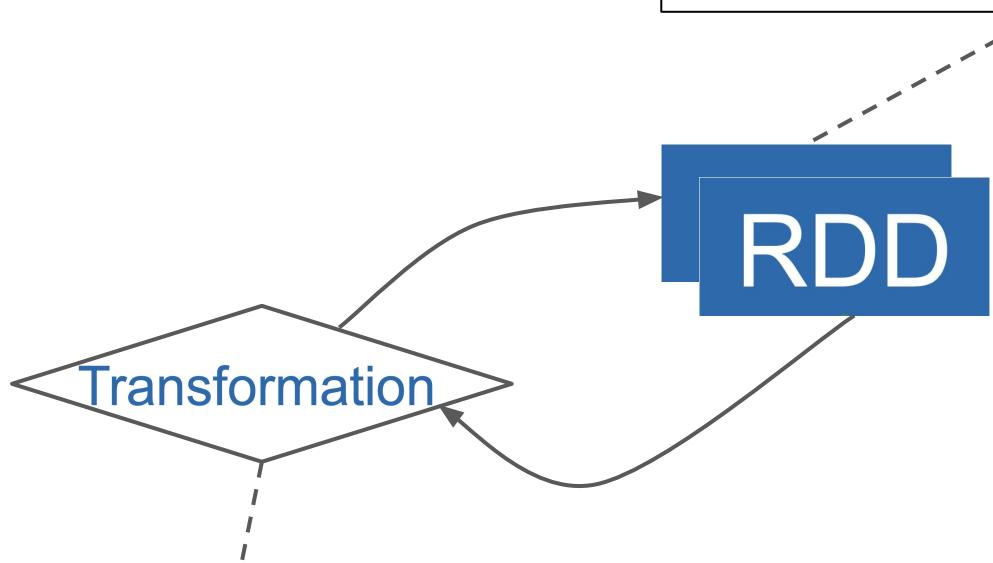
RDD

```
textFile = sc.textFile("hdfs://somefile")
```

# Life of a RDD

- Example

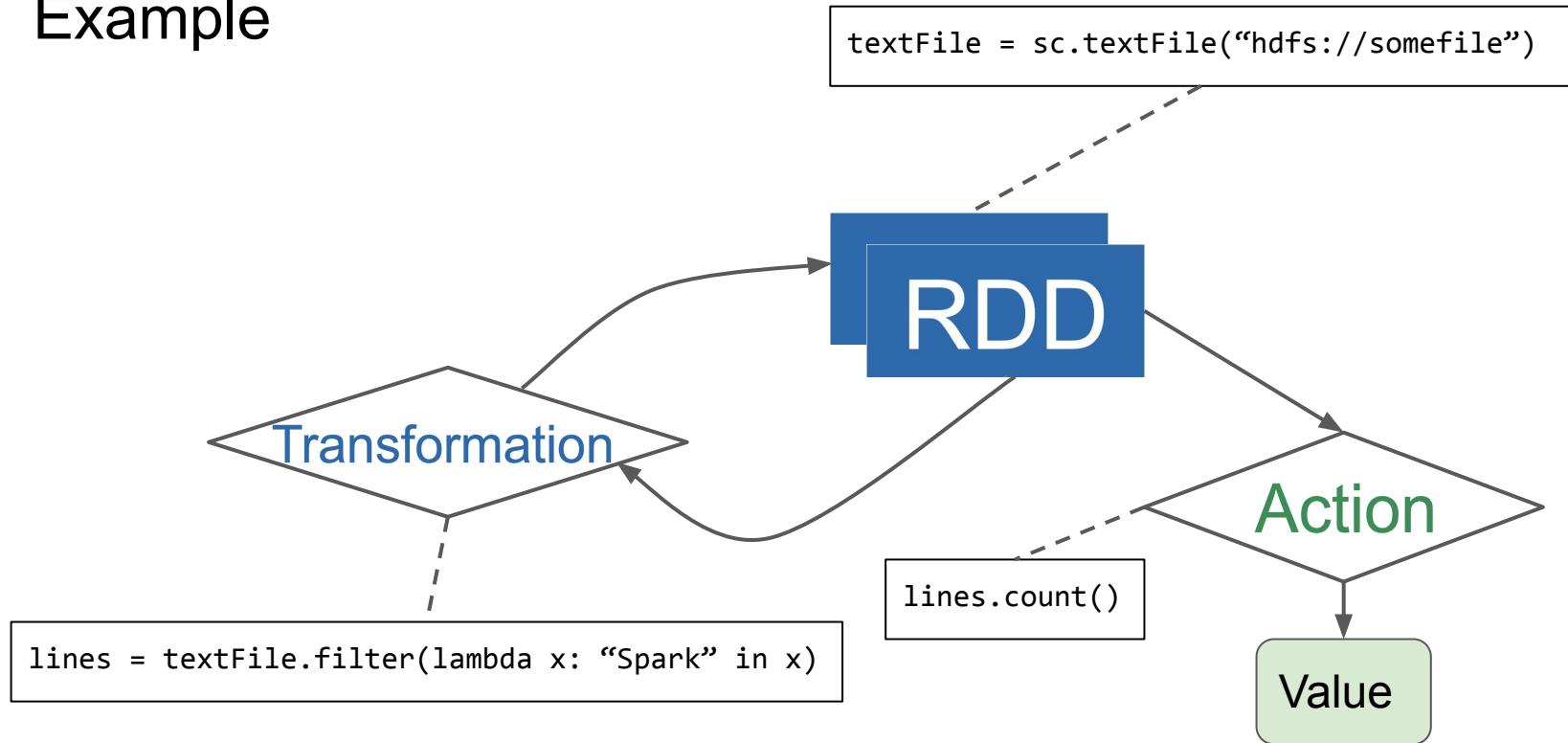
```
textFile = sc.textFile("hdfs://somefile")
```



```
lines = textFile.filter(lambda x: "Spark" in x)
```

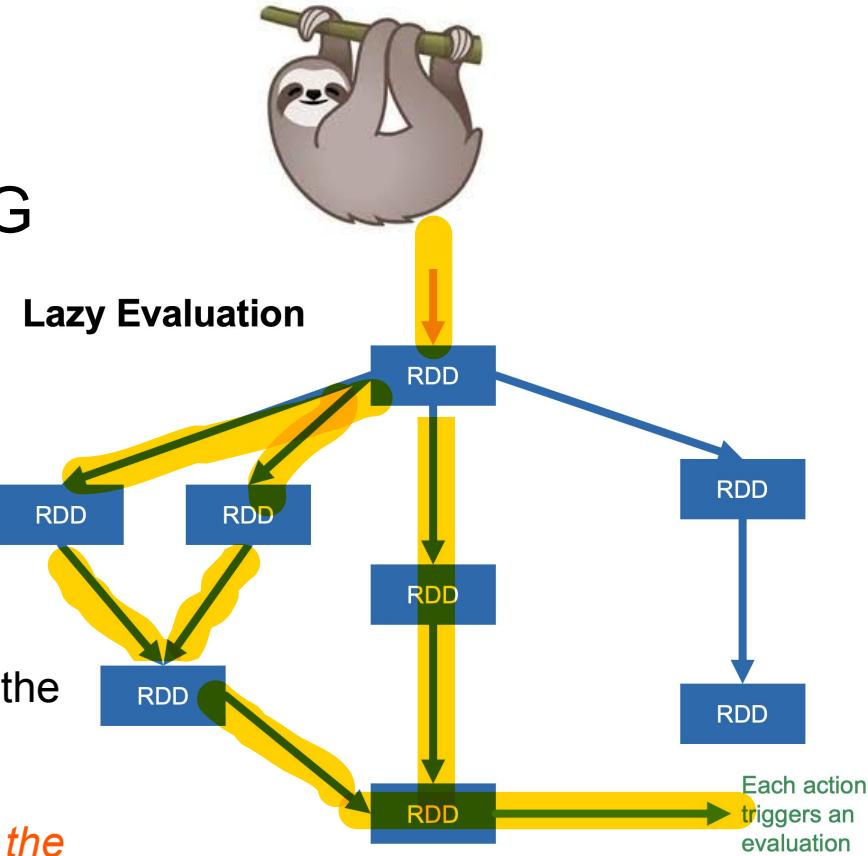
# Life of a RDD

- Example



# Spark Execution

- Transformed RDDs form a DAG
- Lazy evaluation:
  - Does not execute transformation immediately
  - Wait until Action!
    - Because it's **when** you **want** to see the results



*Doesn't apply all the transformation at once (if not the memory consumption will be huge)*

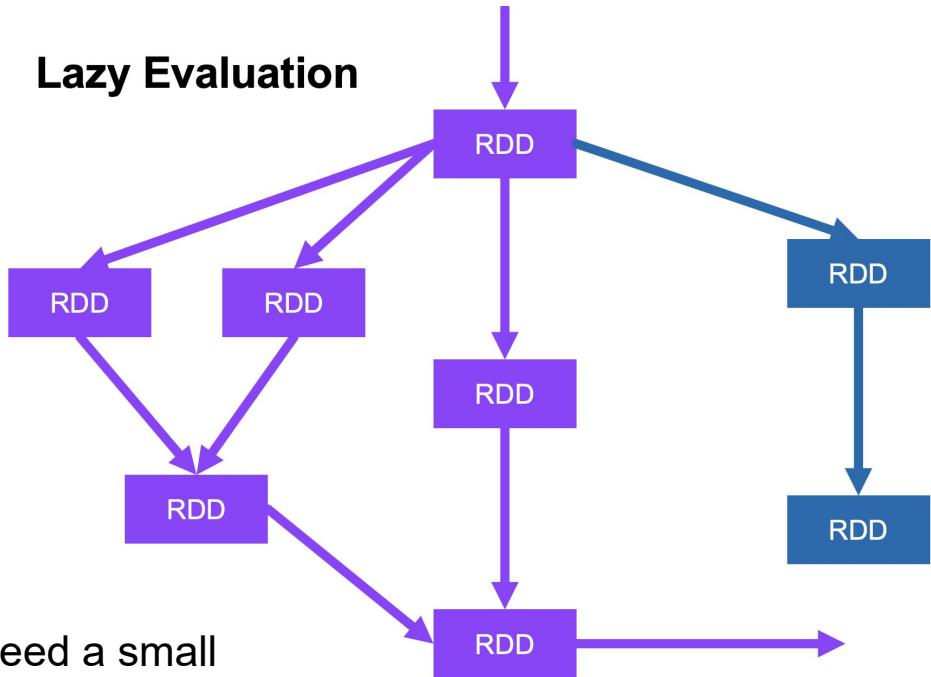
# Spark Execution

- Lazy evaluation

Almost all functional programming languages do this!

Why delaying execution?

Faster, sometimes you only need a small fraction of original RDD



# Spark APIs - Transformation

## Transformation

- Create one RDD from another

### Transformations

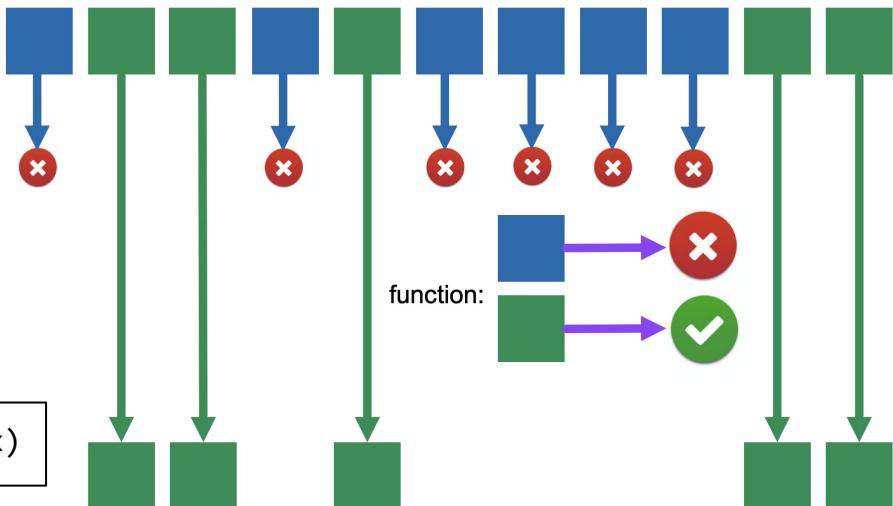


# Transformation

- Filter:
  - Input: RDD, a filter function
  - Output: RDD with element satisfying the function
  - Example:

```
lines = textFile.filter(lambda x: "Spark" in x)
```

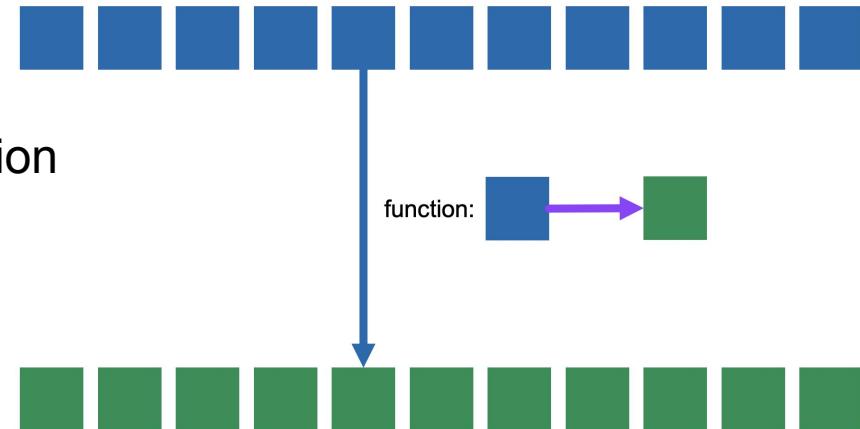
## Transformations: filter



# Transformation

- Map:
  - Input: RDD, map function
  - Output: RDD, where the map function  
applied to every element

Transformations: map



# Transformation

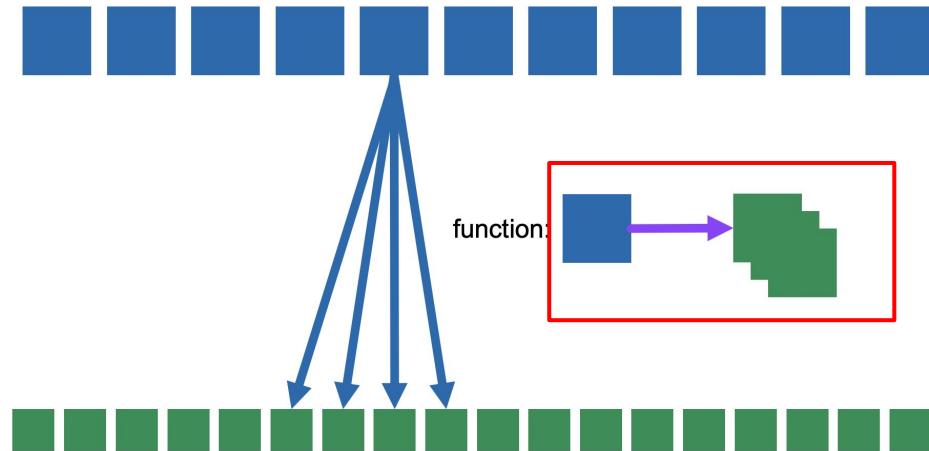
- **flatMap**

- Like map, but
- Input function returns a list
- Output RDD flattens the list

**map:**  
f is one to one  
 $f(x) \rightarrow y$

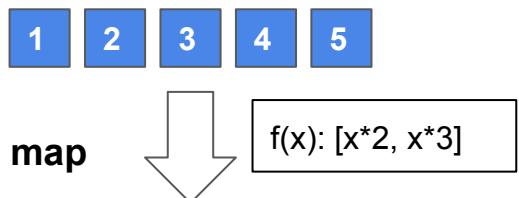
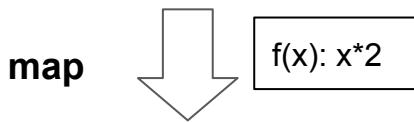
**flatMap:**  
f is one to many  
 $f(x) \rightarrow [y_1, y_2, \dots]$

## Transformations: flatMap

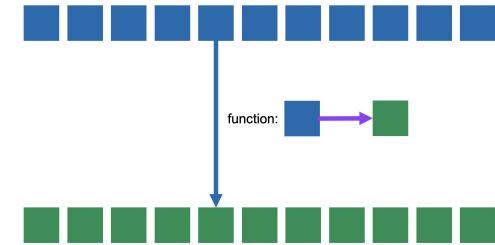


# Transformation

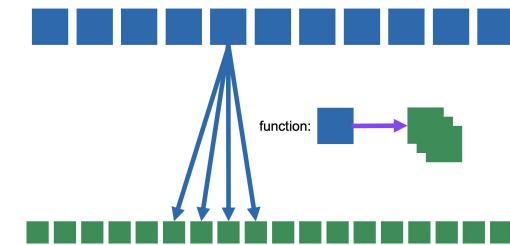
- map vs. flatMap



Transformations: map



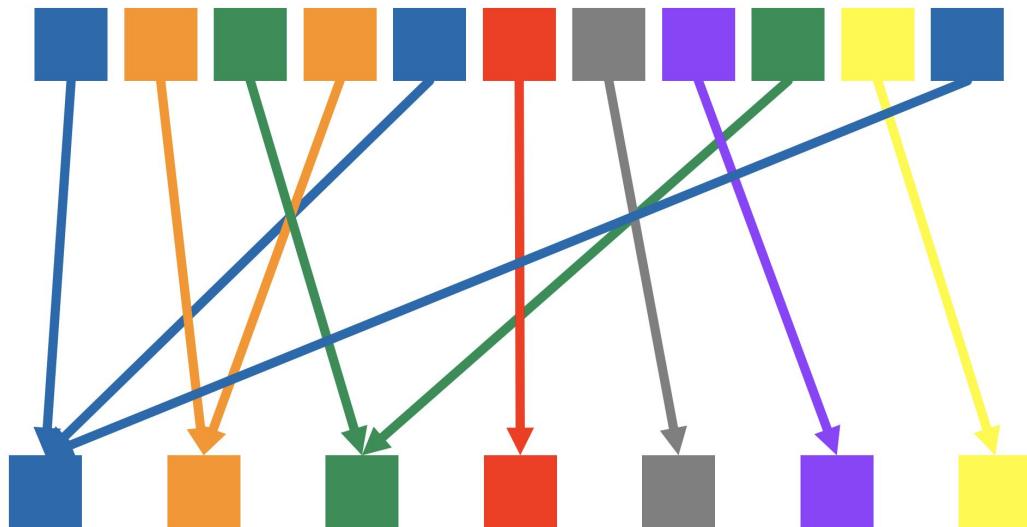
Transformations: flatMap



# Transformation

- Distinct:
  - Input: RDD
  - Output: RDD without duplicates

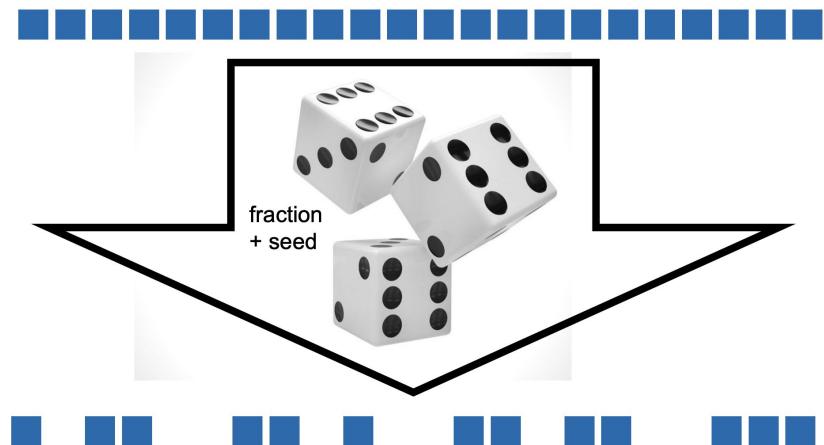
**Transformations: distinct**



# Transformation

- Sample
  - Input: RDD, fraction, seed
  - Output: RDD with sampled element

Transformations: sample

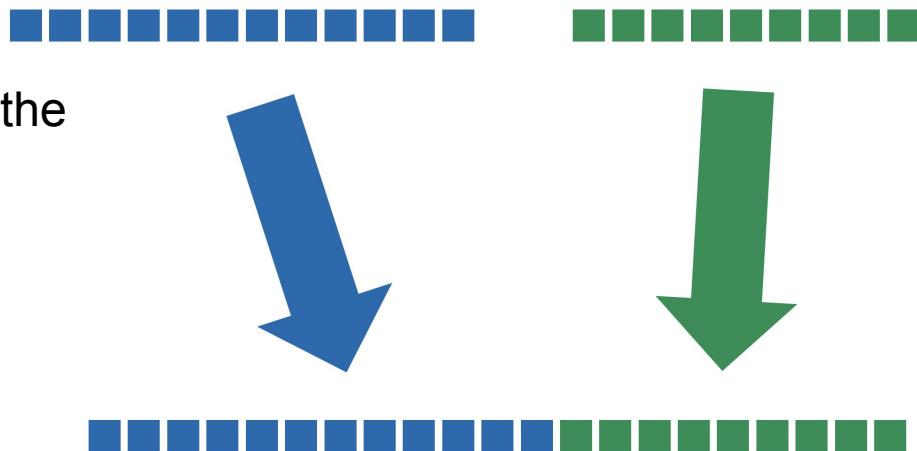


# Transformation

- Union:

- Input: 2 RDDs
- Output: RDD with elements from the input

**Transformations: union**

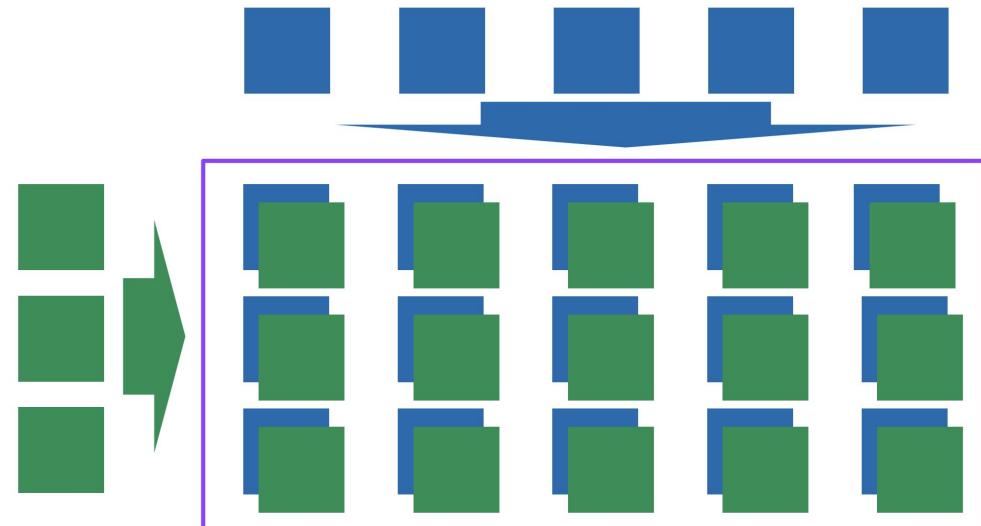


# Transformation

- Cartesian:

- Input: RDD T and U
- Output: RDD of  $(T, U)$  pairs

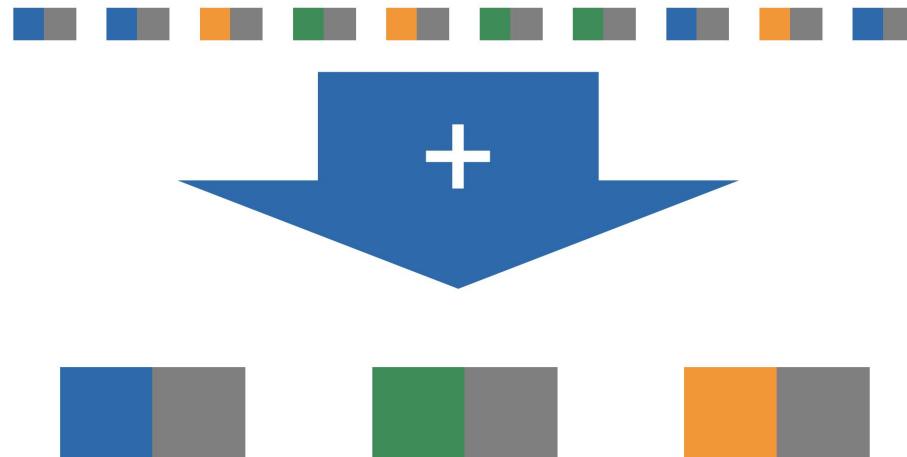
**Transformations: cartesian product**



# Transformation

- `reduceByKey`:
  - *Similar* to Reduce function in MapReduce
  - Input: RDD of key-value pairs, reduce function f
  - Output: RDD of key-value pairs where values for each key is aggregated using f

Transformations: `reduce by key`



# Transformation

- `reduceByKey`

- $f$  is an **accumulator** function

$(k, v_1)$   
 $(k, v_2)$   
 $(k, v_3)$   
...  
 $(k, v_n)$   
 $(k', v_1)$   
...



$f$



$\dots f(f(f(0, v_1), v_2), v_3), \dots$

$f(\text{accumulator}, \text{val}) \rightarrow \text{val}$

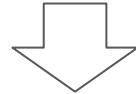
Transformations: reduce by key



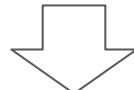
# Transformation

- `reduceByKey`
  - Example

`(1, 2)` `(2, 8)` `(1, 3)` `(4, 2)` `(2, 2)` `(1, 9)`

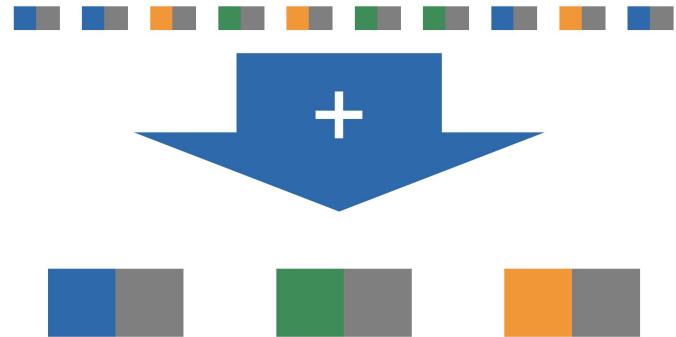


`reduceByKey(lambda acc, v: acc+v)`



`(1, 14)` `(2, 10)` `(4, 2)`

Transformations: reduce by key



# Transformation

- MapReduce's reduce vs. Spark's reduceByKey

## MapReduce

```
reduce(k, [vals]) -> (k', val')
```

## Spark

```
f(acc, val) -> val'
```

(1, 2)      (2, 8)      (1, 3)      (4, 2)      (2, 2)      (1, 9)

*In Spark, the value comes in streams*

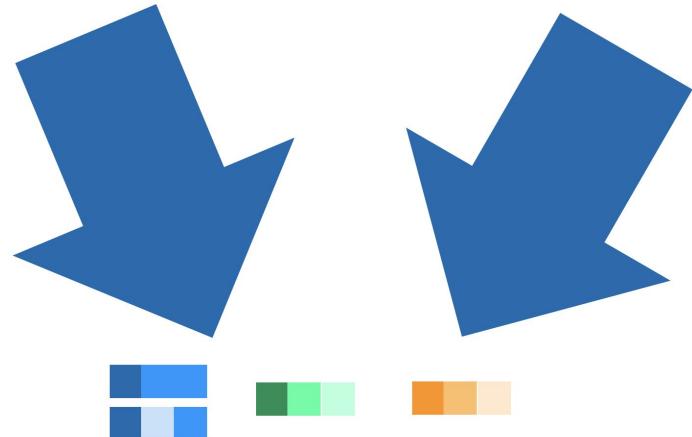
```
reduce(1, [2,3,9])
reduce(2, [8,2])
reduce(4, [2])
```

1, f(f(f(0,2),3),9)
2, f(f(0,8),2)
4, f(0,2)

# Transformation

- Join:
  - Input: two RDDs of key-value pairs
  - Output: joined by key
  - Just like SQL join on 2 tables
    - Each table has 2 columns

**Transformations: join**

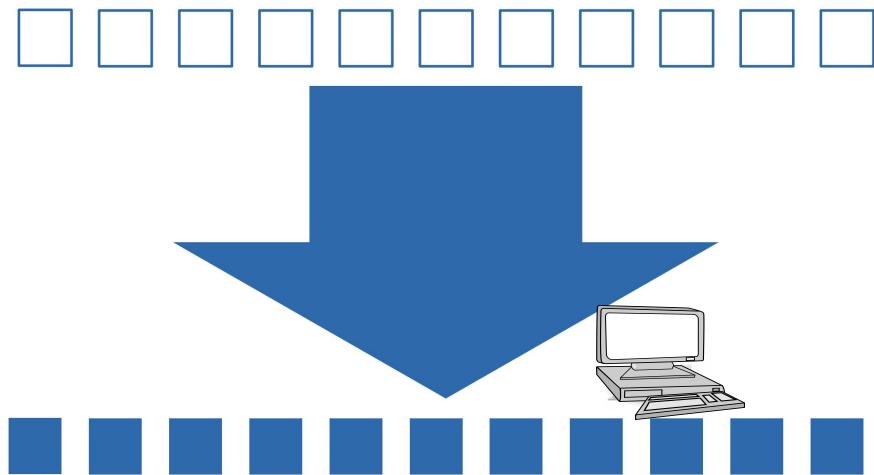


# Action

- Collect
  - See all values

## Spark APIs - Action

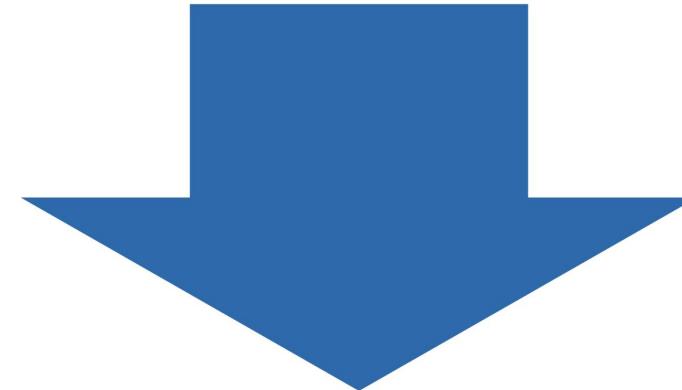
**Actions: collect**



# Action

**Actions: count**

- Count
  - Return number of elements

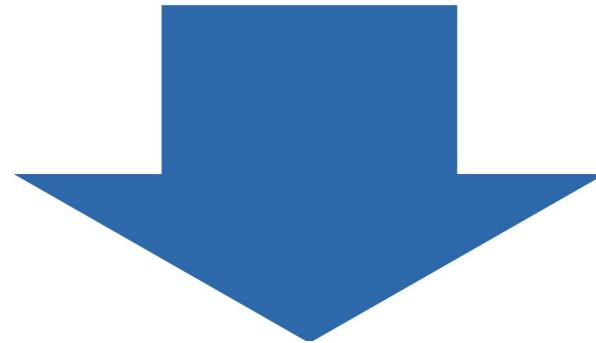


22

# Action

- Count by value
  - How many times a value appears

**Actions: count by value**



12 4 6

# Action

- Count by key

- How many times a key appears

**Actions: count by key**



6 1 4

# Action

- Take
  - Return a given number of elements

**Actions: take**



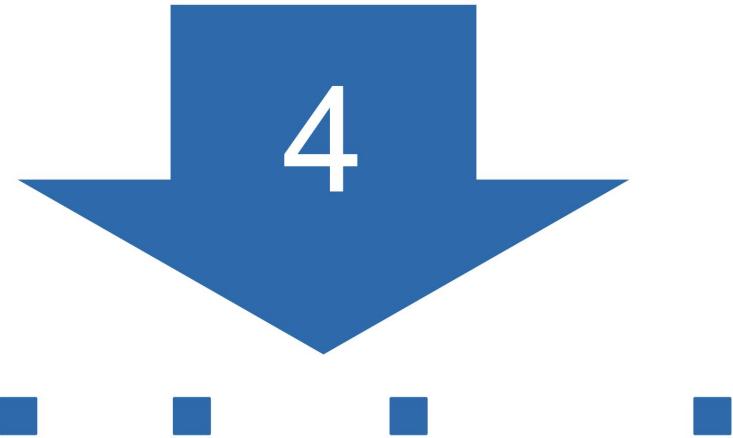
# Action

- **takeSample**

**Actions: takeSample**



- Sample a number of elements



# APIs Summary

- Many more APIs
- See <https://spark.apache.org/docs/latest/rdd-programming-guide.html>

## Transformations

The following table lists some of the common transformations supported by Spark. Refer to the RDD API doc ([Scala](#), [Java](#), [Python](#), [R](#)) and pair RDD functions doc ([Scala](#), [Java](#)) for details.

Transformation	Meaning
<code>map(func)</code>	Return a new distributed dataset formed by passing each element of the source through a function <code>func</code> .
<code>filter(func)</code>	Return a new dataset formed by selecting those elements of the source on which <code>func</code> returns true.
<code>flatMap(func)</code>	Similar to map, but each input item can be mapped to 0 or more output items (so <code>func</code> should return a Seq rather than a single item).
<code>mapPartitions(func)</code>	Similar to map, but runs separately on each partition (block) of the RDD, so <code>func</code> must be of type <code>Iterator&lt;T&gt; =&gt; Iterator&lt;U&gt;</code> when running on an RDD of type <code>T</code> .
<code>mapPartitionsWithIndex(func)</code>	Similar to <code>mapPartitions</code> , but also provides <code>func</code> with an integer value representing the index of the partition, so <code>func</code> must be of type <code>(Int, Iterator&lt;T&gt;) =&gt; Iterator&lt;U&gt;</code> when running on an RDD of type <code>T</code> .
<code>sample(withReplacement, fraction, seed)</code>	Sample a fraction <code>fraction</code> of the data, with or without replacement, using a given random number generator seed.
<code>union(otherDataset)</code>	Return a new dataset that contains the union of the elements in the source dataset and the argument.
<code>intersection(otherDataset)</code>	Return a new RDD that contains the intersection of elements in the source dataset and the argument.
<code>distinct([numPartitions])</code>	Return a new dataset that contains the distinct elements of the source dataset.
<code>groupByKey([numPartitions])</code>	When called on a dataset of (K, V) pairs, returns a dataset of (K, Iterable<V>) pairs. <b>Note:</b> If you are grouping in order to perform an aggregation (such as a sum or average) over each key, using <code>reduceByKey</code> or <code>aggregateByKey</code> will yield much better performance. <b>Note:</b> By default, the level of parallelism in the output depends on the number of partitions of the parent RDD. You can pass an optional <code>numPartitions</code> argument to set a different number

## Actions

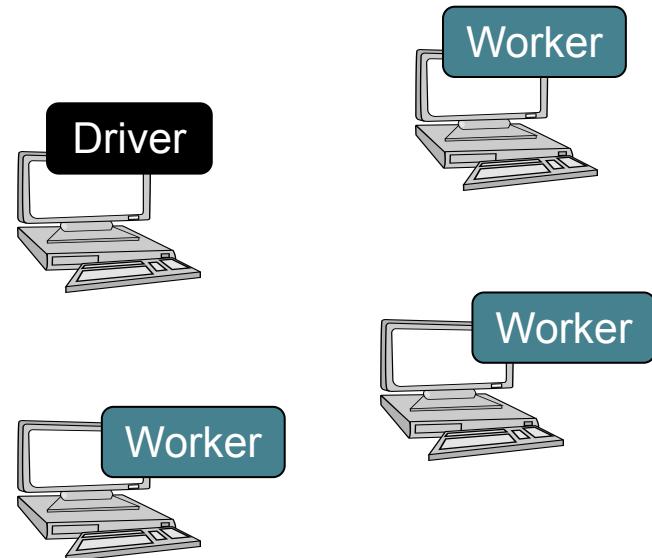
The following table lists some of the common actions supported by Spark. Refer to the RDD API doc ([Scala](#), [Java](#), [Python](#), [R](#)) and pair RDD functions doc ([Scala](#), [Java](#)) for details.

Action	Meaning
<code>reduce(func)</code>	Aggregate the elements of the dataset using a function <code>func</code> (which takes two arguments and returns one). The function should be commutative and associative so that it can be computed correctly in parallel.
<code>collect()</code>	Return all the elements of the dataset as an array at the driver program. This is usually useful after a filter or other operation that returns a sufficiently small subset of the data.
<code>count()</code>	Return the number of elements in the dataset.
<code>first()</code>	Return the first element of the dataset (similar to <code>take(1)</code> ).
<code>take(n)</code>	Return an array with the first <code>n</code> elements of the dataset.
<code>takeSample(withReplacement, num, [seed])</code>	Return an array with a random sample of <code>num</code> elements of the dataset, with or without replacement, optionally pre-specifying a random number generator seed.
<code>takeOrdered(n, [ordering])</code>	Return the first <code>n</code> elements of the RDD using either their natural order or a custom comparator.
<code>saveAsTextFile(path)</code>	Write the elements of the dataset as a text file (or set of text files) in a given directory in the local filesystem, HDFS or any other Hadoop-supported file system. Spark will call <code>toString</code> on each element to convert it to a line of text in the file.
<code>saveAsSequenceFile(path)</code> (Java and Scala)	Write the elements of the dataset as a Hadoop SequenceFile in a given path in the local filesystem, HDFS or any other Hadoop-supported file system. This is available on RDDs of key-value pairs that implement Hadoop's Writable interface. In Scala, it is also available on types that are implicitly convertible to Writable (Spark includes conversions for basic types like Int, Double, String, etc).
<code>saveAsObjectFile(path)</code> (Java and Scala)	Write the elements of the dataset in a simple format using Java serialization, which can then be loaded using <code>SparkContext.objectFile()</code> .

# Spark Example

## Example

- Task: log mining
  - Given a large number of log records
  - Filter out ERROR message
  - Interactively look for various pattern



# Example

Spark program / shell

Base  
RDD

```
lines = spark.textFile("hdfs://...")
```

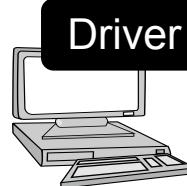
Nothing  
happens  
yet!

Driver

Worker

Worker

Worker



# Example

Transformed  
RDD

Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))
```

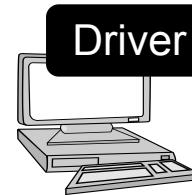
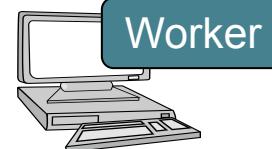
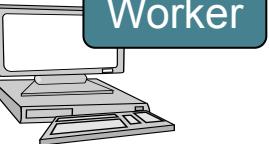
Nothing  
happen yet!

Worker

Driver

Worker

Worker



# Example

Transformed  
RDD

Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])
```

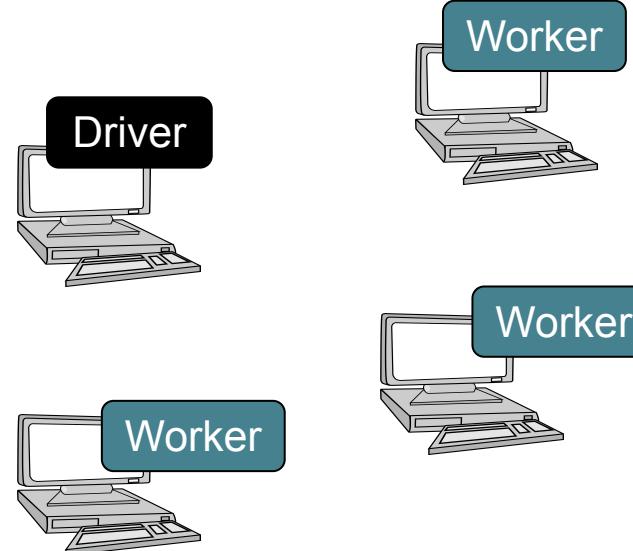
Nothing  
happen yet!

Driver

Worker

Worker

Worker



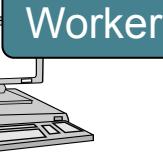
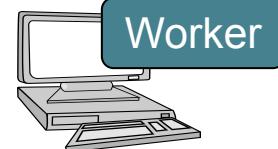
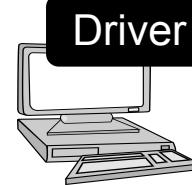
# Example

Next lecture

## Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])  
  
messages.persist()
```

Nothing  
happen yet!



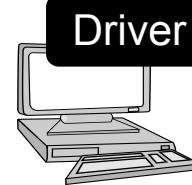
# Example

Transformed  
RDD

Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])  
  
messages.persist()  
  
mysqlMsg = messages.filter(lambda x: "mysql" in x)
```

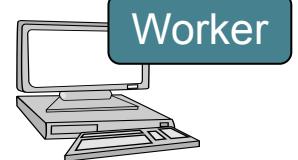
Nothing  
happen yet!



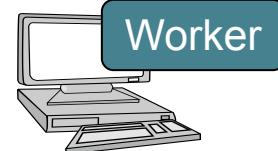
Driver



Worker



Worker



Worker

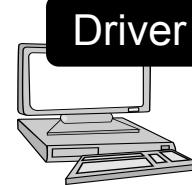
# Example

## Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])  
  
messages.persist()  
  
mysqlMsg = messages.filter(lambda x: "mysql" in x)  
  
mysqlMsg.count()
```

Action!

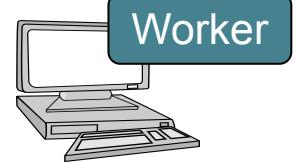
Something  
ABOUT to  
happen!



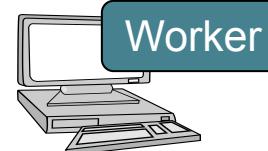
Driver



Worker



Worker



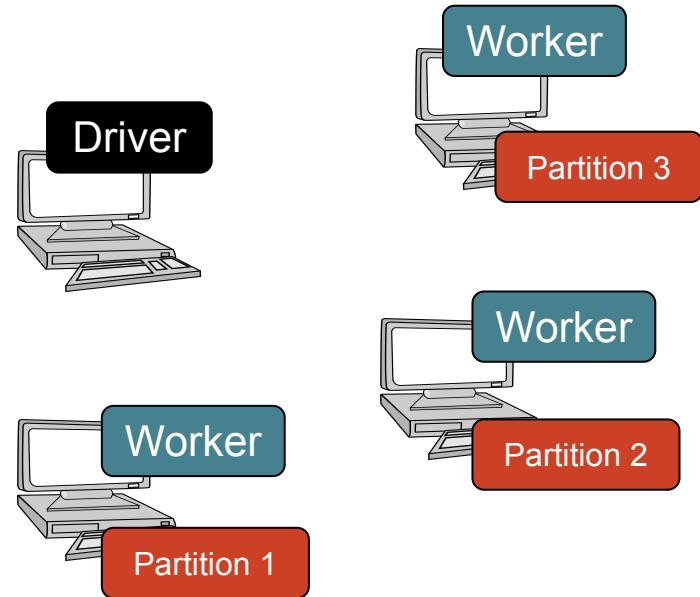
Worker

# Example

## Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])  
  
messages.persist()  
  
mysqlMsg = messages.filter(lambda x: "mysql" in x)  
  
mysqlMsg.count()
```

Action!

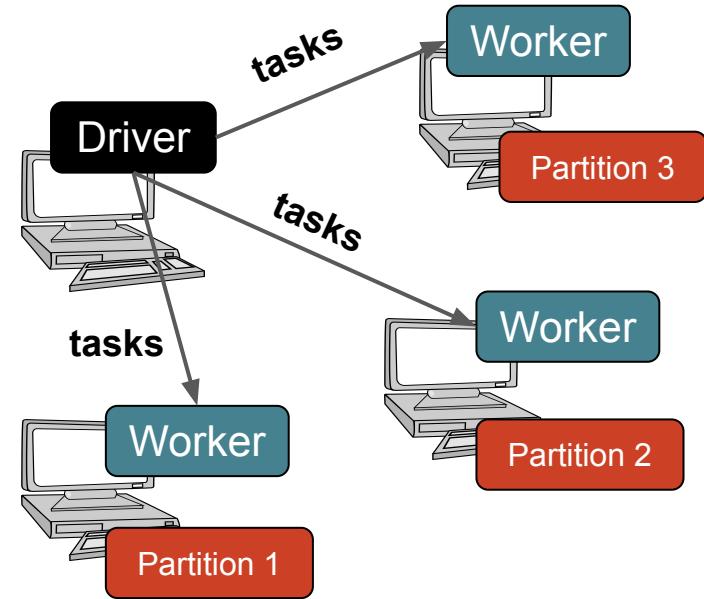


# Example

## Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])  
  
messages.persist()  
  
mysqlMsg = messages.filter(lambda x: "mysql" in x)  
  
mysqlMsg.count()
```

Action!

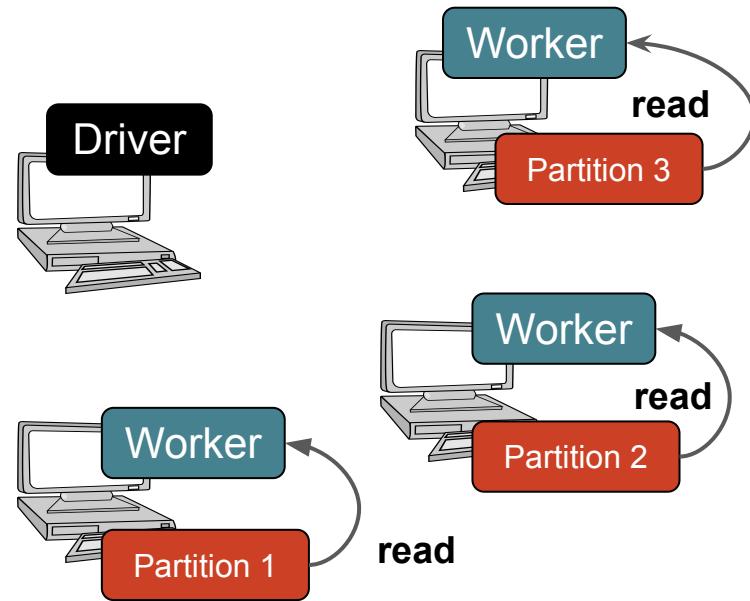


# Example

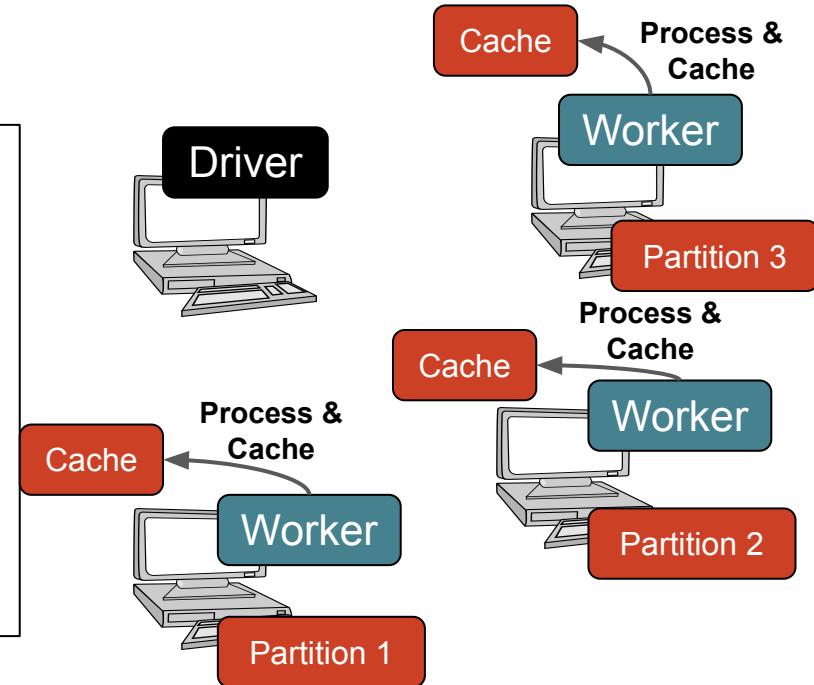
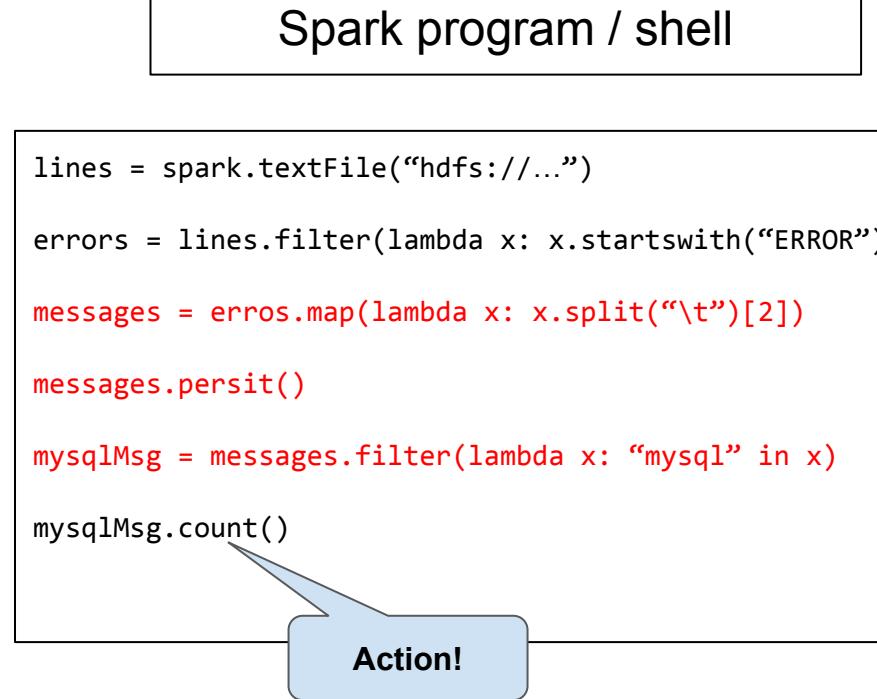
## Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])  
  
messages.persist()  
  
mysqlMsg = messages.filter(lambda x: "mysql" in x)  
  
mysqlMsg.count()
```

Action!



# Example

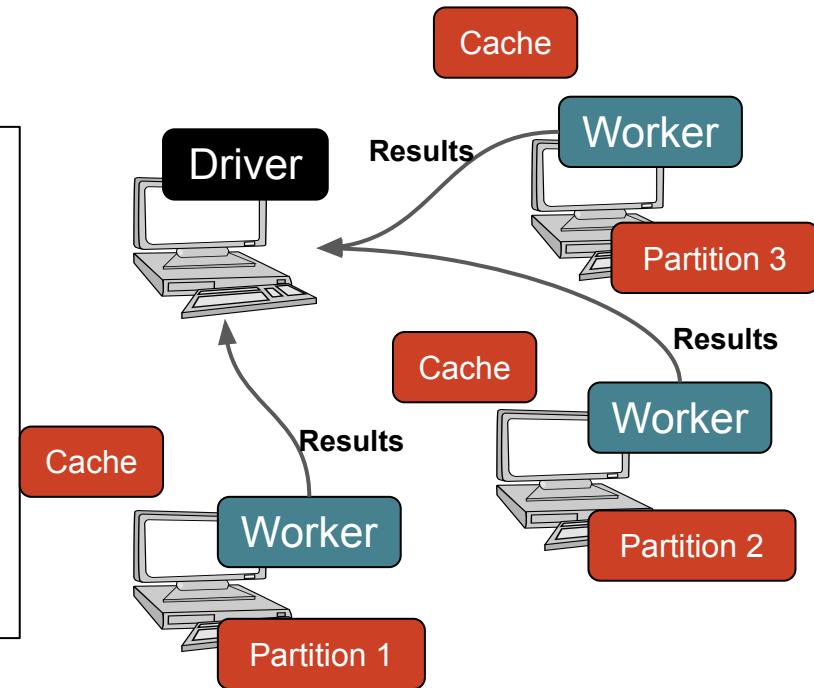


# Example

## Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])  
  
messages.persist()  
  
mysqlMsg = messages.filter(lambda x: "mysql" in x)  
  
mysqlMsg.count()
```

Action!

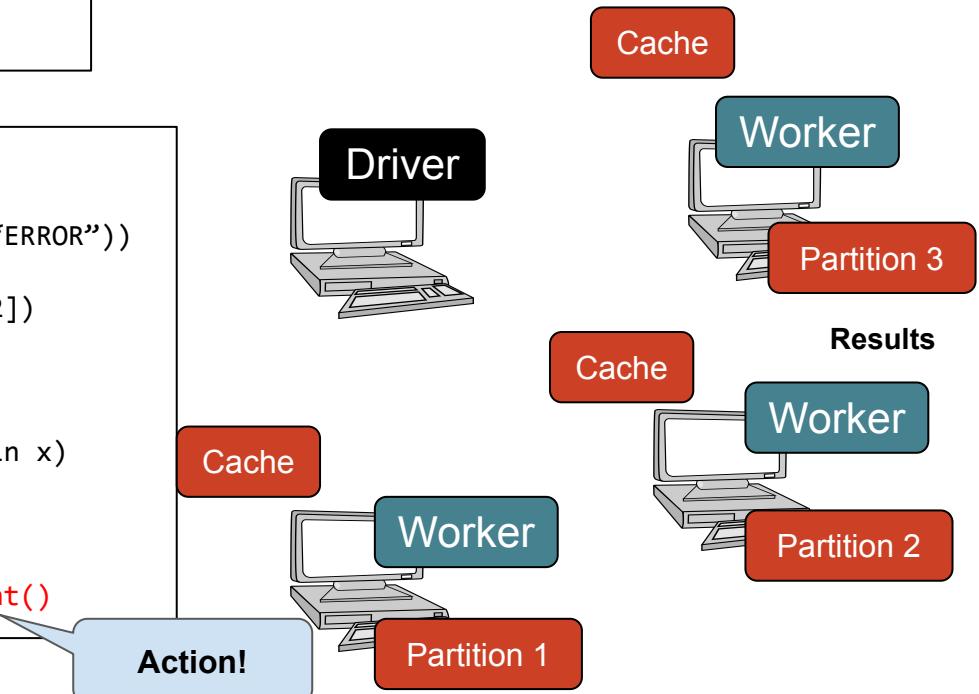


# Example

## Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])  
  
messages.persist()  
  
mysqlMsg = messages.filter(lambda x: "mysql" in x)  
  
mysqlMsg.count()  
  
messages.filter(lambda x: "mongodb" in x).count()
```

Action!

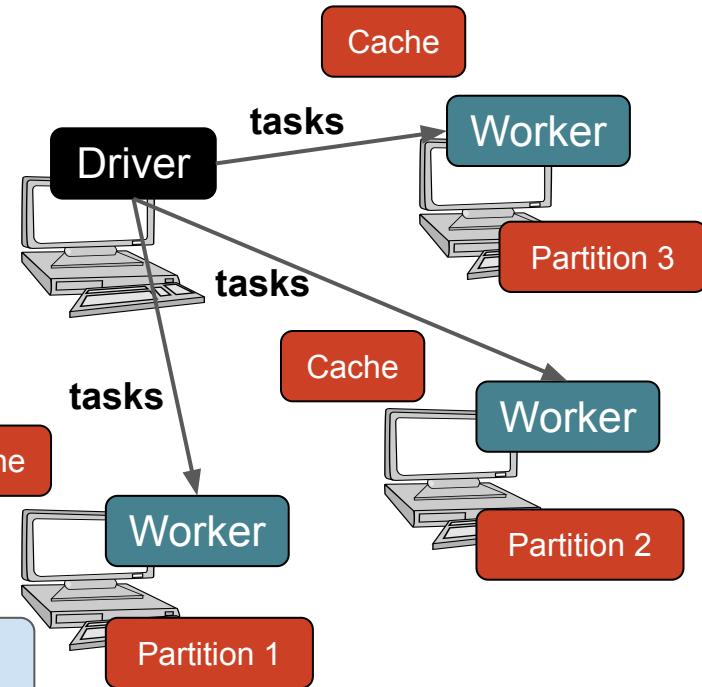


# Example

## Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])  
  
messages.persist()  
  
mysqlMsg = messages.filter(lambda x: "mysql" in x)  
  
mysqlMsg.count()  
  
messages.filter(lambda x: "mongodb" in x).count()
```

Action!

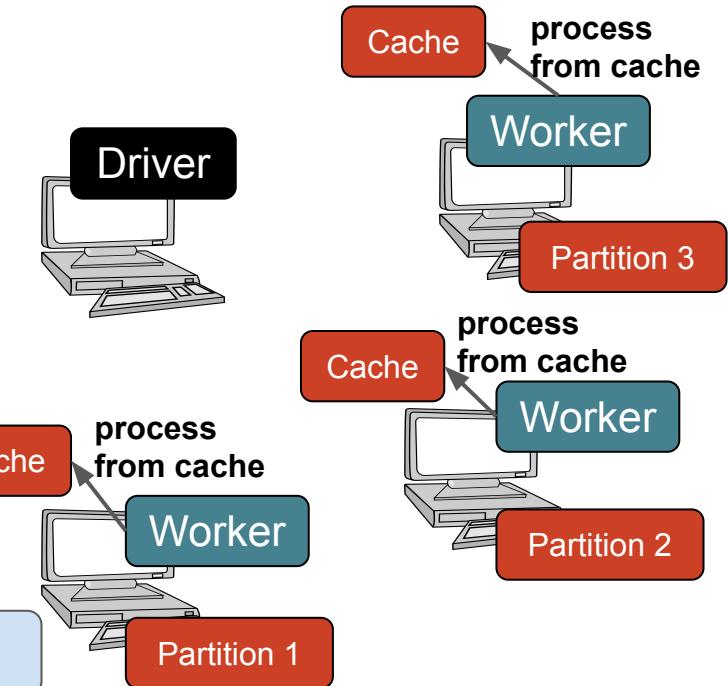


# Example

## Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])  
  
messages.persist()  
  
mysqlMsg = messages.filter(lambda x: "mysql" in x)  
  
mysqlMsg.count()  
  
messages.filter(lambda x: "mongodb" in x).count()
```

Action!

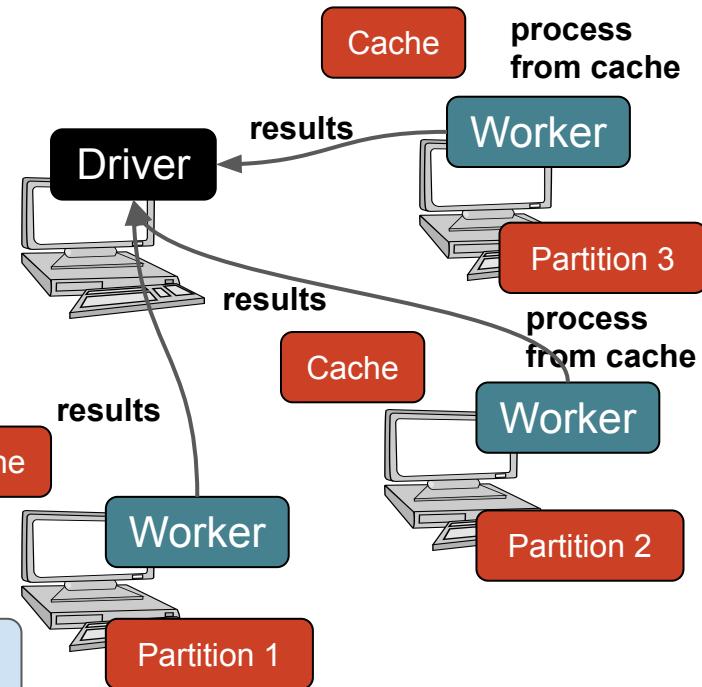


# Example

## Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])  
  
messages.persist()  
  
mysqlMsg = messages.filter(lambda x: "mysql" in x)  
  
mysqlMsg.count()  
  
messages.filter(lambda x: "mongodb" in x).count()
```

Action!



# Example

Cache your data → Faster Results

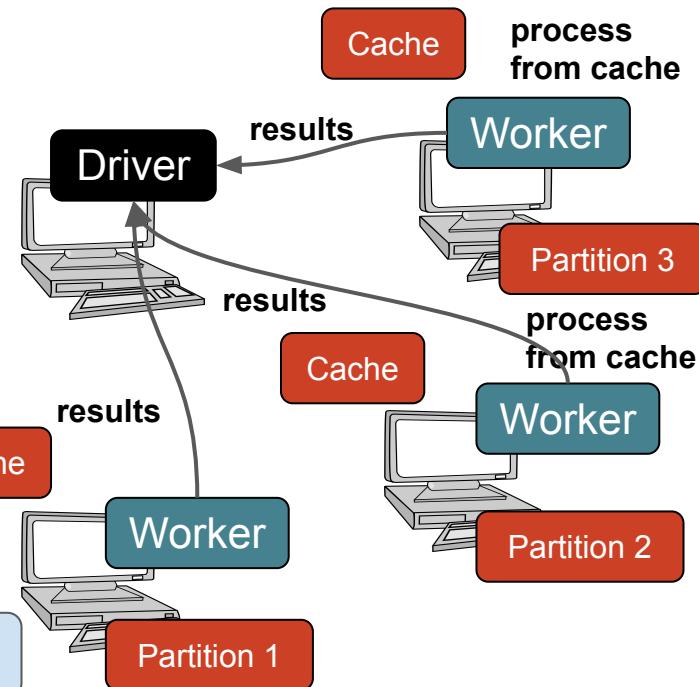
*Full-text search of Wikipedia*

- 60GB on 20 EC2 machines
- 0.5 sec from mem vs. 20s for on-disk

## Spark program / shell

```
lines = spark.textFile("hdfs://...")  
  
errors = lines.filter(lambda x: x.startswith("ERROR"))  
  
messages = errors.map(lambda x: x.split("\t")[2])  
  
messages.persist()  
  
mysqlMsg = messages.filter(lambda x: "mysql" in x)  
  
mysqlMsg.count()  
  
messages.filter(lambda x: "mongodb" in x).count()
```

Action!



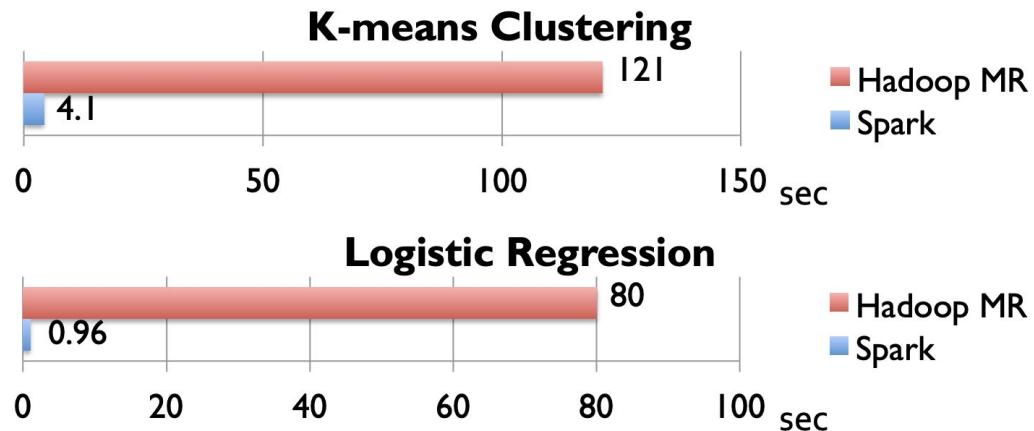
# Spark Demo

## Demo

- 3 nodes Spark cluster
  - 0.9GB text file
- Wordcount
- Text search
  - Search for lines with word “excellent”
  - Search for other words within those lines *find lines with both excellent and good  
.filter( ... ).filter( ... )  
(lambda x: 'excellent' in x)*

# Summary

- Spark much faster than MapReduce
- Especially for
  - Iterative computation
  - Interactive computation
- Key techniques:
  - Caching of RDDs
  - Lazy evaluation



# Question 1

- Why are RDDs called *immutable* if there are transformation operations on RDDs?

# Question 2

- Is Spark strictly better than MapReduce in all use cases?
  - In other words, why do people still use MapReduce?

## - RAM

*if you don't have enough RAM, Spark is not as good because it takes advantage of RAM to speed up computation.*

## - One Pass Jobs (More fundamental)

*If your job only takes one pass over data, still need to load all the data from disk. Spark has no fundamental advantage over MapReduce. MapReduce is designed specifically for one pass. (e.g. indexing)*

*(If more than one pass, Spark is strictly better because they cache the results, faster)*