

With preliminaries out of the way, let's dive into Spark in Scala. We choose Scala because it's the language Spark was written in. As a result, it poses the cleanest interface to the Spark environment.

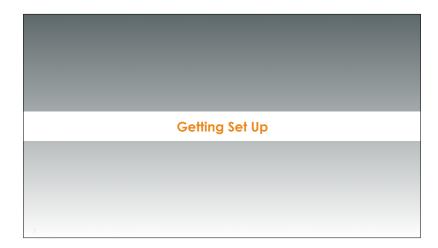


What We'll Cover



- The Spark execution model and Spark Sessions
- The Spark shell
- Scala immutabiles and mutables
- Resilient Distributed Datasets (RDDs)
- External datasets
- RDD operations: Transformations and Actions
- Function literals (aka closures or anonymous functions)
- Persistence
- Some Spark examples

We'll begin by covering these fundamental topics.



Before we begin, though, we should get you set up.

How to upgrade Spark software using our flash drive



Vagrant Version

We included some upgrade scripts as part of our file distribution that should now be in / vagrant. You should:

- 1.Log into your edge node
- 2. cd /vagrant
- ${\it 3. source ./spark-setup-edge.sh}$
- 4. Log into your control node
- 5.cd /vagrant
- 6.source ./spark-setup-control.sh

At this point, you should be up and ready to go with Spark 2.2.0 on your cluster. Your SPARK_HOME variable has been set to /vagrant/latest-spark.

If you are using a local two-node VirtualBox cluster created using Vagrant, you'll want to execute a couple of scripts to get your environment set up with the latest version of Spark (Spark 2.2.0)

No Upgrade Required



EMR Version

Your Amazon EMR instance should already be configured with version 2.2.0 of Spark.

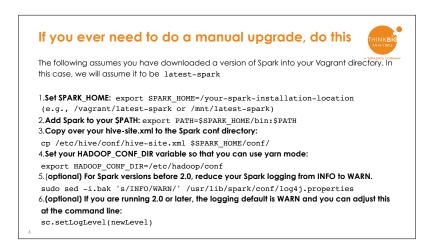
Please note that Spark has been installed in /usr/lib/spark.

The shell variable \$SPARK_HOME has not been set by EMR. Should you need to set the variable SPARK_HOME in later exercises, you will want to use the value /usr/lib/spark.

At this point, you should be up and ready to go with Spark 2.2.0 on your cluster.

5

If you are using Amazon Web Services Elastic Map Reduce instances, your instance should be all set up to use Spark 2.2.0.



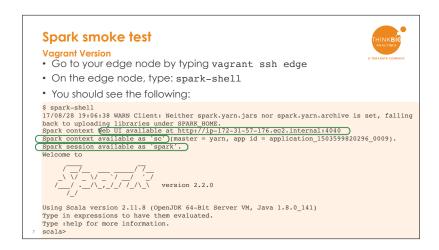
You can always set up a new spark path and manually

One of the best features of Spark is that it is easy to switch to using a different version than the one installed on your cluster. The process really consists of only 4 essential steps:

- 1. Setting your SPARK_HOME shell variable to point at the directory where your Spark installation lives.
- 2. Adding \$SPARK_HOME/bin to your execution path.
- 3. Copying over your hive-site.xml file to your Spark configuration directory
- 4. Setting the shell variable \$HADOOP_CONF_DIR to point at your Hadoop configuration directory.

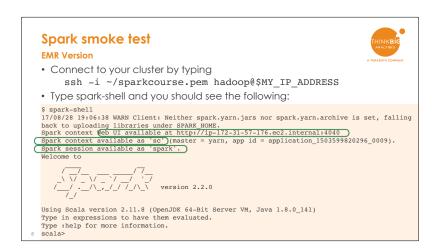
If you are using a Spark version prior to version 2, the default is for Spark to log all INFO messages to the console, which can definitely interfere with seeing what your program is doing. You can avoid this issue by setting your logging level to WARN by editing your log4j.properties file, as shown in step 5.

If you are running Spark 2.0 or later, the default logging level is WARN. You can change it to other levels, say ERROR, using a simple spark-shell command, sc.setLogLevel("ERROR").



If you're using vagrant, Go to your edge node by typing vagrant ssh edge. On the edge node, type: spark-shell. You should see the following.

Note that the spark-shell provides you with the Web address of the Spark user interface. It also has provided you with both a spark context (sc) and a spark session (spark). You'll need those to issue commands to the spark cluster.



If you're using EMR, Go to your edge node by typing

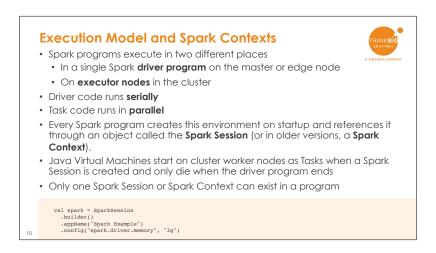
ssh -i ~/sparkcourse.pem hadoop@\$MY_IP_ADDRESS

On your cluster node, type spark-shell. You should see the following.

Note that the spark-shell provides you with the Web address of the Spark user interface. It also has provided you with both a spark context (sc) and a spark session (spark). You'll need those to issue commands to the spark cluster.



With setup done, let's jump into the fundamental components of basic Spark: RDDs, transforms, and actions.



• This has you configure exactly how much memory for your spark shell if you are writing your own spark shell

Spark programs execute in two different places

- 1. In a single Spark driver program on the master or edge node
- 2. On executor nodes in the cluster

Driver code runs serially, while Task code runs in parallel. Every Spark program creates this environment on startup and references it through an object called the Spark Session (or in older versions, a Spark Context).

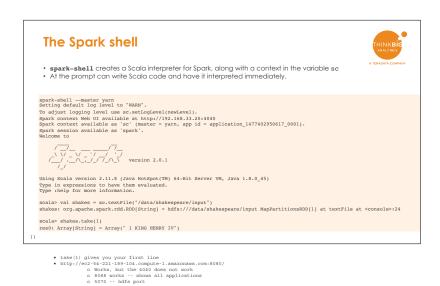
Java Virtual Machines start on cluster worker nodes as Tasks when a Spark Session is created and only die when the driver program ends

Only one Spark Session or Spark Context can exist in a program.

We can create a spark session using the following Scala chain:

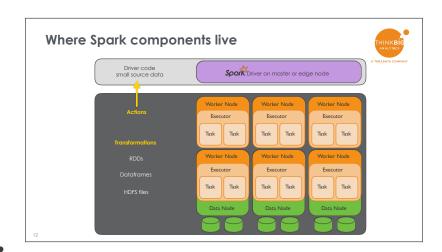
```
val spark = SparkSession
    .builder()
    .appName("Spark Example")
    .config("spark.driver.memory", "3g")
```

The builder method creates the spark session, appName then names that session, and the config method then sets various spark variables. In this case, we're setting the driver memory size to 3 gigabytes. We could similarly set the executor memory sizes using another config statement.



The spark shell is a REPL, or Read Evaluate Print Loop program. It reads what you type and then interprets it as Scala code.

The spark-shell is one of the biggest differences between the Spark Scala programming model and a traditional Java one. Scala has a REPL to ease development; at the time of this writing Java has none (although one is promised in Java 9).



If we look at this visually, the light gray box represents the edge node, depending on your cluster configuration. That's where the Spark driver lives. It runs serially there on one node.

In the dark gray box, we have the rest of the cluster data and worker nodes. That's where HDFS files are distributed over the disks connected to the data nodes. It's also where Spark DataFrames and RDDs exist.

Transformations are Spark methods that execute in parallel on the cluster. They run in parallel.

Actions are Spark methods that take data from the cluster and bring that data back to the driver. Because they must bring the data back to the driver, they force all the nodes involved to synchronize their efforts. They create a momentary serialization.

```
Walkthrough: 01-spark-basics-lab.{scala,md}
```

- Location of these items
- /mnt/sparkclass/exercises/spark-in-scala

Let's go hands-on with Spark with our first Scala lab. Two versions exist: one in Markdown and one in pure Scala code. We recommend you bring up one of those versions in a regular text editor (i.e., Notepad on Windows PCs, TextEdit or your favorite development text editor on a Mac) to view the text while having your terminal window along side it for typing commands.

Hands-On Lab: First Scala Program: Lab 02-spark-first-program.md



We saw in our walkthrough that we count the number of stock quotes for AAPL using code that looks like this:

val rdd = sc.textFile("hdfs:///data/stocks-flat/input")
val aapl = rdd.filter(line => line.contains("AAPL"))
aapl.count

In using that pattern to count the number of lines in the file /data/shakespeare/input that contain the word "king". You'll want to use the toLowerCase function before the contains function to ensure all the text is lower case.

A solution is on the next slide or in the lab script

- Parallelize RDD for Hadoop
- rdd.first is like the take(1) function

Hands-On Lab: First Scala Program

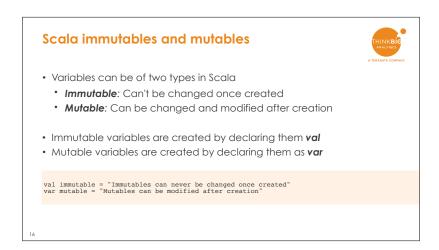


One solution looks like this

val rdd = sc.textFile("hdfs:///data/shakespeare/input")
val kings = rdd.filter(line => line.toLowerCase().contains("king"))
kings.count

You should have gotten 4773 lines that reference a king in all of Shakespeare. That's almost 3% of the 175,376 lines in all of Shakepeare's plays and poems.

15



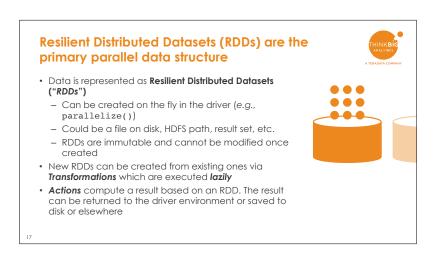
Variables can be of two types in Scala

Immutable: Can't be changed once created

Mutable: Can be changed and modified after creation

This is probably the most profound addition to cluster programming in the 21st century. By explicitly allowing the programmer to tell Scala which variables are immutable, Scala knows it can copy those variables throughout the cluster safely without affecting the validity of the computation. That means frequently used variables are likely to be local to a node and avoids excessive network traffic fetching variables. This is one of the fundamental insights within Spark: Immutables allow cluster programs to run faster because of greater data locality.

Immutable variables are created by declaring them **val** Mutable variables are created by declaring them as **var**



Data is represented as Resilient Distributed Datasets ("RDDs")

Can be created on the fly in the driver (e.g., parallelize())

Could be a file on disk, HDFS path, result set, etc.

RDDs are immutable and cannot be modified once created

New RDDs can be created from existing ones via Transformations which are executed lazily Actions compute a result based on an RDD. The result can be returned to the driver environment or saved to disk or elsewhere

Parallelized collections create RDDs



- \bullet The transformation parallelize turns a local collection into an RDD
- Once in an RDD, that parallelized collection can be operated on in parallel
- You can control the number of partitions that parallelize uses through a second argument, e.g., sc.parallelize(data, 10)
- Default number of partitions used is automatically set by your cluster

```
val data = Array(1, 2, 3, 4, 5)
val distData = sc.parallelize(data)
distData.take(3)
res4: Array[Int] = Array(1, 2, 3)
// Action: take first 3 elements
// result
```

The transformation parallelize turns a local collection into an RDD Once in an RDD, that parallelized collection can be operated on in parallel You can control the number of partitions that parallelize uses through a second argument, e.g., sc.parallelize(data, 10)

Default number of partitions used is automatically set by your cluster

More commonly, file inputs create RDDs



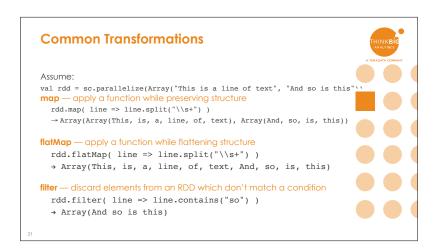
- Spark can create RDDs from any storage source supported by Hadoop
 - · Local file system, HDFS, Cassandra, HBase, Amazon S3
 - If using local file systems, the file must be replicated or shared on all worker nodes
 - Spark happily reads directories, compressed files, and paths specified by regular expressions
 - Like parallelize, textFile and its friends accept a second argument with the number of partitions to parallelize the input across

scala> val shakes = sc.textFile("hdfs:///data/shakespeare/input")
shakes: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[3] at textFile at <console>:21
scala> shakes.take(1)
resl: Array[String] = Array(" 1 KING HENRY IV")

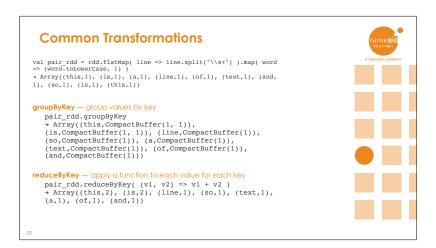
Spark can create RDDs from any storage source supported by Hadoop Local file system, HDFS, Cassandra, HBase, Amazon S3
If using local file systems, the file must be replicated or shared on all worker nodes
Spark happily reads directories, compressed files, and paths specified by regular expressions
Like parallelize, textFile and its friends accept a second argument with the number of partitions to parallelize the input across

RDD Operations: Transformations and Actions **Transformations convert an RDD into another RDD **Actions often turn an RDD into something else **Only actions** cause evaluation val lines = sc.textFile("hdfs:///data/shakespeare/input") // Transformation val lineLengths = lines.map(s => s.length) // Transformation val totalLength = lineLengths.reduce((a, b) => a + b) // Action

Transformations convert an RDD into another RDD Actions often turn an RDD into something else Only actions cause evaluation



- Filter acts like the where in SQL
- map/flatMap



• Key value pairs -

•

More Transformations



For RDDs For PairRDDs:

groupByKey reduceByKey aggregateByKey sortByKey join cogroup keys values repartition

map filter flatMap mapPartitions mapPartitionsWithIndex sample union intersection distinct distinct cartesian

repartition repartitionAndSortWithinPartitions

pipe coalesce

see http://spark.apache.org/docs/2.2.0/programming-guide.html#transformations

Common Spark Actions

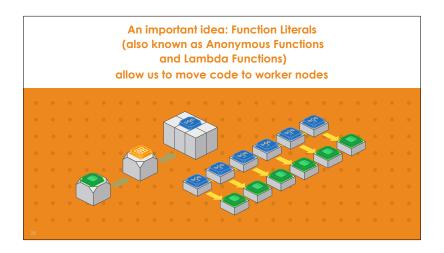
- collect gather results from nodes and return
- first return the first element of the RDD
- take(N) return the first N elements of the RDD
- saveAsTextFile write the RDD as a text file
- saveAsSequenceFile write the RDD as a SequenceFile
- **count** count elements in the RDD
- countByKey count elements in the RDD by key
- **foreach** process each element of an RDD
 - (e.g., rdd.collect.foreach(println))

See http://spark.apache.org/docs/2.2.0/programming-guide.html#actions



24





• Anonymous functions vs. functional literals

Function definitions are pretty common in most languages



- A function definition lets us name a stored operation that takes typed arguments and returns a typed result
- Scala can infer result types and even argument types

At least two ways exist to define functions and methods



- => is a definition or *transformation* operator; it takes its left hand arguments and transforms them using the right hand expression
- Types can be explicit or implied

```
/ Many ways exist to define functions and methods. Here are several val even = (i: Int) \Rightarrow { i % 2 == 0 } // explicit long form val even: (Int) \Rightarrow Boolean = i \Rightarrow { i % 2 == 0 } val even: Int \Rightarrow Boolean = i \Rightarrow (i % 2 == 0 } val even: Int \Rightarrow Boolean = i \Rightarrow (i % 2 == 0 ) // means first argument // implicit result approach // implicit result approach val add = (x: Int, y: Int) \Rightarrow X + y \Rightarrow Val add = (x: Int, y: Int) \Rightarrow X + y \Rightarrow Val add = (x: Int, y: Int) \Rightarrow X + y \Rightarrow Val add (Int, Int) \Rightarrow Int = (x,y) \Rightarrow { x + y } val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow Int = (x,y) \Rightarrow X + y \Rightarrow Val add: (Int, Int) \Rightarrow X + Y \Rightarrow Val add: (Int, Int) \Rightarrow X + Y \Rightarrow Val add: (Int, Int) \Rightarrow X + Y \Rightarrow Val add: (Int, Int) \Rightarrow X + Y \Rightarrow Val add: (Int, Int) \Rightarrow X + Y \Rightarrow Val add: (Int, Int) \Rightarrow X + Y \Rightarrow Val add: (Int, Int) \Rightarrow X + Y \Rightarrow Val add: (Int, Int) \Rightarrow X + Y \Rightarrow Val add: (Int, Int) \Rightarrow X + Y \Rightarrow Val add: (I
```

28

Anonymous functions allow us to do without the



- Function literals (a term from Jason Schwartz who wrote the O'Reilly Learning Scala book) are functions defined and passed in-line without giving them a name
- Function literals also go by the names of **anonymous functions** or **lambda functions**

```
// a function with a name greeting

val greeting = (x: String) => "Hello" + x

val names = List("Joe", "Mary", "Barbara")

names.nap(greeting)

// should get
// res2: List[String] = List(Hello Joe, Hello Mary, Hello Barbara)

// Now let's get rid of the name greeting

names.nap((x: String) => "Hello" + x)

// should get the same answer:
// res2: List[String] = List(Hello Joe, Hello Mary, Hello Barbara)
```

• Mapping a greeting to each name

Scala provides shorthand representations



- => allows us to define a function quickly
- _ provides a quick argument reference by appearance
- First _ means first argument, second _ means second argument

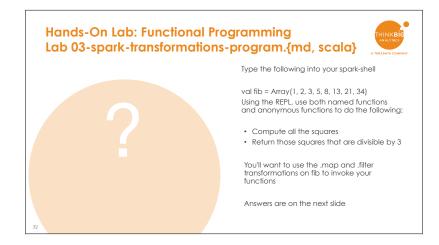
ullet _ defines the first argument

Underscores imply anonymous function arguments

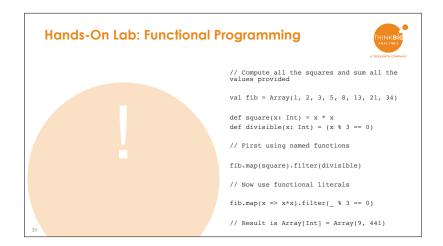
- Assuming the compiler can infer the argument types, you can use an underscore to both define an anonymous function and stand for its first argument
- If there is more than one argument, the arguments are assigned sequentially
- Use with caution; _ means many things in Scala and is very sensitive to context

```
/// Because we can infer the types, we can use the underscore shortcut
// for this function definition
names.map("Hello" + _)
// We can also use this for our even number tester
val nums = Array(1, 2, 3, 4, 5)
nums.map(_ 1 2 = 0)
// Should produce
// Array(Boolean] = Array(false, true, false, true, false)
// This may be harder to understand; it's a reducing function
// similar to that used in MapReduce
nums.reduceLeft(_ + _) // first plus second argument
int = 15
```

• Reduce is the action. The is for two different args



- val sets the values
- The def <something> defines the functions
- Anonymous functions like this can be sent to where the data lives



Persistence lets you reuse RDDs without recomputing them



- Executor memory is limited, so Spark carefully manages it using a least-recently-used (LRU) algorithm.
- Spark caches not only RDDs themselves, but the transformations that created them
- Because it remembers how to recreate every RDD, Spark can destroy them at any time to reclaim memory
- This can require a lot of re-computation if you need to use an RDD more than once
- Spark allows you to mark RDDs that you expect to reuse as persistent
- This is very important when doing iterative algorithms

rdd.persist()

Caching is one of the secrets of creating high-performing Spark applications. The cache or persist function (both work) allows you to specify which of your datasets should be retained in memory.

This technique is most useful in iterative programs that reference the same data structure many many times.

However, you should always make sure you time your program with and without explicit caching. Often Spark will do a better job of managing its memory than you will because it can optimize all aspects of the DAG, not just the piece you are looking at.

Also, you should be aware that because of lazy evaluation, you will gain no performance benefit the first reference after your persist command; the first persist is the one that puts the RDD into memory and saves it. However, the second and subsequent references should run faster.

Persistence Comes In Several Forms ATBRAGAGO ATBRAGAG	
Storage Level	Meaning
MEMORY_ONLY	Store RDD as deserialized Java objects in the JVM. If the RDD does not fit in memory, some partitions will not be cached and will be recomputed on the fly each time they're needed. This is the default level.
MEMORY_AND_DISK	Store RDD as deserialized Java objects in the JVM. If the RDD does not fit in memory, store the partitions that don't fit on disk, and read them from there when they're needed.
MEMORY_ONLY_SER	Store RDD as serialized Java objects (one byte array per partition). This is generally more space-efficient than deserialized objects, especially when using a fast serializer, but more CPU-intensive to read.
MEMORY_AND_DISK_SER	Similar to MEMORY_ONLY_SER, but spill partitions that don't fit in memory to disk instead of recomputing them on the fly each time they're needed.
DISK_ONLY	Store the RDD partitions only on disk.
MEMORY_ONLY_2, MEMORY_AND_DISK_2, etc.	Same as the levels above, but replicate each partition on two cluster nodes.
OFF_HEAP (experimental)	Store RDD in serialized format in Tachyon. Compared to MEMORY_ONLY_SER, OFF_HEAP reduces garbage collection overhead and allows executors to be smaller and to share a pool of memory, making it attractive in environments with large heaps or multiple concurrent applications. Furthermore, as the RDDs reside in Tachyon is discardable. Thus, Tachyon does not attempt to reconstruct a block that it evicts from memory. If you plan to use Tachyon as the off heap store, Sank's is compatible with Tachyon used-of-the-box.

You have many options in how to persist your RDDs. MEMORY_ONLY is the default and is the most common one. In those cases where your dataset is quite large, however, you may wish to specify MEMORY_ONLY_SER. Most of the other options are for special cases or where you have long-running computations that you can't afford to re-run in case of a node failure.

Spark Lab: 04-spark-functions-lab.{scala,md}

Wordcount Operations, MapReduce Style



- Split input line on whitespace
- Construct (word, 1) key-value pairs
- Group by key
- Sort by key
- Merge and sort by key ("merge-sort")
- Group values by key
- Sum values

37

Finally, A common mistake some programmers make is simply copying the functions they used in a MapReduce program over into their Spark code. For example, if we were simply emulating MapReduce wordcount, we'd be tempted to write our code as shown on the next slide.

Wordcount Operations, MapReduce Style

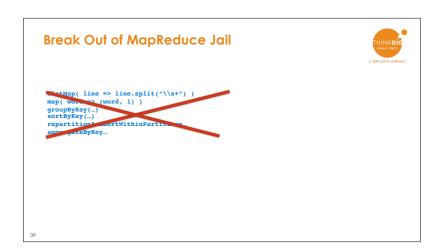


If we translated that directly to Spark we'd get....

```
flatMap( line => line.split("\\s+") )
map( word => (word, 1) )
groupByKey(...)
sortByKey(...)
repartitionAndSortWithinPartitions...
aggregateByKey...
```

20

Don't do this.



Don't do this.

```
Spark's Execution Model Is General Purpose;
Take Advantage Of It

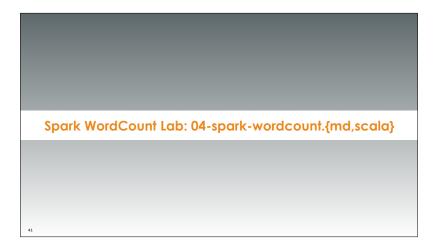
""Map( line => line.split("\\s+") )
map( WD. > (word, 1) )
groupByRey(...)
repartition wortWithinParticle
and setaByKey...

text.flatMap( line => line.split("\\\+") ).
map( word => (word.toLowerCase(), 1) ).
reduceByKey( (v1, v2) => v1 + v2 )

// or

text.flatMap( _.split("\\\+") ).
map(word => (word.toLowerCase(), 1) ).
reduceByKey( _...)
```

Don't be afraid to take advantage of Spark's flexibility and rewrite your algorithm with a more flexible DAG.



This lab is optional -- we touched on Wordcount already in lab 3, but if you want to step through it in detail, this is where to do it.

Summary



- Spark provides a new full-stack cluster development environment
- Resilient distributed datasets (RDDs) provide in-memory storage of data across the cluster instead of spilling to disk
- Transformation operations don't execute until an action is called
- Functional programming allows us to push code to the cluster nodes dynamically

42

This has been a quick introduction to Spark at the RDD level. It's a new way of thinking about cluster computing and it has huge uptake in the Hadoop community. However, the best is yet to come, as we'll see when we begin working with Spark SQL.