



# Developing Spark Applications

## Slide Guide

---

Version 5.1 – Summer 2016

*For use with the following courses:*

- DEV 3600 – Developing Spark Applications
- DEV 360 – Apache Spark Essentials
- DEV 361 – Build and Monitor Apache Spark Applications
- DEV 362 – Create Data Pipeline Applications Using Apache Spark

This Guide is protected under U.S. and international copyright laws, and is the exclusive property of MapR Technologies, Inc.

© 2016, MapR Technologies, Inc. All rights reserved. All other trademarks cited here are the property of their respective owners.





# DEV 360 - Apache Spark Essentials

---

## Slide Guide

Version 5.1 – Summer 2016

This Guide is protected under U.S. and international copyright laws, and is the exclusive property of MapR Technologies, Inc.

© 2016, MapR Technologies, Inc. All rights reserved. All other trademarks cited here are the property of their respective owners.



 MAPR Academy

## Apache Spark Essentials

Getting Started with Apache Spark Essentials

© 2015 MapR Technologies  MAPR

1

Welcome to Apache Spark Essentials. This introductory course enables developers to get started with Apache Spark. It introduces the benefits of Apache Spark for developing big data applications. In the first part of the course, you will use Spark's interactive shell to load and inspect data. You will then go on to build and launch a standalone Spark application. The concepts are taught using scenarios that also form the basis of hands-on labs.





## Learning Goals

- ▶ Describe Features of Apache Spark
- ▶ Define Spark Components
- ▶ Load data into Spark
- ▶ Apply dataset operations to Resilient Distributed Datasets
- ▶ Use Spark DataFrames for simple queries
- ▶ Define different ways to run your application
- ▶ Build and launch a standalone application

At the end of this course, you will be able to:

- Describe Features of Apache Spark
- Define Spark Components
- Load data into Spark
- Apply dataset operations to Resilient Distributed Datasets
- Use Spark DataFrames for simple queries
- Define different ways to run your application
- Build and launch a standalone application





## Prerequisites

### Required

- Basic to intermediate Linux knowledge, including:
  - The ability to use a text editor, such as vi
  - Familiarity with basic command-line options such as mv, cp, ssh, grep, cd, useradd
- Knowledge of application development principles
- A Linux, Windows or MacOS computer with the MapR Sandbox installed (On-demand course)
- Connection to a Hadoop cluster via SSH and web browser (for the ILT and vILT course)

### Recommended

- Knowledge of functional programming
- Knowledge of Scala or Python
- Beginner fluency with SQL
- HDE 100 - Hadoop Essentials

The prerequisites for this course include a basic to intermediate knowledge of Linux, basic knowledge of application development principles. For the on-demand course, you will also need the MapR Sandbox installed. Knowledge of functional programming, Scala or Python and SQL is recommended.





## Course Materials

- DEV360 SlideGuide.pdf
- DEV360 LabGuide.pdf
- DEV360 LabFiles.zip
- DEV360DATA.zip
- Connect to MapR Sandbox or AWS.pdf

The following course materials are provided:

- The slide guide provides the slides and notes.
- The Lab guide contains lab instructions for the lab activities.
- LabFiles.zip contains the files and solutions for the labs.
- DEV360DATA.zip contains the data files that you need for the labs.
- Connect to MapR Sandbox or AWS.pdf provides connection instructions.





Next Steps



## Lesson 1

Introduction to Apache  
Spark

© 2015 MapR Technologies  MAPR

5

You are now ready to start this course. Proceed to Lesson 1:  
Introduction to Apache Spark.



The MAPR Academy logo, consisting of the "MAPR" logo in red followed by the word "Academy" in a larger, grey, sans-serif font.

## Apache Spark Essentials

Lesson 1: Introduction to Apache Spark

© 2015 MapR Technologies The MAPR logo, located at the bottom right of the slide.

1

Welcome to Apache Spark Essentials, Lesson 1 – Introduction to Apache Spark. This lesson describes Apache Spark, lists the benefits of using Spark and defines the Spark components.



 Learning Goals

- ▶ Describe Features of Apache Spark
  - How Spark fits in Big Data ecosystem
  - Why Spark & Hadoop fit together
- ▶ Define Spark Components

© 2015 MapR Technologies  MAPR®

2

At the end of this lesson, you will be able to:

- Describe the features of Apache Spark, such as:
  - How Spark fits in the Big Data ecosystem, and
  - Why Spark & Hadoop fit together
- Also, you will be able to define the Spark Components

In this section, we will take a look at the features of Apache Spark, how it fits in the Big Data ecosystem



2



## What is Apache Spark?



- Cluster computing platform on top of storage layer
- Extends MapReduce with support for more components
  - Streaming
  - Interactive analytics
- Runs in memory

© 2015 MapR Technologies  MAPR®

3

- Apache Spark is a cluster computing platform on top of a storage layer.
- It extends MapReduce with support for more types of components such as streaming and interactive analytics.
- Spark offers the ability to run computations in memory, but is also more efficient than MapReduce for running on disk



 Why Apache Spark?A large, bold number '1' is centered within a white circle, which is set against a light blue square background.**Fast**

- 10x faster on disk
- 100x in memory

- Spark provides reliable in-memory performance. Iterative algorithms are faster as data is not being written to disk between jobs.
- In-memory data sharing across DAGs make it possible for different jobs to work with the same data quickly.
- Spark processes data 10 times faster than MapReduce on disk and 100 times faster in memory.



 Why Apache Spark?

### Ease of Development

- Write programs quickly
- More operators
- Interactive Shell
- Less code

- You can build complex algorithms for data processing very quickly in Spark.
- Spark provides support for many more operators than MapReduce such as Joins, reduceByKey, combineByKey.
- Spark also provides the ability to write programs interactively using the Spark Interactive Shell available for Scala and Python.
- You can compose non-trivial algorithms with little code.



 Why Apache Spark?

### Deployment Flexibility

- Deployment
  - Mesos
  - YARN
  - Standalone
  - Local
- Storage
  - HDFS
  - S3

© 2015 MapR Technologies  MAPR.

6

- You can continue to use your existing big data investments.
- Spark is fully compatible with Hadoop.
  - It can run in YARN, and access data from sources including HDFS, MapR-FS, HBase and HIVE.
- In addition, spark can also use the more general resource manager Mesos.



 Why Apache Spark?A large, bold, dark blue number '4' is centered within a white circle, which is itself centered within a larger light blue square.

### Unified Stack

Builds applications combining different processing models

- Batch
- Streaming
- Interactive Analytics

© 2015 MapR Technologies  MAPR®

7

Spark has an integrated framework for advanced analytics like Graph processing, advanced queries, stream processing and machine learning.

You can combine these libraries into the same application and use a single programming language through the entire workflow





## Why Apache Spark?



### Multi-language support

- Scala
- Python
- Java
- SparkR

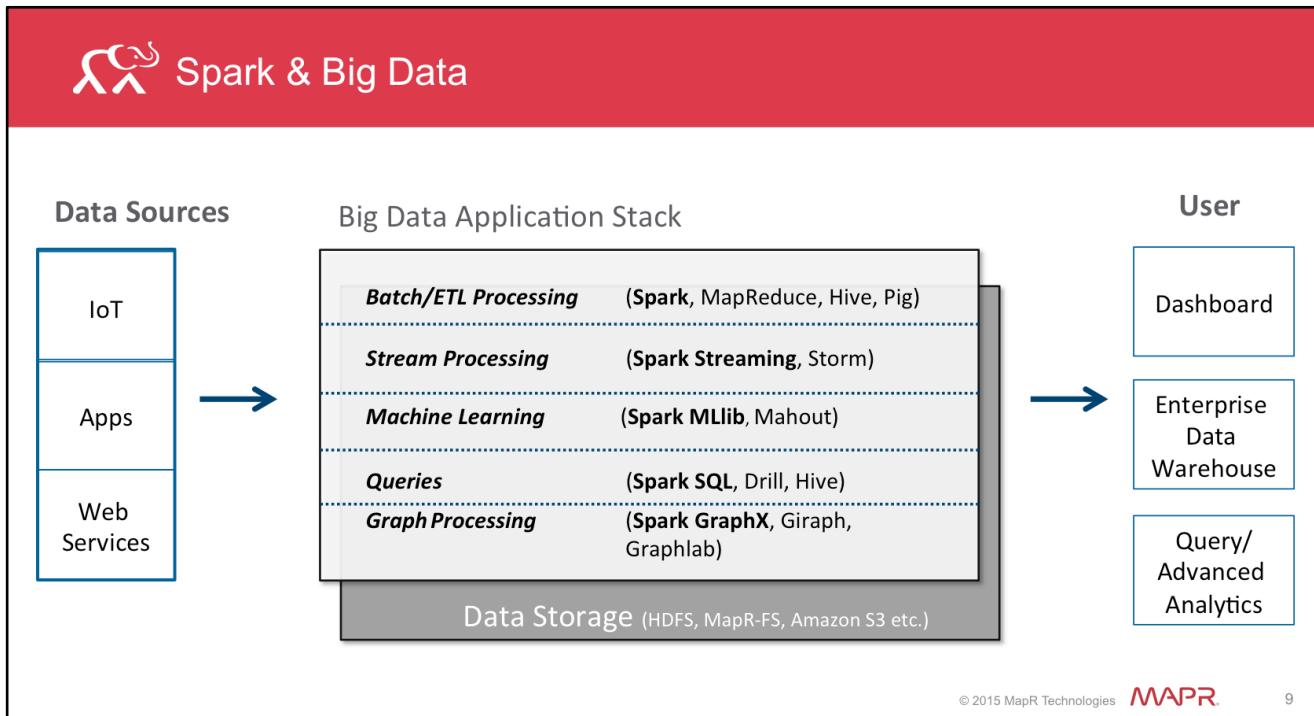
© 2015 MapR Technologies  MAPR®

8

Developers have the choice of using Scala, a relatively new language, Python, Java and as of 1.4.1, SparkR.



8



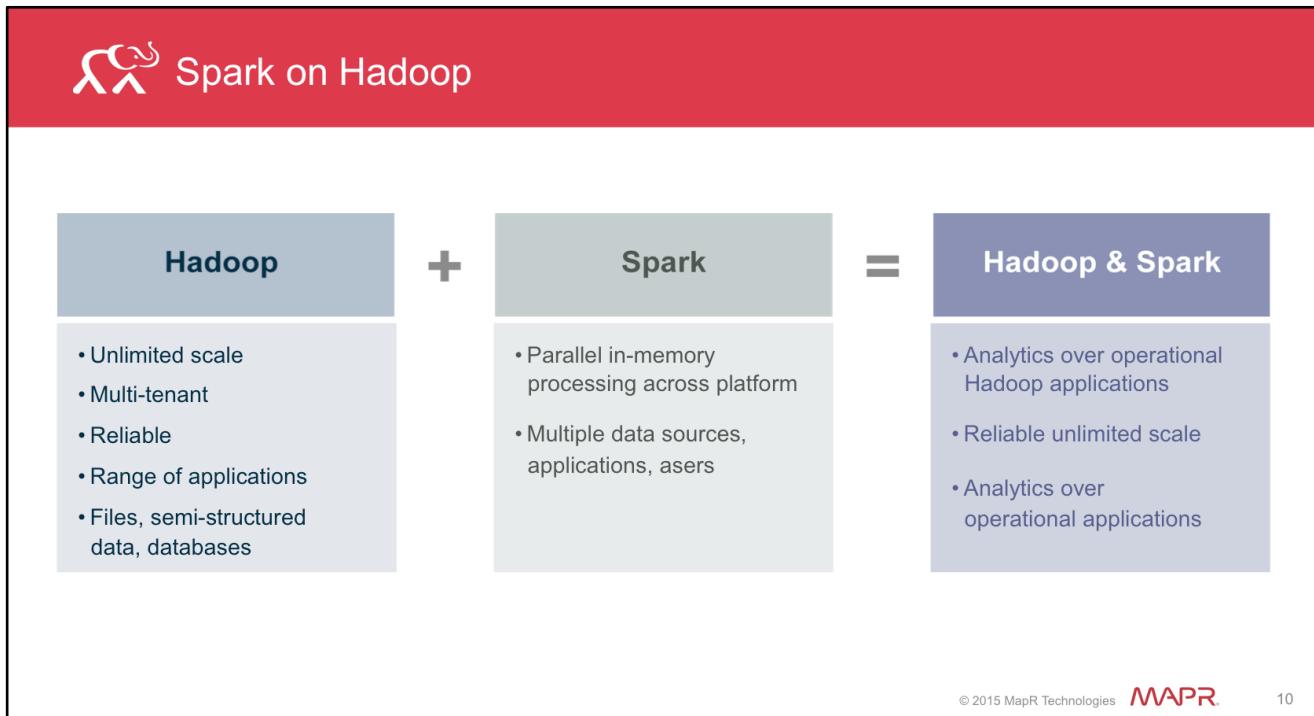
This graphic depicts where Spark fits in the Big Data Application stack.

On the left we see different data sources. There are multiple ways of getting data in using different industry standards such as NFS or existing Hadoop tools.

The stack in the middle represents various Big Data processing workflows and tools that are commonly used. You may have just one of these workflows in your application, or a combination of many. Any of these workflows could read/write to or from the storage layer.

As you can see here, with Spark, you have a unified stack. You can use Spark for any of the workflows.

The output can then be used to create real-time dashboards and alerting systems for querying and advanced analytics; and loading into an enterprise data warehouse.



## Hadoop

Hadoop has grown into a multi-tenant, reliable, enterprise-grade platform with a wide range of applications that can handle files, databases, and semi-structured data.

## Spark

Spark provides parallel, in-memory processing across this platform, and can accommodate multiple data sources, applications and users.

## Hadoop + Spark

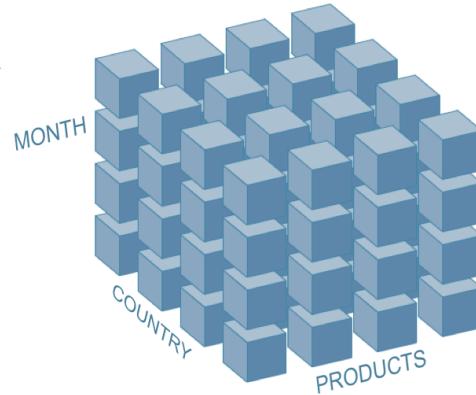
The combination of Spark and Hadoop takes advantage of both platforms, providing reliable, scalable, and fast parallel, in-memory processing. Additionally, you can easily combine different kinds of workflows to provide analytics over your Hadoop and other operational applications.



 Use Case

## OLAP Analytics

- Service provider using MapR & Spark delivers real-time multi-dimensional OLAP analytics
- Must accept data of any type/format
- Perform rigorous analytics across large datasets



© 2015 MapR Technologies 

11

A service provider offers services in customer analytics, technology and decision support to their clients. They are providing real-time multi-dimensional OLAP analytics. They should be able to accept data of any type/format and perform rigorous analytics across datasets that can go up to 1-2 TBs in size.



 Use Case

### Operational Analytics

- Health Insurance company uses MapR to store patient information
- Spark computes patient re-admittance probability
- Real-time analytics over NoSQL
- Spark with MapR-DB



© 2015 MapR Technologies 

12

A health insurance company is using MapR to store patient information, which is combined with clinical records. Using real time analytics over NoSQL, they are able to compute re-admittance probability. If this probability is high enough, additional services, such as in home nursing, are deployed.





## Complex Data Pipelining

- Pharmaceutical company uses Spark on MapR for gene sequencing analysis
- Spark reduces processing from weeks to hours
- Complex machine learning without MapReduce



© 2015 MapR Technologies **MAPR**.

13

A leading pharmaceutical company uses Spark to improve gene sequencing analysis capabilities, resulting in faster time to market. Before Spark, it would take several weeks to align chemical compounds with genes. With ADAM running on Spark, gene alignment only takes a matter of a few hours.



 Knowledge Check**What are the advantages of using Apache Spark?**

1. Compatible with Hadoop
2. Ease of development
3. Fast
4. Multiple Language support
5. Unified stack
6. All of the above

Answer: 6



 Learning Goals

## Describe Features of Apache Spark

- How Spark fits in Big Data ecosystem
  - Why Spark & Hadoop fit together
- ▶ Define Spark Components

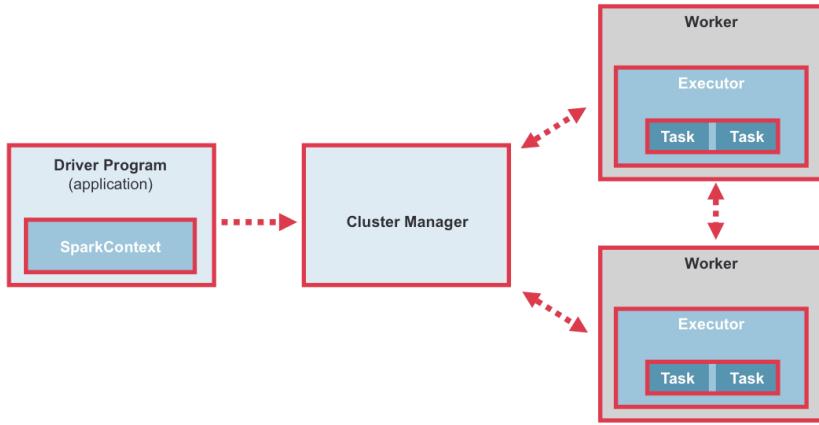
© 2015 MapR Technologies  MAPR®

15

In this section, we take a look at Spark components and libraries.



## Spark Components



© 2015 MapR Technologies **MAPR**

16

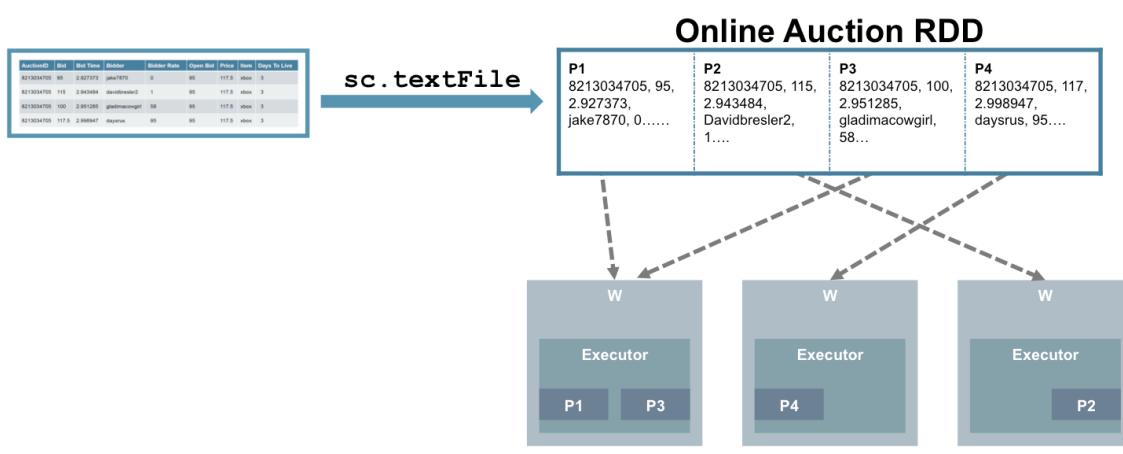
A Spark cluster consists of two processes, a driver program and multiple workers nodes each running an executor process. The driver program runs on the driver machine, the worker programs run on cluster nodes or in local threads.

1. The first thing a program does is to create a `SparkContext` object. This tells Spark how and where to access a cluster
2. `SparkContext` connects to cluster manager. Cluster manager allocates resources across applications
3. Once connected, Spark acquires executors in the worker nodes (an executor is a process that runs computations and stores data for your application)
4. Jar or python files passed to the `SparkContext` are then sent to the executors.
5. `SparkContext` will then send the tasks for the executor to run.
6. The worker nodes can access data storage sources to ingest and output data as needed.





## Spark Resilient Distributed Datasets



© 2015 MapR Technologies

17

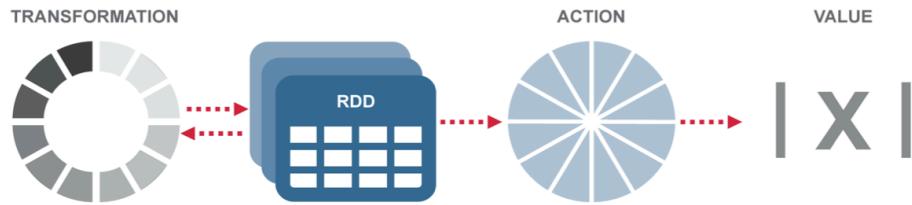
Resilient distributed datasets, or RDD, are the primary abstraction in Spark. They are a collection of objects that is distributed across nodes in a cluster, and data operations are performed on RDD.

- Once created, RDD are immutable.
- You can also persist, or cache, RDDs in memory or on disk.
- Spark RDDs are fault-tolerant. If a given node or task fails, the RDD can be reconstructed automatically on the remaining nodes and the job will complete.





## Spark Resilient Distributed Datasets



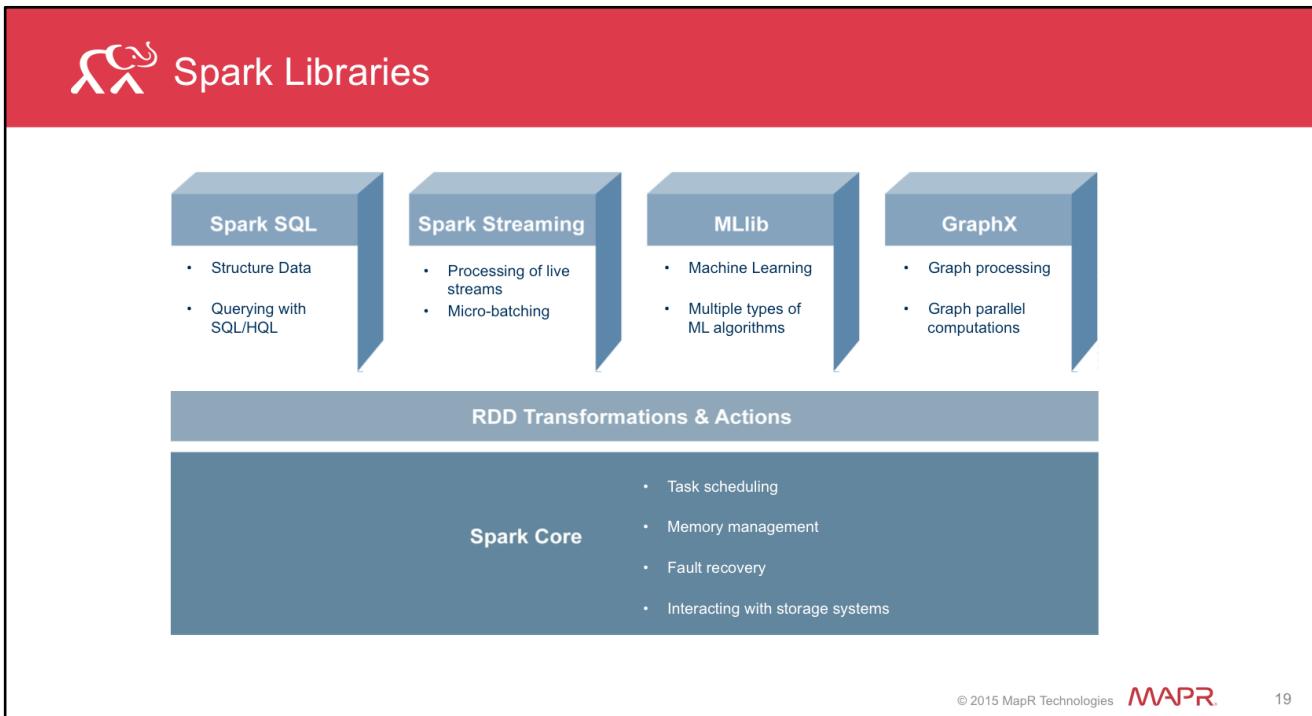
© 2015 MapR Technologies **MAPR**.

18

There are two types of data operations you can perform on an RDD, transformations and actions.

- A transformation will return an RDD. Since RDD are immutable, the transformation will return a new RDD.
- An action will return a value.





The **Spark core** is a computational engine that is responsible for task scheduling, memory management, fault recovery and interacting with storage systems. The Spark core contains the functionality of Spark. It also contains the APIs that are used to define RDDs and manipulate them.

The higher level components are integrated into this stack.

- Spark SQL can be used for working with structured data. You can query this data via SQL or HiveQL. Spark SQL supports many types of data sources such as structured Hive tables and complex JSON data.
- Spark streaming enables processing of live streams of data and doing real-time analytics.
- MLlib is a machine learning library that provides multiple types of machine learning algorithms such as classification, regression, clustering.
- GraphX is a library for manipulating graphs and performing graph-parallel computations.





## Knowledge Check

**Match the following:**

- A. Responsible for task scheduling, memory management
- B. Tells Spark how & where to access a cluster
- C. Collection of objects distributed across many nodes in a cluster

1

**SparkContext**

2

**RDD**

3

**Spark Core**

1-b, 2-c, 3-a



 Next Steps

## Lesson 2

Loading & Inspecting Data

© 2015 MapR Technologies  MAPR®

21

Congratulations! You have completed Lesson 1: Introduction to Apache Spark. Go onto Lesson 2 to learn about loading and inspecting data.



 MAPR Academy

## Apache Spark Essentials

Lesson 2: Load & Inspect Data

© 2015 MapR Technologies  MAPR.

1

In this lesson, we look at loading data into Spark using Resilient Distributed Datasets (RDDs) and DataFrames. We will use transformations and actions to explore and process this data



 Learning Goals

- ▶ Describe the different data sources and formats

Create & use Resilient Distributed Datasets (RDD)

Apply operations to RDDs

Cache intermediate RDD

Create & use DataFrames

© 2015 MapR Technologies  MAPR.

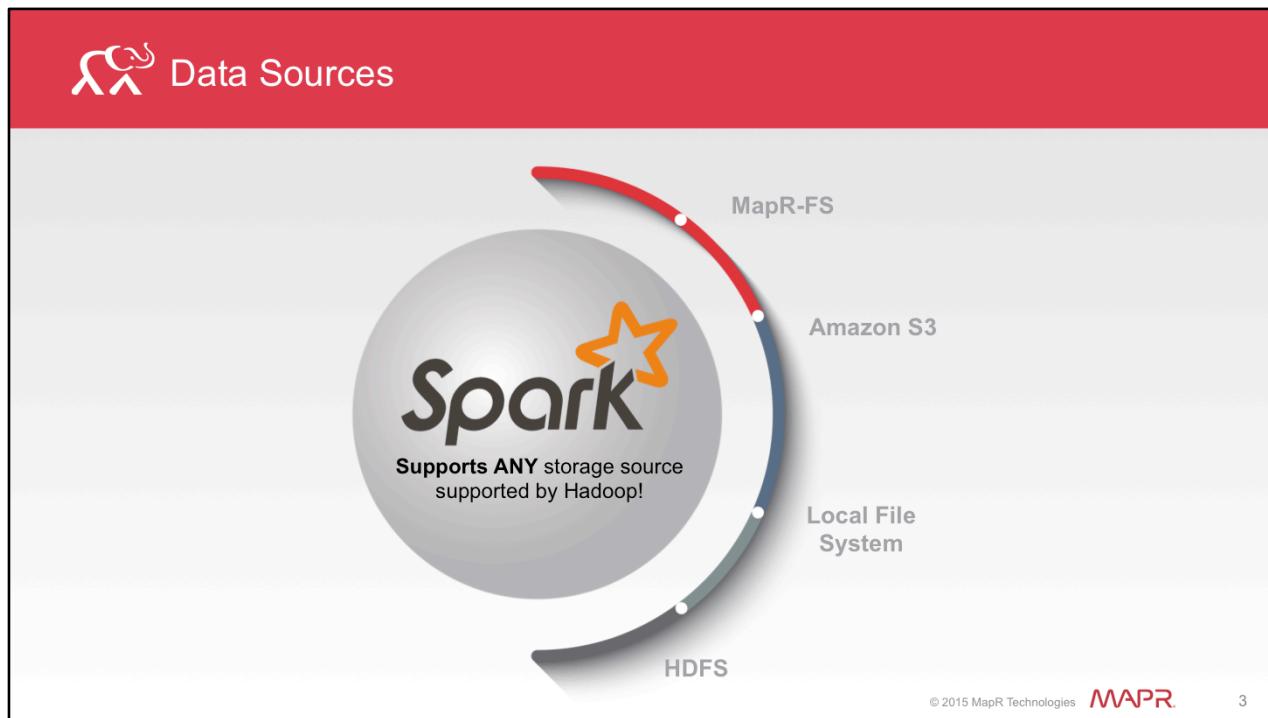
2

At the end of this lesson, you will be able to

- Describe the different data sources and formats available to use with Apache Spark
- Create and use RDDs
- Apply dataset operations to RDDs
- Cache intermediate RDDs
- Use Apache Spark DataFrames

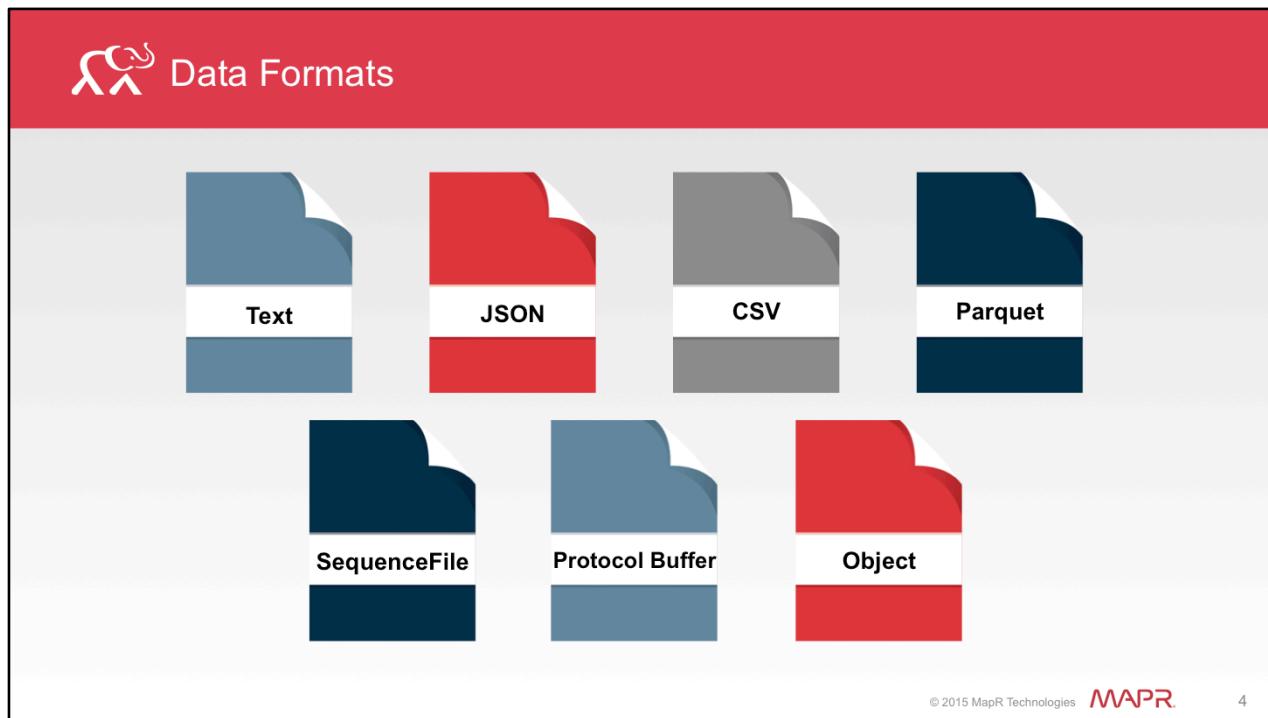
In the first section, we will look at the different data sources available to use with Spark.





Before we load data into Spark, take a look at the data sources and data that Spark supports.

You can load data from any storage source that is supported by Hadoop. You can upload data from the local file system, HDFS, MapR-FS, Amazon S3, Hive, HBase or MapR-DB, JDBC databases, and any other data source that you are using with your Hadoop cluster.



Spark provides wrappers for text files, JSON, CSV, Parquet files, SequenceFiles, protocol buffers and object files. In addition, Spark can also interact with any Hadoop-supported formats.

 Learning Goals

Describe the different data sources and formats

▶ **Create & use Resilient Distributed Datasets (RDDs)**

Apply operations on RDDs

Cache intermediate RDD

Use Spark DataFrames for simple queries

© 2015 MapR Technologies  MAPR.

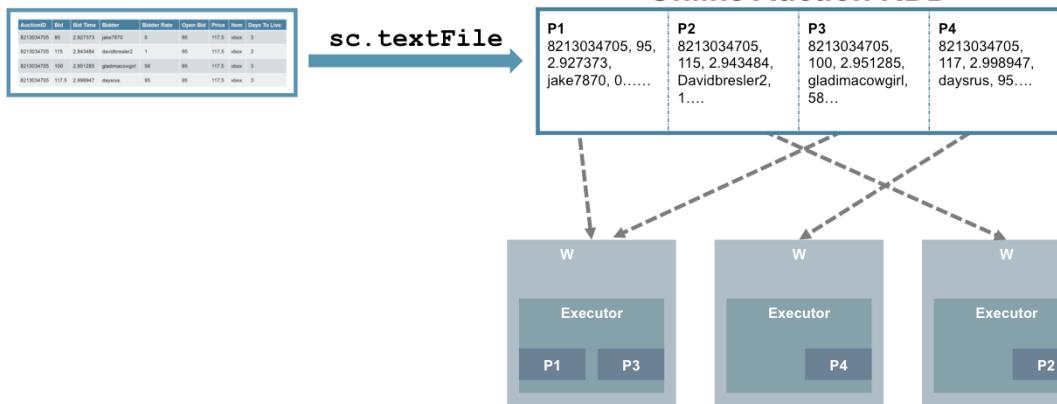
5

In this section, we will take a look at RDDs and load data into a Spark RDD.





## Review: Spark Resilient Distributed Datasets



© 2015 MapR Technologies

6

Resilient distributed datasets are the primary abstraction in Spark. They are a collection of objects that is distributed across many nodes in a cluster. They are immutable once created and cannot be changed.

There are two types of data operations you can perform on an RDD-- transformations and actions.

Transformations are lazily evaluated i.e. They're not computed immediately and will only execute when an action runs on it.

You can also persist, or cache, RDDs in memory or on disk.

Spark RDDs are fault-tolerant. If a given node or task fails, the RDD can be reconstructed automatically on the remaining nodes and the job will complete.





## Scenario: Online Auction Data

AuctionID	Bid	Bid Time	Bidder	Bidder Rate	Open Bid	Price	Item	Days To Live
8213034705	95	2.927373	jake7870	0	95	117.5	xbox	3
8213034705	115	2.943484	davidbresler2	1	95	117.5	xbox	3
8213034705	100	2.951285	gladimacowgirl	58	95	117.5	xbox	3
8213034705	117.5	2.998947	daysrus	95	95	117.5	xbox	3

© 2015 MapR Technologies 

7

We are working with online Auction data that is in a CSV file in the local file system.

The data contains three types of items – xbox, cartier and palm.

Each row in this file represents an individual bid.

Every bid has the auctionID, the bid amount, the bid time in days from when the auction started, the bidder userid, bidder rating, the opening bid for that auction, the final selling price, the item type, and the length of time for the auction in days.





## Scenario: Online Auction Data



**Seller:**  
What should be my starting price to get the highest selling price?



© 2015 MapR Technologies **MAPR**.

8



**Buyer:**  
How can I predict the selling price, so I know what to bid?



These are some common questions to ask of auction data.

- One use case is ad-targeting. Advertising companies auction ad space and companies vie to have their ads seen.
- Another use case is telephony bandwidth where spectrums are auctioned off in a closed bid.

In this course, we are going to load and inspect the data. Later, we will use machine learning code to answer such questions.





You can create an RDD from an existing collection or from external data sources.

When you load the data into Spark, you are creating an RDD.

The first RDD you create is usually called the **input RDD** or **base RDD**.

We are going to use the Spark Interactive Shell to load the auction data into Spark.





## Spark Interactive Shell

- Write programs **INTERACTIVELY!**
- Scala or Python
- **INSTANT** feedback as code is entered
- **SparkContext** initialized on Shell start up



The screenshot shows a terminal window titled "Spark Shell" with the command "spark-shell" entered. The output shows the version of Python being used (2.6.6) and the creation of a SparkContext named "sc". It then reads a file from HDFS, counts the lines, and filters for lines containing the word "holiday".

```
Using Python version 2.6.6 (r266:14292, Sep 11 2012, 00:34:23)
...
>>> sc.textFile("hdfs://ip-172-31-51-234/user/etl-2/cnblogs_20071209-200808_02")
...
>>> file.count()
...
>>> file.filter(lambda line: "holiday" in line).count()
...
0
```

© 2015 MapR Technologies  10

- The Spark interactive shell allows you to write programs interactively.
- It uses the Scala or Python REPL.
- The Interactive shell provides instant feedback as code is entered.
- When the shell starts, SparkContext is initialized and is then available to use as the variable sc. As a reminder, the first thing a program does is create a SparkContext. The SparkContext tells Spark where and how to access the Spark cluster.



 Creating RDD for Auction Data

## 1. Start Spark Interactive Shell



**Scala** /opt/mapr/spark/spark-<version>/bin/spark-shell



**python™** /opt/mapr/spark/spark-<version>/bin/pyspark-shell

## 2. Load data into Spark using `SparkContext.textFile`

```
val auctionRDD = sc.textFile("path to file/auctiondata.csv")
```

© 2015 MapR Technologies 

11

1. First we start the interactive shell. This can be done by running spark-shell from bin in the Spark install directory. You can also use the Python interactive shell.

**NOTE:** During the course of the lesson, we will demonstrate the Scala code. However, in the lab activities, you have the option to use Python, and solutions are provided for both Python and Scala.

2. Once we have started our shell, we next load the data into Spark using the `SparkContext.textFile` method. Note that “sc” in the code refers to `SparkContext`.





## Creating RDD from Different Data Sources

- **Text Files (returns one record per line)**
  - `SparkContext.textFile()`

- **SequenceFiles**
  - `SparkContext.sequenceFile[K,V]`

- **Other Hadoop InputFormats**

- `SparkContext.hadoopRDD`

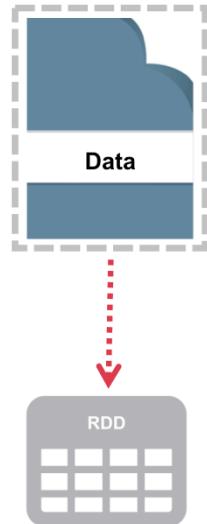
© 2015 MapR Technologies  MAPR.

12

The `textFile` method returns an RDD that contains one record per line. Apart from text files, Spark's Scala API also supports several other data formats:

- `SparkContext.wholeTextFiles` lets you read a directory containing multiple small text files, and returns each of them as (filename, content) pairs. This is in contrast with `textFile`, which would return one record per line in each file.
- For [SequenceFiles, use `SparkContext's sequenceFile\[K, V\]` method where K and V are the types of key and values in the file.](#)
- For other Hadoop InputFormats, you can use the `SparkContext.hadoopRDD` method, which takes an arbitrary `JobConf` and input format class, key class and value class. Set these the same way you would for a Hadoop job with your input source.



 Lazy Evaluation© 2015 MapR Technologies 

13

RDDs are lazily evaluated. We have defined the instructions to create our RDD and defined what data we will use. The RDD has not actually been created yet.

The RDD will be created, the data loaded into it and the results returned only when Spark encounters an action.





## Knowledge Check

**Which of the following is true of the Spark Interactive Shell?**

1. Initializes SparkContext and makes it available
2. Available in Java
3. Provides instant feedback as code is entered
4. Allows you to write programs interactively

Answers: 1,3,4



 Learning Goals

Describe the different data sources and formats

Create & use Resilient Distributed Datasets (RDDs)

▶ **Apply operations on RDDs**

Cache intermediate RDD

Create & use DataFrames

© 2015 MapR Technologies  MAPR.

15

Now that we have created an RDD, we can apply some operations on it. In this section, we will discuss the two types of operations: transformations and actions that can be applied to RDDs.



The slide has a red header bar with the title "Dataset Operations" and a logo. Below the header are two main sections: "Transformations" and "Actions". Each section contains a bulleted list of points. The "Transformations" section includes a circular icon with a gradient and a semi-circular icon below it. The "Actions" section includes a circular icon with radiating lines and a semi-circular icon below it. The slide is numbered 16 in the bottom right corner.

## Dataset Operations

### Transformations

- Creates new dataset from existing one
- Transformed RDD executed only when action runs on it
- Examples: filter(), map(), flatMap()

### Actions

- Return a value to driver program after computation on dataset
- Examples: count(), reduce(), take(), collect()

16

There are two types of operations you can perform on an RDD:

- transformations*, which create a new dataset from an existing one
- and *actions*, which return a value to the driver program after running a computation on the dataset

Transformations are lazily evaluated, which means they are not computed immediately. A transformed RDD is executed only when an action runs on it.

Some examples of Transformations are filter and map

Examples of Actions include count() and reduce().

Let us talk about transformations first.





## Commonly Used Transformations

<b>map</b>	Returns new RDD by applying func to each element of source
<b>filter</b>	Returns new RDD consisting of elements from source on which function is true
<b>groupByKey</b>	Returns dataset (K, Iterable<V>) pairs on dataset of (K,V)
<b>reduceByKey</b>	Returns dataset (K, V) pairs where value for each key aggregated using the given reduce function
<b>flatMap</b>	Similar to map, but function should return a sequence rather than a single item
<b>distinct</b>	Returns new dataset containing distinct elements of source

© 2015 MapR Technologies 

17

Here is a list of some of the commonly used transformations. You can visit the link provided in the resources included with this course for more transformations.





## Scenario: Want all Bids on XBox

```
8214889177,61,6.359155,714ark,15,0.01,90.01,xbox,7  
8214889177,76,6.359294,rjdorman,1,0.01,90.01,xbox,7  
8214889177,90,6.428738,baylorjeep,3,0.01,90.01,xbox,7  
8214889177,88,6.760081,jasonjasonparis,18,0.01,90.01,xbox,7  
8214889177,90.01,6.988831,gpgtpse,268,0.01,90.01,xbox,7  
1638893549,175,2.230949,schadenfreud,0,99,177.5,cartier,3  
1638893549,100,2.600116,chuik,0,99,177.5,cartier,3  
1638893549,120,2.60081,kiwisstuff,2,99,177.5,cartier,3  
1638893549,150,2.601076,kiwisstuff,2,99,177.5,cartier,3
```

© 2015 MapR Technologies  MAPR.

18

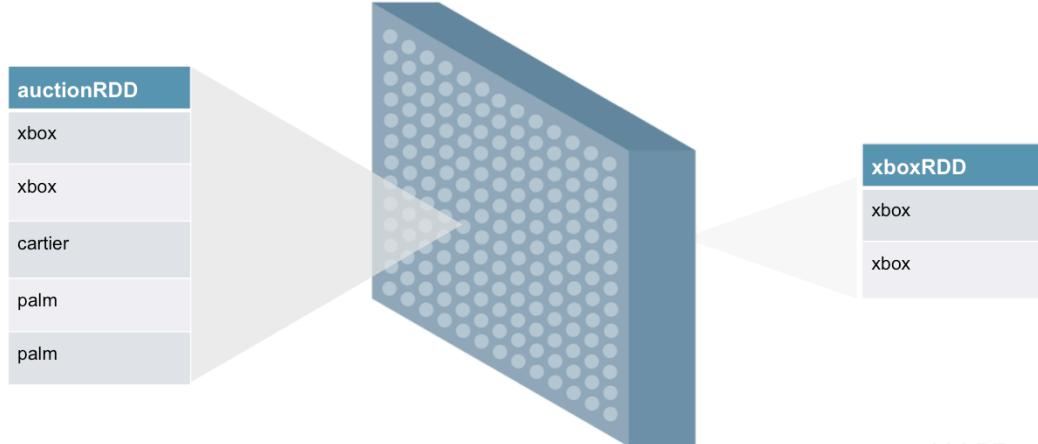
A sample of data in auctiondata.csv is shown here. Each line represents a bid on a particular item such as xbox, cartier, palm, etc. We only want to look at the bids on xbox.

Q. What transformation could we use to get bids only on Xboxes?



 Transformation: filter()

```
val xboxRDD = auctionRDD.filter(line=>line.contains("xbox"))
```



© 2015 MapR Technologies 

19

To do this we can use the filter transformation on the auctionRDD. Each element in the auctionRDD is a line. So we apply the filter to each line of the RDD.

The filter transformation filters out based on the specified condition. We are applying the filter transformation to the auctionRDD where the condition checks to see if the line contains the word “xbox”. If the condition is true, then that line is added to the resulting RDD, in this case, xboxRDD. The filter transformation is using an **anonymous function**.



 Anonymous Function

```
val xboxRDD = auctionRDD.filter(line=>line.contains("xbox"))
```

**"=>"**

Anonymous Syntax



© 2015 MapR Technologies  MAPR.

20

The filter transformation is applying an anonymous function to each element of the RDD which is a line.

This is the Scala syntax for an anonymous function, which is a function that is not a named method.



20

 Anonymous Function

```
val xboxRDD = auctionRDD.filter(line=>line.contains("xbox"))
```

**"line"**  
Input Variable

© 2015 MapR Technologies  MAPR.

21

The `line => line.contains` syntax means that we are applying a function where the input variable is to the left of the `=>` operator. In this case, the input is `line`,



 Anonymous Function

```
val xboxRDD = auctionRDD.filter(line=>line.contains("xbox"))
```

line.contains  
("xbox")  
Condition

© 2015 MapR Technologies  MAPR.

22

The filter will return the result of the code to the right of the function operator. In this example the output is the result of calling line.contains with the condition, does it contain “xbox”.

Anonymous functions can be used for short pieces of code.

We will now take a look at applying actions.





## Commonly Used Actions

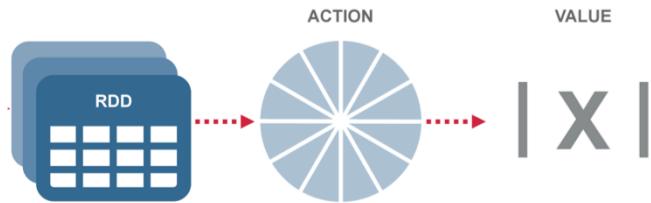
<b>count()</b>	Returns the number of elements in the dataset
<b>reduce(func)</b>	Aggregate elements of dataset using function func
<b>collect()</b>	Returns all elements of dataset as an array to driver program
<b>take(n)</b>	Returns an array with first n elements
<b>first()</b>	Returns the first element of the dataset
<b>takeOrdered(n, [ordering])</b>	Return first n elements of RDD using natural order or custom operator

© 2015 MapR Technologies  MAPR®

23

Here is a list of some of the commonly used actions. You can visit the link provided in the resources included with this course for more actions.



 Actions on RDD

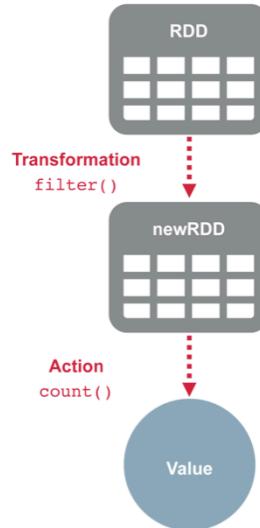
© 2015 MapR Technologies  MAPR.

24

An action on an RDD returns values to the driver program, after running the computation on the dataset.

As mentioned earlier, transformations are lazy. They are only computed when an action requires a result to be returned to the Driver program.



 Actions on RDD

© 2015 MapR Technologies 

25

Here the baseRDD and new RDD are not materialized till we add an action – count().

When we run the command for the action, Spark will read from the text file and create the base RDD. Then the filter transformation is applied to the baseRDD to produce the newRDD. Then count runs on the newRDD sending the value of the total number of elements in the RDD to the driver program.



 Actions on RDD

© 2015 MapR Technologies  MAPR.

26

Once the action has run, the data is no longer in memory.





## Knowledge Check

Match the statements to the appropriate box:

**Actions**

2,4

OR

**Transformations**

1,3,5

1. Returns RDD
2. Returns a value
3. Computed lazily
4. Examples include count, take
5. Examples include filter, map

Actions: 2,4

Transformations: 1,3,5





## Knowledge Check

Match the action to the query:

Query	Action
1. Total # of all bids	A. xboxRDD.take(10)
2. Total # of all bids on Xboxes	B. auctionRDD.count()
3. First 10 bids on Xboxes	C. xboxRDD.count()

© 2015 MapR Technologies  MAPR.

28

1 → b  
2 → c  
3 → a





## Inspecting the Online Auction Data

- How many items were sold?
- How many bids per item type?
- How many different kinds of item types?
- What was the minimum number of bids?
- What was the maximum number of bids?
- What was the average number of bids?

© 2015 MapR Technologies  MAPR.

29

Now that we have looked at the definitions transformations and actions, let us apply these to the base RDD, to inspect the data. We want to find answers to the questions listed here. We will discuss some of these questions in this section, and you will answer the remaining questions in a Lab activity.





## 1. Set up Variables to Map Input

```
val auctionid = 0  
val bid = 1  
val bidtime = 2  
val bidder = 3  
val bidderrate = 4  
val openbid = 5  
val price = 6  
val itemtype = 7  
val daystolive = 8
```

© 2015 MapR Technologies  MAPR.

30

We are setting up the variables to map the input based on our auction data. This is one way to load the data into an RDD in a 2-D array, which makes it easier to reference individual “columns”.



 2. Load the Data

```
val auctionRDD = sc.textFile("/user/user01/data/  
auctiondata.csv").map(_.split(","))  
  
line=>line.split(",")
```

© 2015 MapR Technologies 

31

As we have seen before, we are using the SparkContext method `textFile` to load the data in `.csv` format from the local file system. We use the `map` transformation that applies the `split` function to each line to split it on the `,`.





## Inspect Data: How Many Items Were Sold?

```
val items_sold = auctionRDD  
  .map(bid=>bid(auctionid))  
  .distinct  
  .count
```

### Answer:

628 items were sold

© 2015 MapR Technologies  MAPR.

32

The first question that we want to ask of our data, is how many items were sold.

Each auctionID represents an item for sale. Every bid on that item will be one complete line in the baseRDD dataset, therefore we have many lines in our data for the same item.

The map transformation shown here will return a new dataset that contains the auctionID from each line. The distinct transformation is run on the “auctionID” dataset, and will return another new dataset, which contains the distinct auctionIDs. The count action is run on the “distinct” dataset, and will return the number of distinct auctionIDs.

Each of the datasets that are created by a transformation is temporary, and stored in memory until the final action is complete.





## Inspecting Data: How Many Bids Per Item Type?

```
val bids_item = auctionRDD  
  .map(bid=>(bid(itemtype), 1))  
  .reduceByKey((x,y)=>x+y)  
  .collect()
```

**Answer:**

Array((palm,5917), (cartier,1953), (xbox,2811))

© 2015 MapR Technologies  MAPR.

33

Question 2 – How many bids are there per item type? The item type in our data could be xbox, cartier or palm. Each item type, may have a number of different items.

This line of code is explained in more detail next.

A: Array((palm,5917), (cartier,1953), (xbox,2811))





## Walkthrough: itemtype()

```
val bids_item = auctionRDD  
.map(bid=>(bid(itemtype), 1))
```



TIP:  
To see what the data looks like at this point, use take(1)  
bids\_item.take(1)

**Results:**

```
Array[(String, Int)] = Array((xbox,1), (xbox,1), (xbox,1),  
(xbox,1), (xbox,1).....(cartier,1),(cartier,  
1).....(palm,1),(palm,1).....(palm,1))
```

The function in the first map transformation is mapping each record (bid) of the dataset to an ordered pair (2-tuple) consisting of (itemtype, 1). If we want to see what the data looks like at this point, apply take(1). i.e. bids\_item.take(1).



 Walkthrough: `reduceByKey()`

```
val bids_item = auctionRDD
    .map(bid=>(bid(itemtype),1))
    .reduceByKey((x,y)=>x+y)
```

**Results:**

`Array[(String, Int)] = Array((palm,5917),.....)`



**NOTE:**  
**reduceByKey** not the same as reducer in MapReduce.  
**reduceByKey** – receives two values associated with key.  
reducer in MapReduce receives list of values associated with key

© 2015 MapR Technologies 

35

In this example, `reduceByKey` runs on the key-value pairs from the map transformation – (`itemtype,1`) and aggregates on the key based on the function. The function that we have defined here for the `reduceByKey` transformation is sum of the values i.e `reduceByKey` returns key-value pairs which consists of the sum of the values by key.

Note that `reduceByKey` here does not work the same way as the reducer in MapReduce. `reduceByKey` receives two values associated with a key. The reducer in MapReduce receives a list of values associated with the key.





## Walkthrough: collect()

```
val bids_item = auctionRDD  
    .map(x=>(x(item),1))  
    .reduceByKey((x,y)=>x+y)  
    .collect()
```

**Results:**

Array[(String, Int)] = Array((palm,5917), (cartier,1953), (xbox,2811))

© 2015 MapR Technologies 

36

The collect() action, collects the results and sends it to the driver. We see the result as shown here – i.e. the number of bids per item type.

A: Array((palm,5917), (cartier,1953), (xbox,2811))





## Knowledge Check

Match the Data Operation to the result:

Data Operation	Result
1. <code>auctionRDD.first()</code>	A. Returns the total number of bids
2. <code>auctionRDD.map(func)</code>	B. Returns the func applied to each bid
3. <code>auctionRDD.count()</code>	C. Returns the first bid

1 → c  
2 → b  
3 → a





## Lab 2.1: Load & Inspect Online Auction Data



In this lab, you use the interactive shell to load the online auction data. You use transformations and actions to answer the questions about the data.



 Learning Goals

Describe the different data sources and formats

Create & use Resilient Distributed Datasets (RDD)

Apply operations on RDDs

▶ **Cache intermediate RDD**

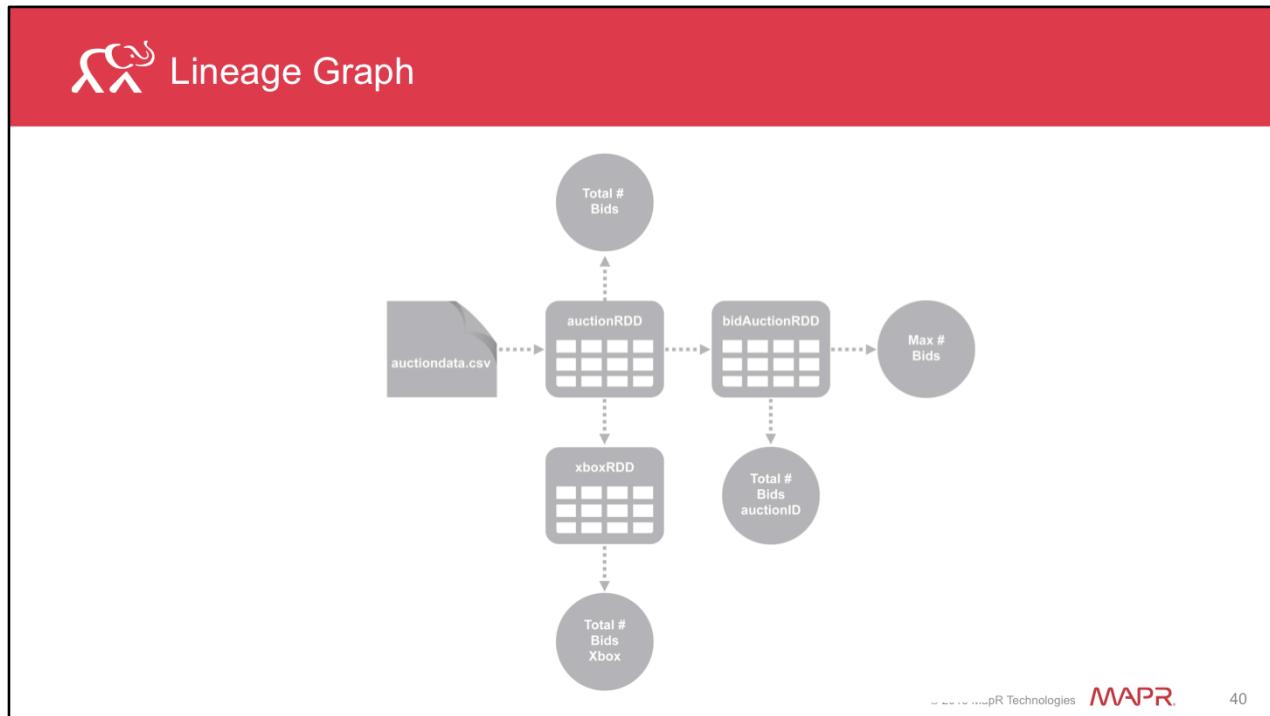
Create & use DataFrames

© 2015 MapR Technologies  MAPR.

39

In the next section, we will discuss reasons for caching an RDD.





© 2014 MapR Technologies MAPR.

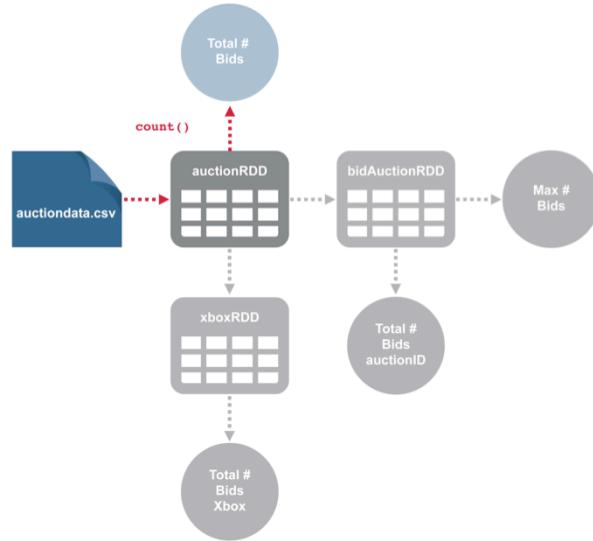
40

As mentioned earlier, RDDs are lazily evaluated. What this means is that even the base RDD (auctionRDD) is not created till we call an action on an RDD.

Every time we call an action on an RDD, Spark will recompute the RDD and all its dependencies.



## Lineage Graph: count()



© 2014 MapR Technologies

MAPR

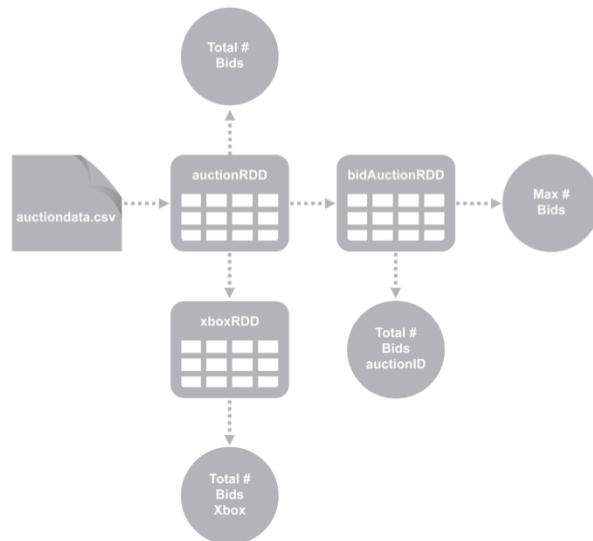
41

For example, when we call the count for the total number of bids, data will be loaded into the auctionRDD and then count action is applied.





## Lineage Graph: count()



© 2014 MapR Technologies

MAPR.

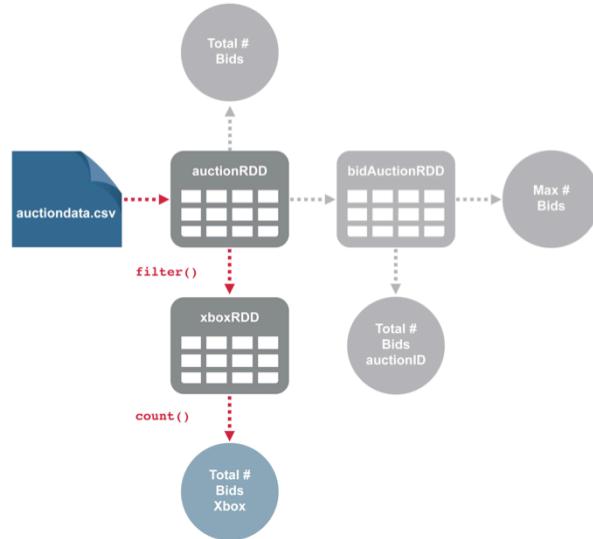
42

The data is no longer in memory.





## Lineage Graph: count()



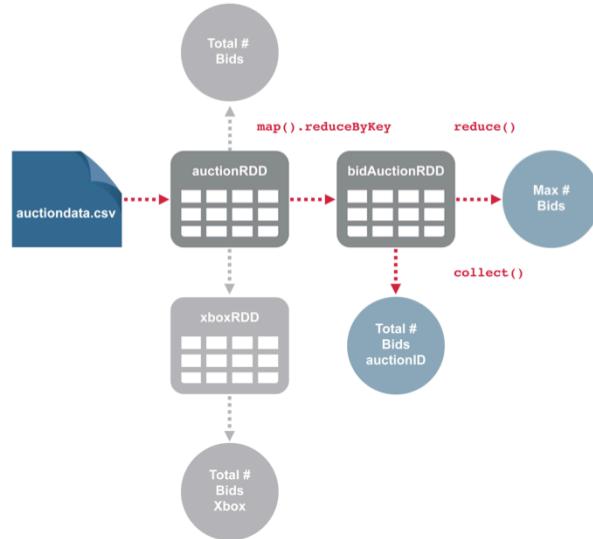
© 2014 MapR Technologies

MAPR

43

Similarly for the count action for total number of bids on xbox will compute the auctionRDD and xboxRDD and then the count.



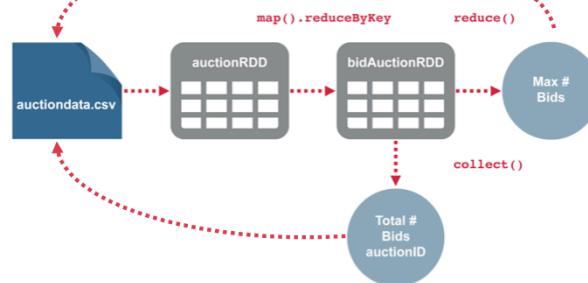
Lineage Graph: `collect()` & `reduce()`

© 2014 MapR Technologies

44

When you call the `collect` for the total bids per auctionid or the `reduce` to find the max number of bids, the data is loaded into the `auctionRDD`, `bidAuctionRDD` is computed and then the respective action operates on the RDD and the results returned. The `reduce` and `collect` will each recompute from the start.



 Why Cache RDD

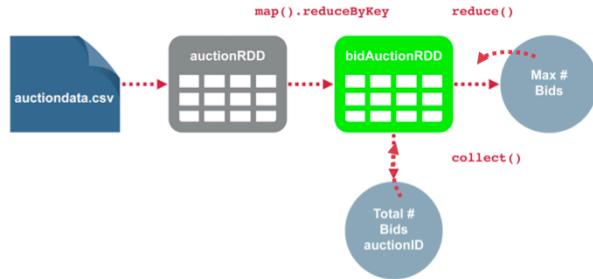
© 2014 MAPR Technologies

MAPR.

45

Each action is called and computed independently. In tabs 3 and 4 in the previous example, when we call `reduce()` or `collect()`, each time the process is completed and the data is removed from memory. When we call one of these actions again, the process starts from the beginning, loading the data into memory. Recomputing every time we run an action can be expensive especially for iterative algorithms.



 Why Cache RDD© 2015 MapR Technologies 

46

When there is a branching in the lineage, and we are using the same RDD multiple times, it is advisable to cache or persist the RDD(s). The RDD has already been computed and the data is already in memory. We can reuse this RDD without using any additional compute or memory resources.

We need to be careful not to cache everything, and only those that are being iterated. The cache behavior depends on the available memory. If the file does not fit in the memory, then the count reloads the data file, and then proceeds as normal.



 Caching the RDD

```
val auctionRDD = sc.textFile("/path/auctiondata.csv")

val bidAuctionRDD = auctionRDD
    .map(x=>(x(auctionid),1))
    .reduceByKey((x,y)=>x+y)

bidAuctionRDD.cache

bidAuctionRDD.collect
```

© 2015 MapR Technologies  MAPR.

47

1. The first line defines the instructions on how to create the RDD, though the file is not, yet read.
2. The second line is says to apply a transformation to the base RDD.
3. The third line says the cache the contents.
4. Again, nothing will happen until you to get the fourth line, the collect action. At this point, the auctiondata.csv is read, the transformation is applied, and the data is cached and collected

The next count or other action, will now just use the data from the cache, rather than re-loading the file and performing the first transformation.





## Knowledge Check

**Which of the following is true of caching the RDD?**

1. When there is branching in lineage, it is advisable to cache the RDD
2. Use rdd.cache() to cache the RDD
3. Cache behavior depends on available memory. If not enough memory, then action will reload from file instead of from cache
4. rdd.persist(MEMORY\_ONLY) is the same as rdd.cache()
5. All of the above

Answers: 5



 Learning Goals

Describe the different data sources and formats available to use with Apache Spark

Create & use Resilient Distributed Datasets (RDDs)

Apply dataset operations on RDDs

Cache intermediate RDD

▶ **Create & use DataFrames**

© 2015 MapR Technologies  MAPR.

49

In this section we are going to look at constructing and using Spark DataFrames.



 What is a Spark DataFrame?

- Programming abstraction in SparkSQL
- Distributed collection of data organized into named columns
- Scales from KBs to PBs
- Supports wide array of data formats & storage systems
- Works in Scala, Python, Java

© 2015 MapR Technologies  MAPR.

50

A DataFrame is distributed collection of data organized into named columns. It scales from KBs to PBs. It supports a wide array of data formats. They can be constructed from structured data files, tables in Hive, external databases or existing RDDs.

The DataFrames API is available in Scala, Python and Java.



 Creating DataFrames

```
//1.Starting point is SQLContext
val sqlContext = new
org.apache.spark.sql.SQLContext(sc)

//2.Used to convert RDD implicitly into a DataFrame
import sqlContext.implicits._

//3. Define schema using a Case class
case class Auction(auctionid: String, bid: Float,
bidtime: Float, bidder: String, biderrate: Int,
openbid: Float, finprice: Float, itemtype: String,
dtl: Int)
```

© 2015 MapR Technologies  MAPR.

51

Spark SQL supports two different methods for converting existing RDDs into DataFrames, reflection and programmatic.

In this example, we are creating a DataFrame from an RDD using the reflection approach.

1. We first need to create a basic SQLcontext, which is the starting point for creating the DataFrame.
2. We then import sqlContext.implicits.\_, since we want to create a DataFrame from an existing RDD.
3. Next, we define the schema for the dataset using a Case class.



 Creating DataFrames**//4. Create the RDD**

```
val auctionRDD=sc.textFile("/user/user01/data/  
ebay.csv").map(_.split(","))
```

**//5. Map the data to the Auctions class**

```
val auctions=auctionRDD.map(a=>Auction(a(0),  
a(1).toFloat, a(2).toFloat, a(3), a(4).toInt,  
a(5).toFloat, a(6).toFloat, a(7), a(8).toInt))
```

© 2015 MapR Technologies  MAPR.

52

4. We create the RDD next using the SparkContext textFile method. In this example, we are creating a DataFrame from an RDD.

5. Once the RDD is created, we map the data to the schema i.e. the Auction class created earlier.



 Creating DataFrames

```
//6. Convert RDD to DataFrame  
val auctionsDF=auctions.toDF()  
  
//7. Register the DF as a table  
auctionsDF.registerTempTable("auctionsDF")
```

© 2015 MapR Technologies  MAPR.

53

6. We then convert the RDD to a DataFrame using the toDF() method.
7. In the last step, we register the DataFrame as a table. Registering it as a table allows us to use it in subsequent SQL statements.

Now we can inspect the data.





## Inspecting Data: Examples

```
//To see the data  
auctionsDF.show()  
  
//To see the schema  
auctionsDF.printSchema()  
  
//Total number of bids  
val totbids=auctionsDF.count()  
  
//Show the columns in the DataFrame  
auctionsDF.columns
```

© 2015 MapR Technologies  MAPR.

54

Here are some examples of using DataFrame functions and actions.





## Commonly Use Actions

<b>collect</b>	Returns an array that contains all of rows in this DataFrame.
<b>count</b>	Returns the number of rows in the DataFrame
<b>describe(cols)</b>	Computes statistics for numeric columns, including count, mean, stddev, min, and max.
<b>first(); head()</b>	Returns first row
<b>show()</b>	Displays the first 20 rows of DataFrame in tabular form
<b>take(n)</b>	Returns the first n rows

© 2015 MapR Technologies  MAPR®

55

This table lists commonly used DataFrame actions. Refer to the link provided for more actions.





## Commonly Used DataFrame Functions

<b>cache()</b>	Cache this DataFrame
<b>columns</b>	Returns all column names as an array
<b>explain()</b>	Only prints the physical plan to the console for debugging purposes
<b>printSchema()</b>	Prints the schema to the console in a tree format
<b>registerTempTable( tableName)</b>	Registers this DataFrame as a temporary table using the given name
<b>toDF()</b>	Returns a new DataFrame

© 2015 MapR Technologies  MAPR.

56

Here are some commonly used DataFrame functions. Refer to the link provided for more functions



 Language Integrated Queries

<b>agg(expr, exprs)</b>	Aggregates on the entire DataFrame without groups
<b>distinct</b>	Returns a new DataFrame that contains only unique rows
<b>filter(conditionExpr)</b>	Filters based on given SQL expression
<b>groupBy(col1, cols)</b>	Groups DataFrame using the specified columns so we can run aggregation on them
<b>select(cols)</b>	Selects a set of columns based on expressions

© 2015 MapR Technologies 

57

This table lists commonly used language integrated queries. Refer to the [link here](#) for more information.



 Inspecting the Data

1. How many bids per item type?
2. What is the max, min and average number of bids per item?
3. How many distinct item types?
4. Get all the bids with closing price over 150?

```
auctionsDF.filter(auctionsDF("price")>150)  
.count()
```

© 2015 MapR Technologies  MAPR.

58

**Can you answer the following?**

For the fourth question: we use filter with a condition (price > 150) and then apply count.





## Knowledge Check

Match the Data Operation to the result:

Query	DataFrame Operation
1. Distinct item types?	A. aDF.count()
2. Total number of bids?	B. aDF.show()
3. See 20 rows in DataFrame?	C. aDF.select("itemtype") .distinct().count()

© 2015 MapR Technologies  59

- 1 → c  
2 → a  
3 → b





## Lab 2.2: Load & Inspect Data - DataFrames



In this lab, you will use the interactive shell to inspect the data in dataframes which is also based on the online auction data.



 Next Steps

## Lesson 3

### Build a Simple Spark Application

© 2015 MapR Technologies  MAPR.

61

Congratulation! You have completed Lesson 2. Go on to Lesson 3 to learn about building a simple Spark application.



 MAPR Academy

## Apache Spark Essentials

Lesson 3: Build a Simple Spark Application

© 2015 MapR Technologies  MAPR.

1

Welcome to Apache Spark Essentials – Lesson 3: Build a Simple Spark Application. In this lesson, we will build a simple standalone Spark application.



 Learning Goals

- ▶ Define the Spark program lifecycle
- ▶ Define the function of SparkContext
- ▶ Describe ways to launch Spark applications
- ▶ Launch a Spark application

© 2015 MapR Technologies  MAPR.

2

At the end of this lesson, you will be able to:

1. Define the Spark program lifecycle
2. Define function of SparkContext
3. Describe ways to launch Spark applications
4. Launch a Spark application





## Lifecycle of a Spark Program

1. Create input RDDs in your driver program
2. Use lazy transformations to define new RDDs
3. Cache() any RDDs that are reused
4. Kick off computations using actions

© 2015 MapR Technologies  MAPR.

3

Create some input RDDs from external data or parallelize a collection in your driver program.

Lazily transform them to define new RDDs using transformations like filter() or map()

Ask Spark to cache() any intermediate RDDs that will need to be reused.

Launch actions such as count() and collect() to kick off a parallel computation, which is then optimized and executed by Spark.





## Lifecycle of a Spark Program – Step 1

1. Create input RDDs in your driver program

```
val auctionRDD = sc  
.textFile("/user/user01/data/auctiondata.csv")  
.map(line=>line.split(","))
```

© 2015 MapR Technologies  MAPR.

4

We did this in the previous lesson where we created an RDD using the SparkContext textFile method to load a csv file.





## Lifecycle of a Spark Program – Step 2

## 2. Use lazy transformations to define new RDDs

```
val auctionRDD = sc  
  .textFile("/user/user01/data/auctiondata.csv")  
  .map(line=>line.split(","))  
val bidsitemRDD = auctionRDD.map(bid=>(bid(itemtype),1))  
.reduceByKey((x,y)=>x+y)
```

© 2015 MapR Technologies  MAPR.

5

We can apply a transformation to the input RDD, which results in another RDD. In this example, we apply a map and a reduceByKey transformation on auctionRDD to create bidsitemRDD





## Lifecycle of a Spark Program – Step 3

## 3. Cache() any RDDs that are reused

```
val auctionRDD = sc  
.textFile("/user/user01/data/auctiondata.csv")  
.map(line=>line.split(","))  
val bidsitemRDD = auctionRDD.map(bid=>(bid(item),1))  
.reduceByKey((x,y)=>x+y)  
bidsitemRDD.cache()
```

© 2015 MapR Technologies  MAPR.

6

When we run the command for an action, Spark will load the data, create the inputRDD, and then compute any other defined transformations and actions.

In the example here, if we are going to apply actions and transformations to the bidsitemRDD, we should cache this RDD so that we do not have to recompute the inputRDD each time we perform an action on bidsitemRDD.





## Lifecycle of a Spark Program – Step 4

## 4. Kick off computations using actions

```
val auctionRDD = sc  
.textFile("/user/user01/data/auctiondata.csv")  
.map(line=>line.split(","))  
val bidsitemRDD = auctionRDD.map(bid=>(bid(item),1))  
.reduceByKey((x,y)=>x+y)  
bidsitemRDD.cache()  
bidsitemRDD.count()
```

© 2015 MapR Technologies  MAPR.

7

Once you cache the intermediate RDD, you can launch actions on this RDD without having to recompute everything from the start.

In this example, we launch the count action on the cached bidsitemRDD.





## Knowledge Check

Place the steps of a Spark application lifecycle in the right order:

1. val burglaries=sfpdRDD.  
filter(line=>line.contains("BURGLARIES"))
2. val sfpdRDD=sc.textFile("/user/user01/data/  
sfpd.csv")
3. val totburglaries=burglaries.count()
4. burglaries.cache()

1

**1,2,4,3**

2

**2,1,4,3**

3

**3,1,4,2**

© 2015 MapR Technologies  MAPR.

8

Answers:



 Learning Goals

Define the Spark program lifecycle

- ▶ Define function of SparkContext

Describe ways to launch Spark applications

Launch a Spark application

In this section, we take a look at the function of SparkContext.

