

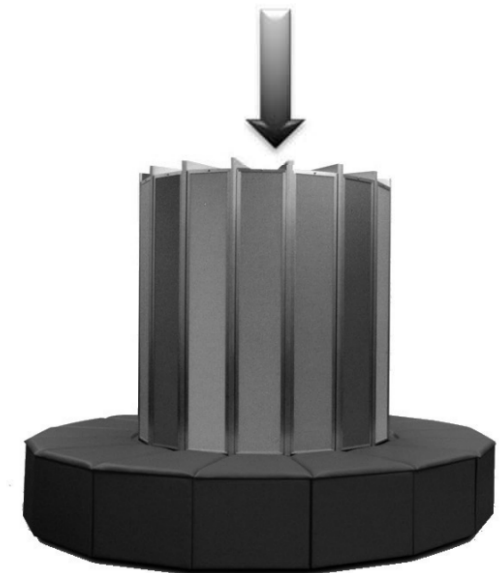
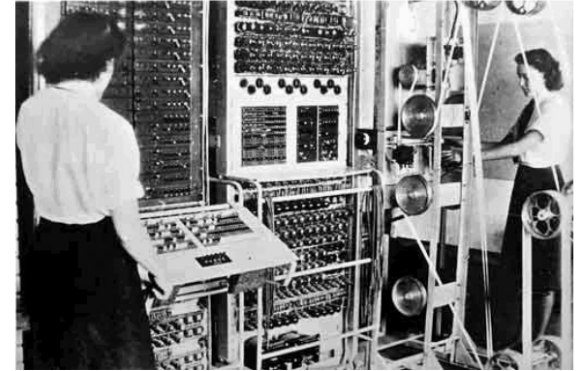
Technologies for Big Data Analysis



PhD. Eng. Iuliana MARIN

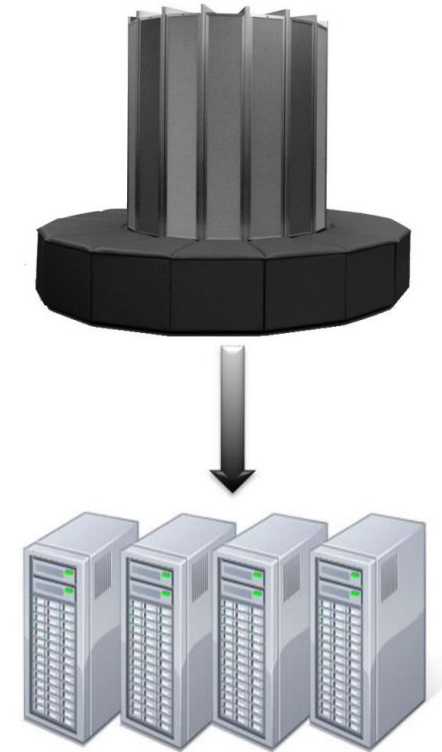
Why is Big Data Needed?

- Traditionally, computation has been processor-bound
 - Relatively small amounts of data
 - A lot of complex processing
- An early solution: bigger computers
 - Faster processor, more memory
 - But it still could not keep up



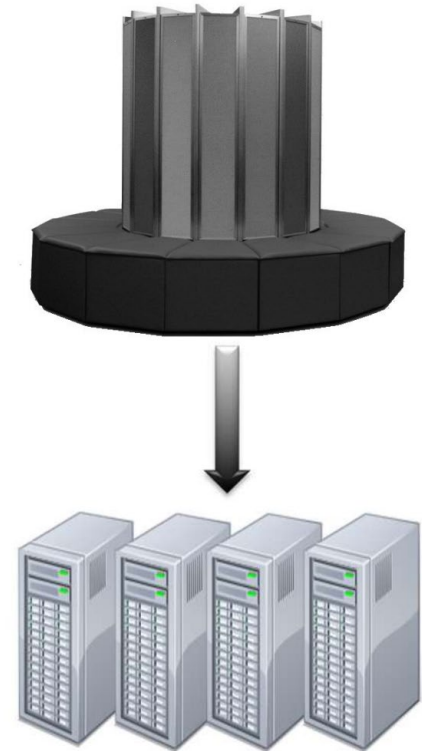
Why is Big Data Needed?

- A better solution: more computers
 - Distributed systems evolved
 - Use multiple machines for a single job
 - MPI (Message Passing Interface) as an example



Why is Big Data Needed?

- Distributed Systems Problems
 - Programming for traditional distributed systems is complex
 - Data exchange requires synchronization
 - Finite bandwidth is available
 - Temporal dependencies are complicated
 - It is difficult to deal with partial failures of the system
 - Ken Arnold, CORBA designer:
 - “Failure is the defining difference between distributed and local programming, so you have to design distributed systems with the expectation of failure”
 - Developers spend more time designing for failure than they do actually working on the problem itself

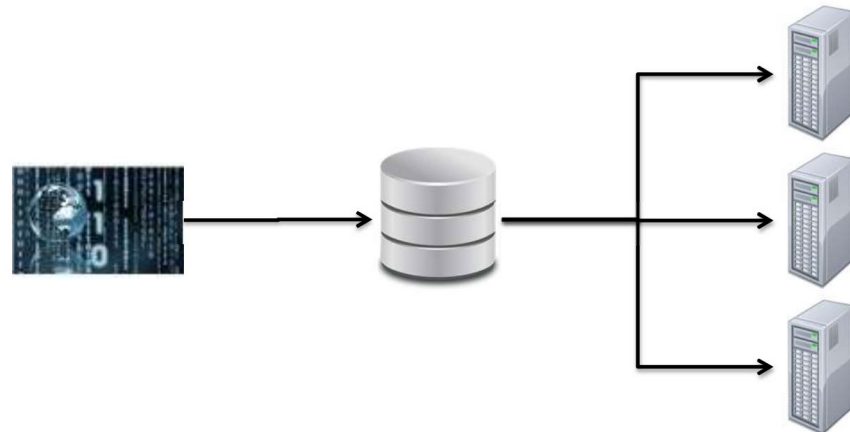


Why is Big Data Needed?

- Moore's Law has held firm for over 40 years
 - Processing power doubles every two years
 - Processing speed is no longer the problem
- Getting the data to the processors becomes the bottleneck
- Quick calculation
 - Typical disk data transfer rate: 75MB/sec
 - Time taken to transfer 100GB of data to the processor: approximately 22 minutes!
 - Assuming sustained reads
 - Actual time will be worse, since most servers have less than 100GB of RAM available
 - A new approach is needed

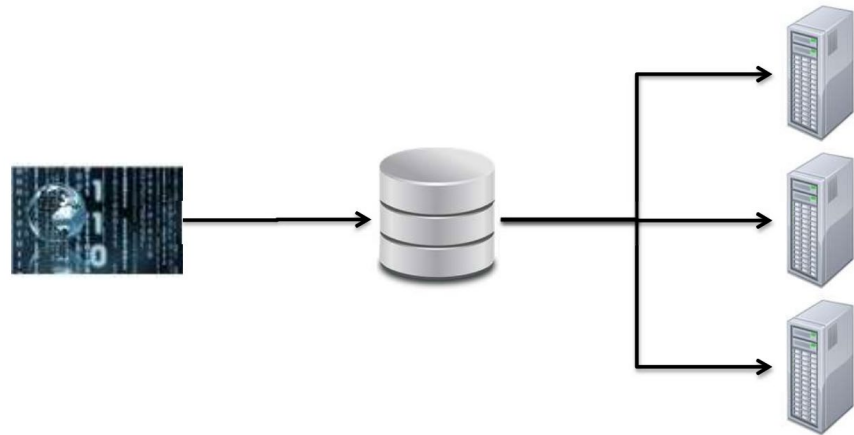
Distributed Systems: The Data Bottleneck

- Traditionally, data is stored in a central location
- Data is copied to processors at runtime
- Fine for limited amounts of data
- Modern systems have much more data
 - terabytes+ per day
 - petabytes+ total
- We need a new approach...



Partial Failure Support

- The system must support partial failure
 - Failure of a component should result in a graceful degradation of application performance
 - Not complete failure of the entire system



Quiz

- A given server has a 1 in 1000 chance of failure on any given day, but you have 500 servers like this in your data center. What is the likely failure of your servers? (choose the best answer.)
 - a) Once in ten days
 - b) Once every two days
 - c) Once every day

Quiz

- A given server has a 1 in 1000 chance of failure on any given day, but you have 500 servers like this in your data center. What is the likely failure of your servers? (choose the best answer.)
 - a) Once in ten days
 - b) Once every two days
 - c) Once every day
- Answer: b. If a given server has a 1 in 1000 chance of failure on any given day, but you have 100 servers like this in your data center, you will likely have a failure once every ten days. And when you “scale out” to 500 servers, the problem actually becomes worse because you’ll have a failure every two days.

Data Recoverability

- If a component of the system fails, its workload should be assumed by still-functioning units in the system
 - Failure should not result in the loss of any data
- If a component of the system fails and then recovers, it should be able to rejoin the system
 - Without requiring a full restart of the entire system
- **Consistency** - Component failures during execution of a job should not affect the outcome of the job
- **Scalability** - Adding load to the system should result in a graceful decline in performance of individual jobs
 - Not failure of the system
 - Increasing resources should support a proportional increase in load capacity

Hadoop's History

- Hadoop is based on work done by Google in the late 1990s/early 2000s
 - Specifically, on papers describing the Google File System (GFS) published in 2003, and MapReduce published in 2004
 - This work takes a radical new approach to the problem of distributed computing
 - Meets all the requirements we have for reliability and scalability
 - Core concept: distribute the data as it is initially stored in the system Individual nodes can work on data local to those nodes
 - No data transfer over the network is required for initial processing

Core Hadoop Concepts

- Applications are written in high-level code
 - Developers need not worry about network programming, temporal dependencies or low-level infrastructure
- Nodes talk to each other as little as possible
 - Developers should not write code which communicates between nodes
 - 'Shared nothing' architecture
- Data is spread among machines in advance
 - Computation happens where the data is stored, wherever possible
 - Data is replicated multiple times on the system for increased availability and reliability

Hadoop: Very High-Level Overview

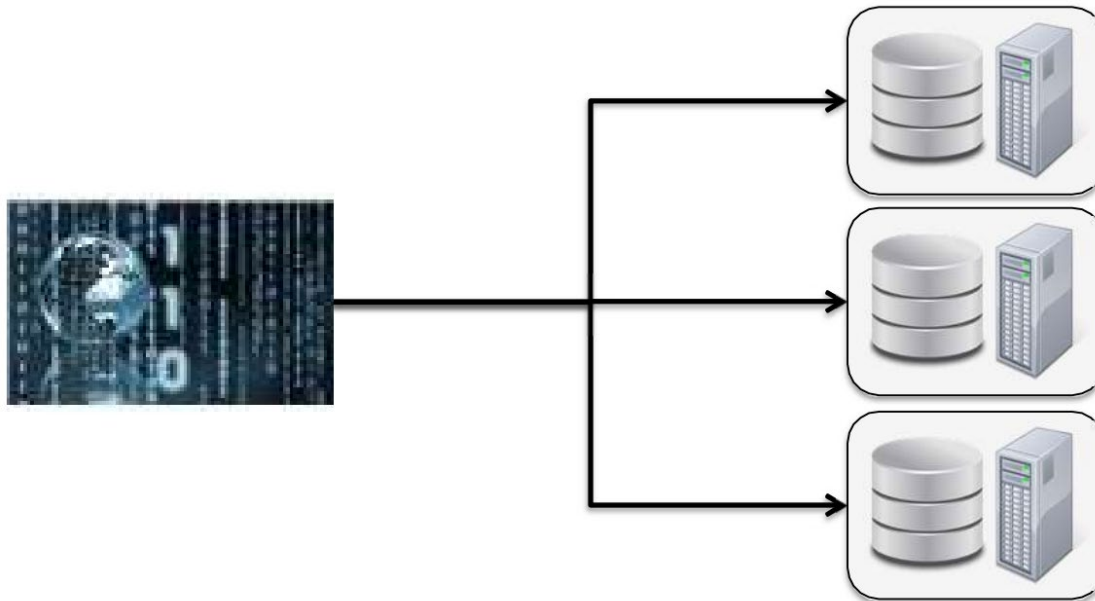
- When data is loaded into the system, it is split into 'blocks'
 - Typically 64MB or 128MB
- Map tasks (the first part of the MapReduce system) work on relatively small portions of data
 - Typically a single block
- A master program allocates work to nodes such that a Map task will work on a block of data stored locally on that node whenever possible
 - Many nodes work in parallel, each on their own part of the overall dataset

Fault Tolerance

- If a node fails, the master will detect that failure and re-assign the work to a different node on the system
- Restarting a task does not require communication with nodes working on other portions of the data
- If a failed node restarts, it is automatically added back to the system and assigned new tasks
- If a node appears to be running slowly, the master can redundantly execute another instance of the same task
 - Results from the first to finish will be used
 - Known as 'speculative execution'

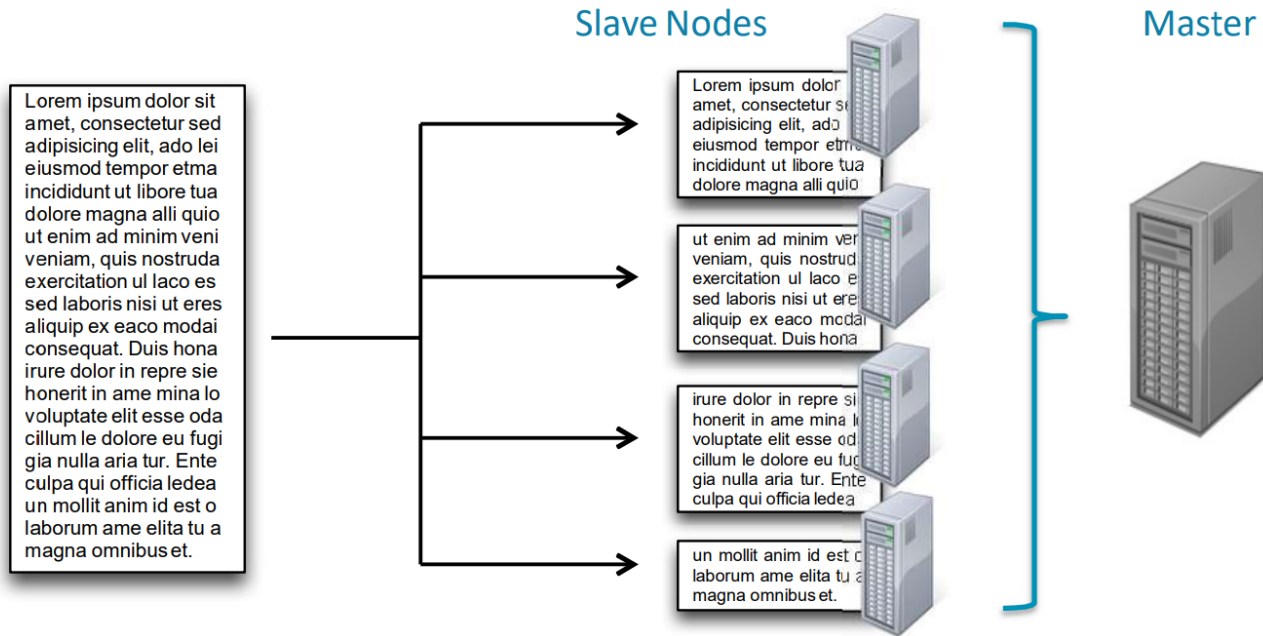
Hadoop

- A radical new approach to distributed computing
 - Distribute data when the data is being stored
 - Run computation where the data is stored



Hadoop: Very High-Level Overview

- Data is split into “blocks” when loaded
- Map tasks typically work on a single block
- A master program manages tasks



Core Hadoop Concepts

- Applications are written in high-level code
- Nodes talk to each other as little as possible
- Data is distributed in advance
 - Bring the computation to the data
- Data is replicated for increased availability and reliability
- Hadoop is scalable and fault-tolerant

Quiz

- If “mean time to failure” for one disk is 3 years = about 1000 days, but you have 100 computers in a cluster with 10 disks per computer, what is the likely disk failure in your cluster? (choose the best answer.)
 - a) Once in ten days
 - b) Once every two days
 - c) Once every day

Quiz

- If “mean time to failure” for one disk is 3 years = about 1000 days, but you have 100 computers in a cluster with 10 disks per computer, what is the likely disk failure in your cluster? (choose the best answer.)
 - a) Once in ten days
 - b) Once every two days
 - c) Once every day

- Answer: c.

100 computers in a cluster

10 disks per computer = 1000 disks

If “mean time to failure” for one disk is 3 years = about 1000 days.

For 1000 disks, Disk failure on average happens once per day.

Big data analysis with Scala and Spark

Apache Spark

- Apache Spark is an Open source analytical processing engine for large scale powerful distributed data processing and machine learning applications.
- Spark is Originally developed at the University of California, Berkeley's, and later donated to Apache Software Foundation.
- In February 2014, Spark became a Top-Level Apache Project and has been contributed by thousands of engineers and made Spark as one of the most active open-source projects in Apache.

Apache Spark

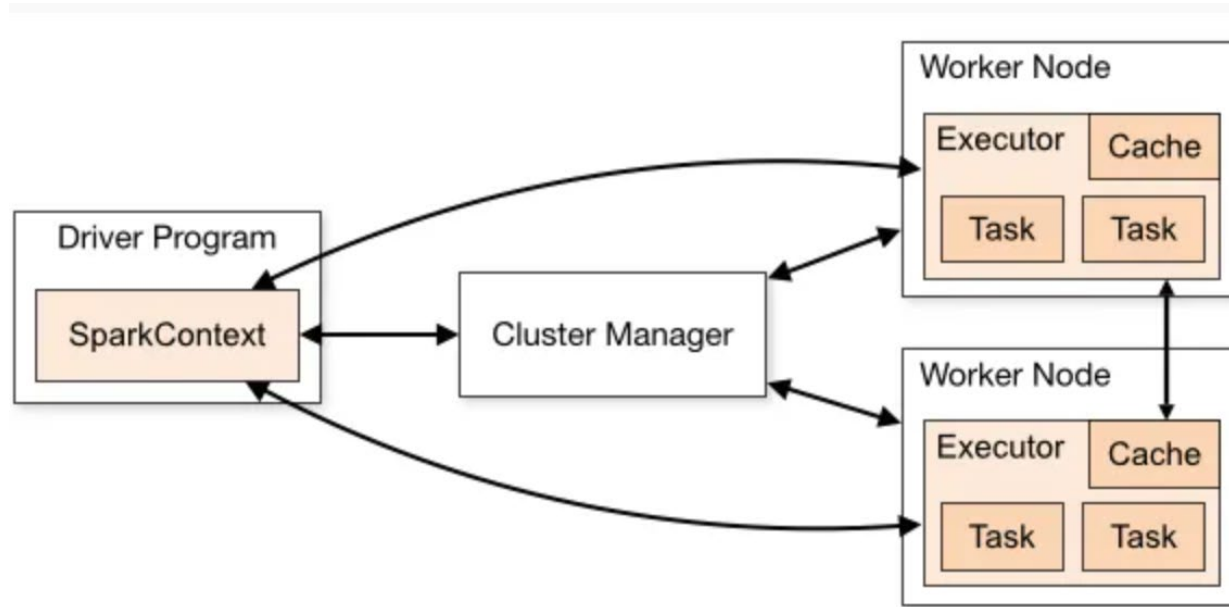
- In-memory computation
- Distributed processing using parallelize
- Can be used with many cluster managers (Spark, Yarn, Mesos)
- Fault-tolerant
- Immutable
- Lazy evaluation
- Cache & persistence
- Inbuild-optimization when using DataFrames
- Supports ANSI SQL

Apache Spark Advantages

- Spark is a general-purpose, in-memory, fault-tolerant, distributed processing engine that allows you to process data efficiently in a distributed fashion.
- Applications running on Spark are 100x faster than traditional systems.
- You will get great benefits using Spark for data ingestion pipelines.
- Using Spark we can process data from Hadoop HDFS, AWS S3, Databricks DBFS, Azure Blob Storage, and many file systems.
- Spark also is used to process real-time data using Streaming and Kafka.
- Using Spark Streaming you can also stream files from the file system and also stream from the socket.
- Spark natively has machine learning and graph libraries.

Apache Spark Architecture

- Apache Spark works in a master-slave architecture where the master is called "Driver" and slaves are called "Workers". When you run a Spark application, Spark Driver creates a context that is an entry point to your application, and all operations (transformations and actions) are executed on worker nodes, and the resources are managed by Cluster Manager.

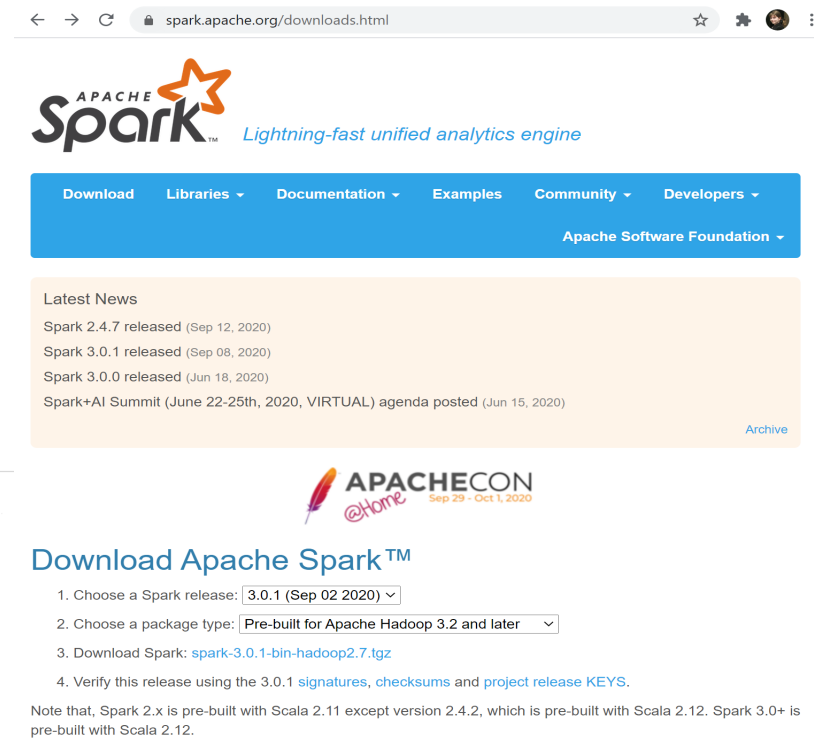
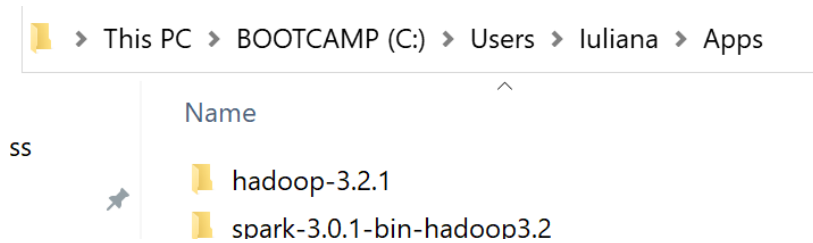


Cluster Manager Types

- Standalone – a simple cluster manager included with Spark that makes it easy to set up a cluster.
- Apache Mesos – Mesos is a Cluster manager that can also run Hadoop MapReduce and Spark applications.
- Hadoop YARN – the resource manager in Hadoop 2. This is mostly used, cluster manager.
- Kubernetes – an open-source system for automating deployment, scaling, and management of containerized applications.

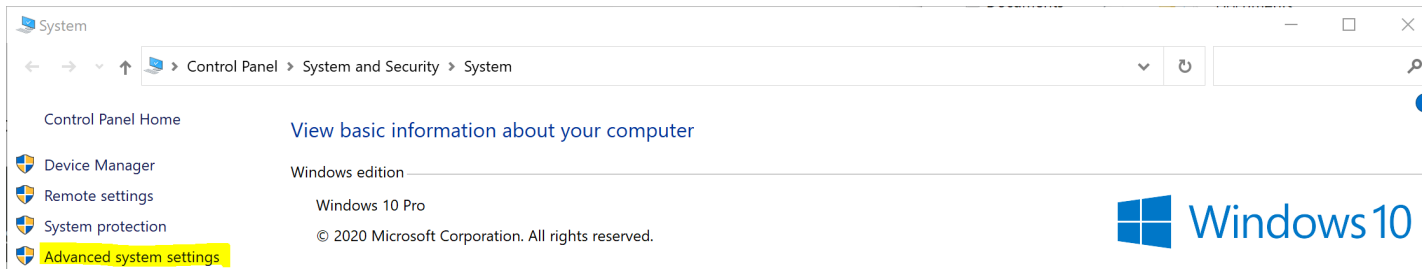
Spark Installation

- Download Apache Spark by accessing Spark Download page and select the link from "Download Spark (point 3)".
- After download, untar the binary using 7zip and copy the underlying folder spark-3.0.1-bin-hadoop3.2 to C:\Users\Iuliana\AppData.
- Install Java JDK 8.



Spark Installation

- You can set the environment variables by going to This PC, right click, selecting Properties and then by clicking on Advanced system settings.



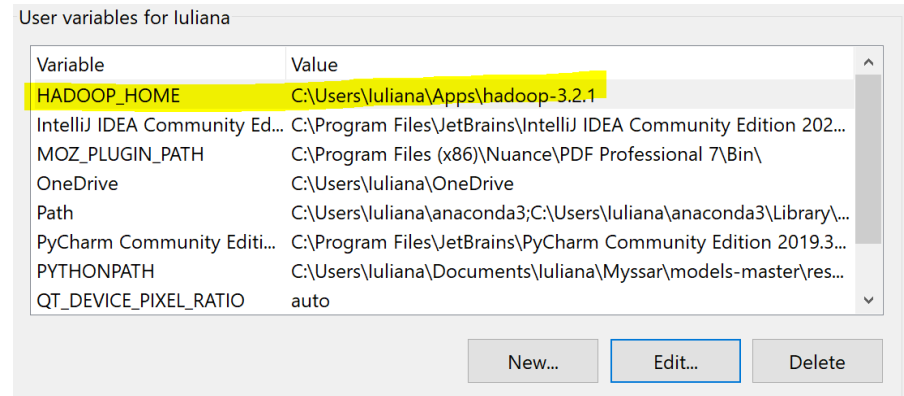
- Click on Environment Variables.
- For System variables, click on New and then add.

Two screenshots illustrating the steps to set environment variables. The left screenshot shows the 'System variables' window with a table of variables. The 'SPARK_HOME' variable is highlighted in blue. The right screenshot shows the 'System Properties' window with the 'Advanced' tab selected and 'Environment Variables...' highlighted in yellow.

Variable	Value
QT_DEVICE_PIXEL_RATIO	auto
SBT_HOME	C:\Program Files (x86)\sbt\
SPARK_HOME	C:\Users\Iuliana\AppData\Local\bin-hadoop3.2;
TEMP	C:\WINDOWS\TEMP
TMP	C:\WINDOWS\TEMP

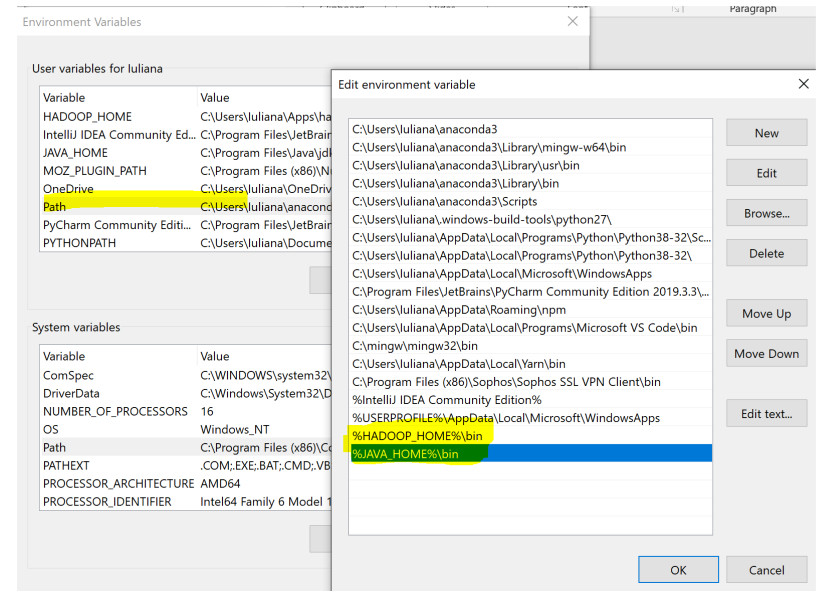
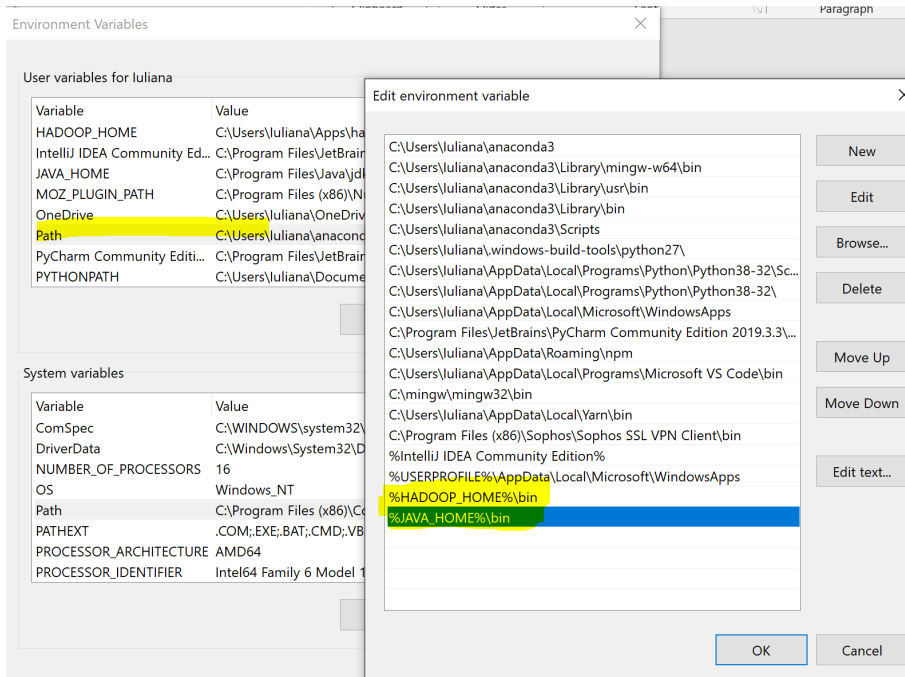
Spark Installation

- Install JDK 8 from Oracle. Download Hadoop 3.2.1 from <https://github.com/cdarlint/winutils/tree/master/hadoop-3.2.1>.
- Unzip it inside the C:\Users\Iuliana\Apps folder. Add the libwinutils.lib, winutils.pdb and winutils files inside the bin folder. Winutils contain the Windows binaries for Hadoop versions.
- Set the environment variable for it.



Spark Installation

- Edit the Path and add the path towards the bin folders for Hadoop and Java:



Spark Installation

- Go to the path where Hadoop was installed and then run winutils.exe.

Command Prompt

```
Microsoft Windows [Version 10.0.19041.685]
(c) 2020 Microsoft Corporation. All rights reserved.

C:\Users\Iuliana>cd C:\Users\Iuliana\AppData\Local\Microsoft\WindowsApps\hadoop-3.2.1\bin

C:\Users\Iuliana\AppData\Local\Microsoft\WindowsApps\hadoop-3.2.1\bin>winutils.exe
Usage: winutils.exe [command] ...
Provide basic command line utilities for Hadoop on Windows.

The available commands and their usages are:

chmod          Change file mode bits.

Usage: chmod [OPTION] OCTAL-MODE [FILE]
       or: chmod [OPTION] MODE [FILE]
Change the mode of the FILE to MODE.

-R: change files and directories recursively
```

Spark Installation

- Test the Hadoop version using the command “hadoop –version”.

```
C:\Users\Iuliana\AppData\Local\hadoop-3.2.1\bin>hadoop -version
java version "1.8.0_241"
Java(TM) SE Runtime Environment (build 1.8.0_241-b07)
Java HotSpot(TM) 64-Bit Server VM (build 25.241-b07, mixed mode)
```

- If an error naming that JAVA_HOME was not found, then change its path into (by modifying Program Files with PROGRA~1):

User variables for Iuliana

Variable	Value
HADOOP_HOME	C:\Users\Iuliana\AppData\Local\hadoop-3.2.1
IntelliJ IDEA Community Ed...	C:\Program Files\JetBrains\IntelliJ IDEA
JAVA_HOME	C:\PROGRA~1\Java\jdk1.8.0_241

Spark Installation

- Start Spark shell, by going to its folder and run the command "spark-shell2".

```
Command Prompt - spark-shell2
Microsoft Windows [Version 10.0.19041.685]
(c) 2020 Microsoft Corporation. All rights reserved.

C:\Users\Iuliana>cd C:\Users\Iuliana\Apps\spark-3.0.1-bin-hadoop3.2
C:\Users\Iuliana\AppData\Local\Programs\Spark\spark-3.0.1-bin-hadoop3.2>cd bin
C:\Users\Iuliana\AppData\Local\Programs\Spark\spark-3.0.1-bin-hadoop3.2\bin>spark-shell2
20/12/16 20:11:10 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
Spark context Web UI available at http://192.168.0.220:4040
Spark context available as 'sc' (master = local[*], app id = local-1608142275687).
Spark session available as 'spark'.
Welcome to

  ____      __
 /    \    /  \
/_    _  /    \
 \___  \/      \
  ___/   \      \
     \___)      /
          \    /
           \  /
            \/

 version 3.0.1

Using Scala version 2.12.10 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_241)
Type in expressions to have them evaluated.
Type :help for more information.
```

Spark Installation

- Press enter inside the Command Prompt and run the following commands to test the environment's functionalities.

```
Command Prompt - spark-shell2
scala> spark
res1: org.apache.spark.sql.SparkSession = org.apache.spark.sql.SparkSession@8872669

scala> sc
res2: org.apache.spark.SparkContext = org.apache.spark.SparkContext@3ca9e72d

scala> val rdd = spark.sparkContext.parallelize(Seq(1,2,3))
rdd: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[0] at parallelize at <console>:23

scala> rdd.count
res3: Long = 3

scala> rdd.foreach(println)
1
2
3

scala>
```

Spark and JSON Files

- Create a JSON file inside the bin folder of Spark. The file will contain several employees.

The screenshot shows a Windows File Explorer window with the address bar set to `C:\Users\Iuliana\AppData\Local\Microsoft\WindowsApps\spark-3.0.1-bin-hadoop3.2\bin`. The file list contains:

Name	Date modified	Type
beeline	8/28/2020 12:22 PM	File
beeline	8/28/2020 12:22 PM	Windows Command ...
docker-image-tool	8/28/2020 12:22 PM	Shell Script
employee	12/16/2020 11:13 PM	JSON File

Below the File Explorer, a Notepad++ window is open to `C:\Users\Iuliana\AppData\Local\Microsoft\WindowsApps\spark-3.0.1-bin-hadoop3.2\bin\employee.json`. The content of the file is:

```
1 {"id" : "1201", "name" : "ion", "age" : "25"}
2 {"id" : "1202", "name" : "dana", "age" : "28"}
3 {"id" : "1203", "name" : "ana", "age" : "39"}
4 {"id" : "1204", "name" : "bogdan", "age" : "23"}
5 {"id" : "1205", "name" : "david", "age" : "23"}
```

SQL Management

- Use the following command to create the SQLContext. Use the following command to read the JSON document named employee.json.

```
scala> val sqlContext = new org.apache.spark.sql.SQLContext(sc)
```

- The data is shown as a table with the fields – id, name, and age.

```
scala> val dfs = sqlContext.read.json("employee.json")  
dfs: org.apache.spark.sql.DataFrame = [age: string, id: string ... 1 more field]
```

```
scala> dfs.show()  
-----+  
|age| id| name|  
-----+  
| 25|1201| ion|  
| 28|1202| dana|  
| 39|1203| ana|  
| 23|1204| bogdan|  
| 23|1205| david|  
-----+
```

SQL Management

- The structure (schema) of the DataFrame can be seen using the following command:

```
scala> dfs.printSchema()
root
 |-- age: string (nullable = true)
 |-- id: string (nullable = true)
 |-- name: string (nullable = true)
```

- Get the values of the name column which belong to the DataFrame.

```
scala> dfs.select("name").show()
+-----+
| name |
+-----+
| ion  |
| dana|
| ana  |
| bogdan|
| david|
+-----+
```

SQL Management

- A filter can be applied based on age to find the employees who have the age greater than 24.

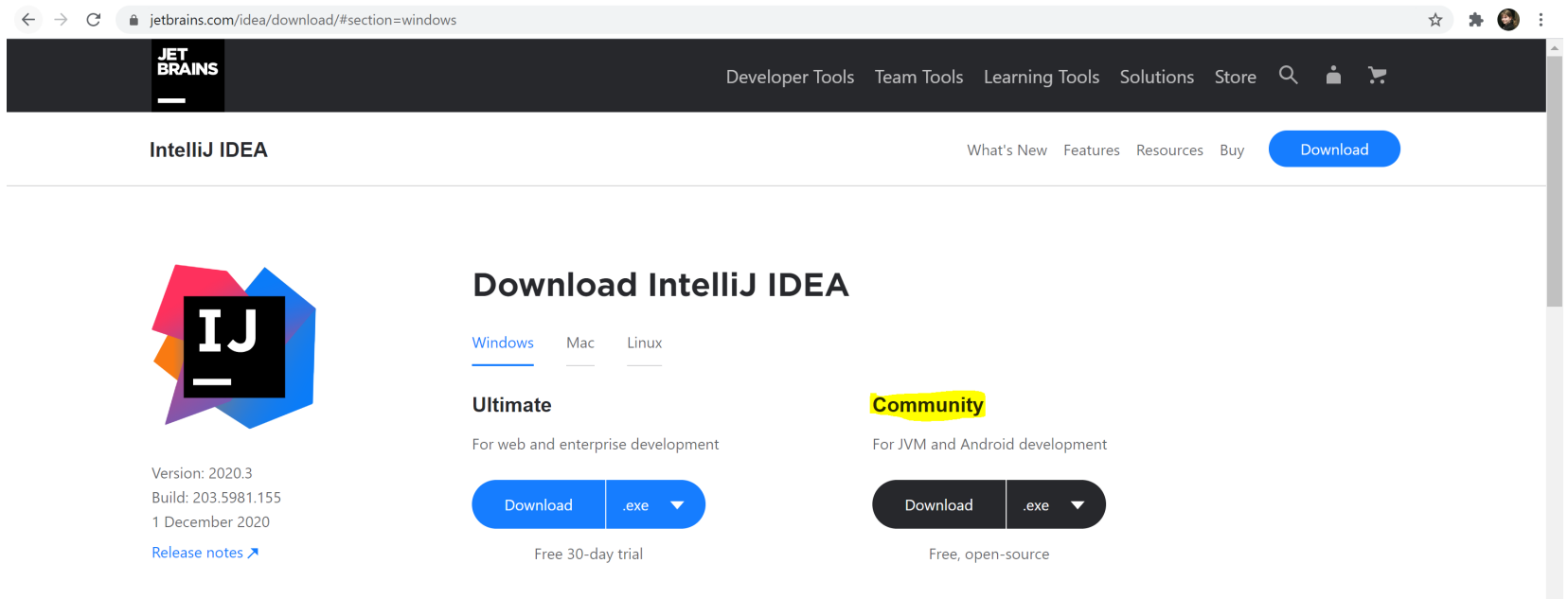
```
scala> dfs.filter(dfs("age") > 24).show()
+---+-----+---+
|age|  id|name|
+---+-----+---+
| 25|1201| ion|
| 28|1202| dana|
| 39|1203| ana|
+---+-----+---+
```

- GroupBy is used to count the number of employees who have the same age.

```
scala> dfs.groupBy("age").count().show()
+---+-----+
|age|count|
+---+-----+
| 28|    1|
| 23|    2|
| 25|    1|
| 39|    1|
+---+-----+
```

IntelliJ IDEA Community Installation

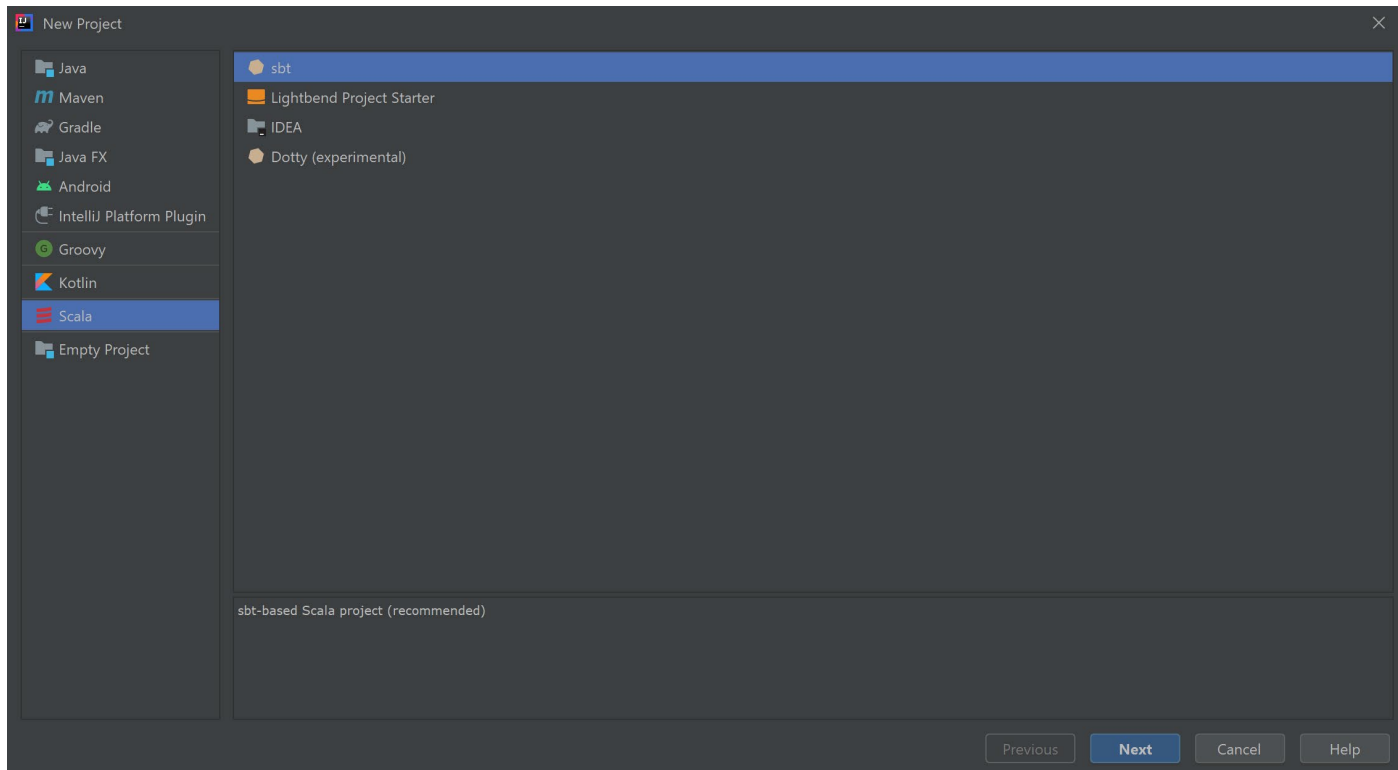
- Install IntelliJ IDEA Community edition to run Spark applications written in Scala due to its good Scala compatibility.



The screenshot shows the JetBrains website's download page for IntelliJ IDEA. The browser address bar displays 'jetbrains.com/idea/download/#section=windows'. The page header includes the JetBrains logo and navigation links for Developer Tools, Team Tools, Learning Tools, Solutions, and Store. A 'Download' button is visible in the top right. The main content area features the IntelliJ IDEA logo on the left, with version information: Version: 2020.3, Build: 203.5981.155, and release date: 1 December 2020. Below the logo is a 'Release notes' link. The central section is titled 'Download IntelliJ IDEA' and offers three operating system options: Windows (selected), Mac, and Linux. Two editions are presented: 'Ultimate' (for web and enterprise development, available as a free 30-day trial) and 'Community' (for JVM and Android development, available as free, open-source). Each edition has a 'Download .exe' button.

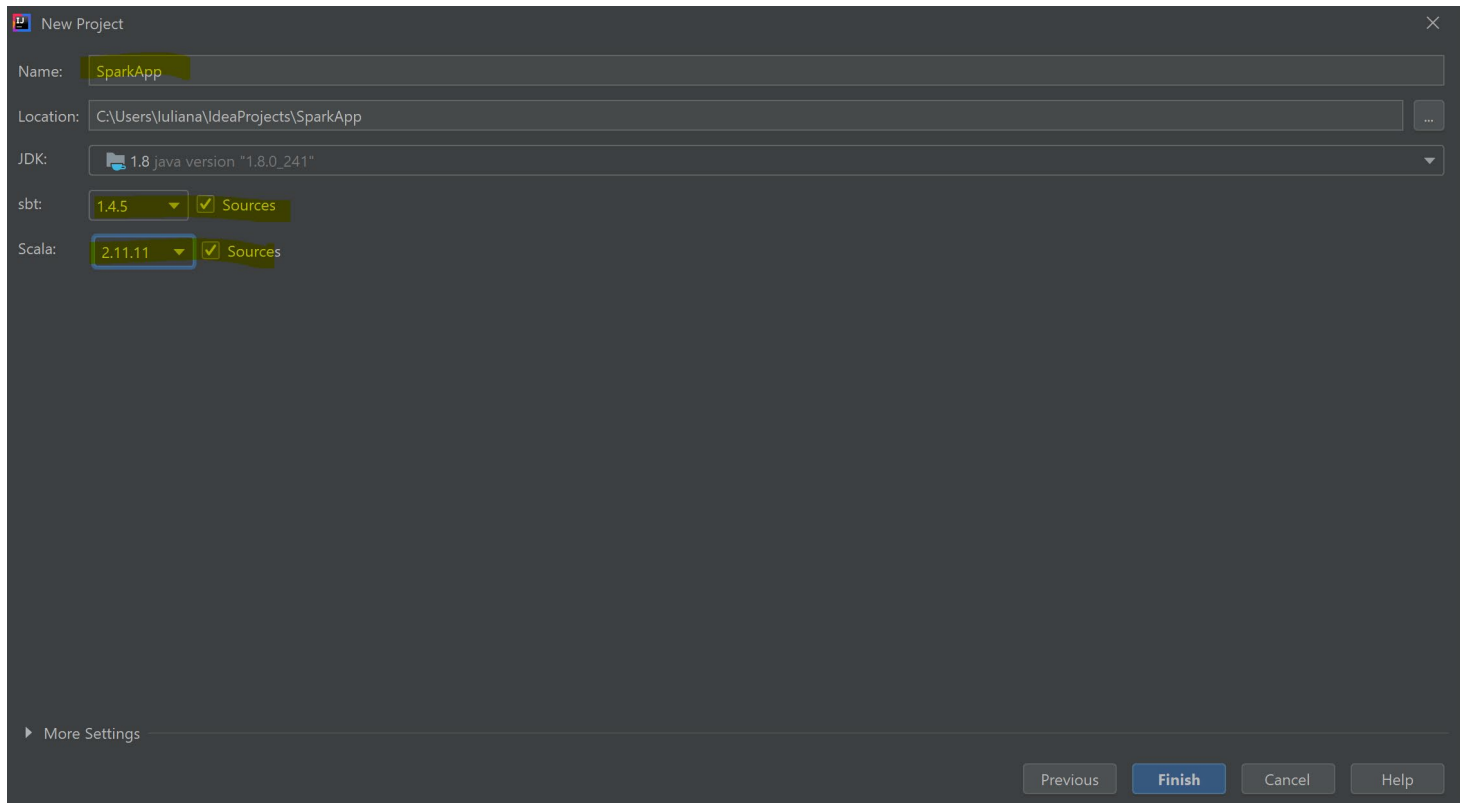
Spark Installation using SBT

- Install sbt using the official website.
- Create a new project by going to File -> New Project and the select Scala and sbt. Afterwards, click on Next.



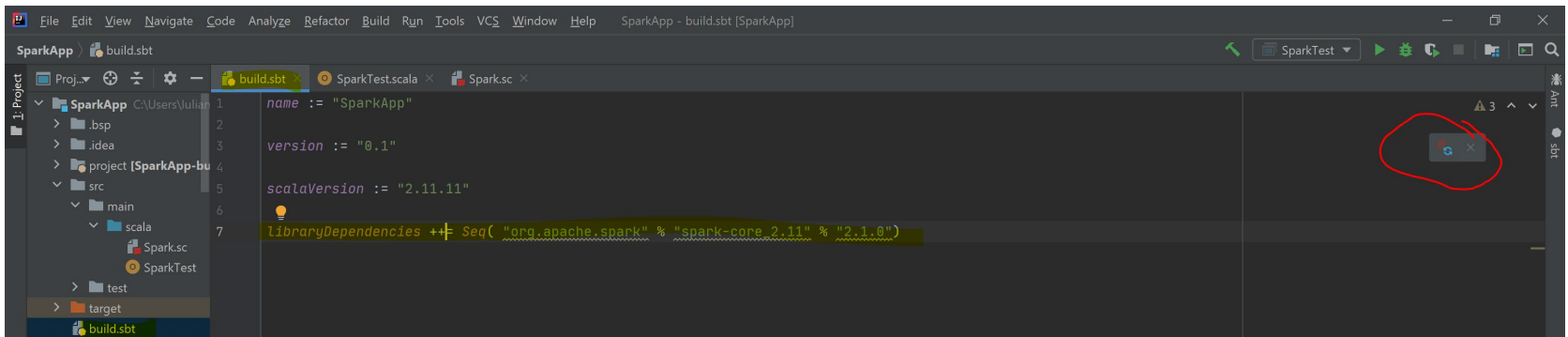
Spark Installation using SBT

- Type the name of the project and click the sbt checkbox for sources, as well as the one for Scala.



Spark Installation using SBT

- Open the build.sbt file and type the last row to include library dependencies for Spark.
- Click on the logo which is depicted in the picture with a red circle and let the program install the needed library.



MapReduce in Scala

- Open the src folder until you reach the scala folder. Right click on it, select New and choose Scala Worksheet. A file with the extension .sc will be created.
- Add the following MapReduce code which will display the top 10 words that appear in the copyright file of Java JDK.

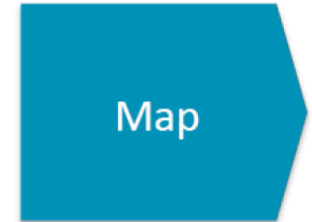
```
import org.apache.log4j.{Level, Logger}
import org.apache.spark.SparkContext

Logger.getRootLogger.setLevel(Level.INFO)
val sc = new SparkContext("local[*]", "SparkDemo")
val lines = sc.textFile("C:\\Program Files\\Java\\jdk1.8.0_241\\COPYRIGHT");
val words = lines.flatMap(line => line.split(' '))
val wordsKVRdd = words.map(x => (x,1))
val count = wordsKVRdd.reduceByKey((x,y) => x + y).map(x =>
(x._2,x._1)).sortByKey(false).map(x => (x._2, x._1)).take(10)
count.foreach(println)
```

Hadoop Components: MapReduce

■ The Mapper

- Each Map task (typically) operates on a single HDFS block
- Map tasks (usually) run on the node where the block is stored



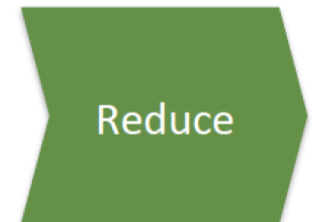
■ Shuffle and Sort

- Sorts and consolidates intermediate data from all mappers
- Happens as Map tasks complete and before Reduce tasks start



■ The Reducer

- Operates on shuffled/sorted intermediate data (Map task output)
- Produces final output



Hadoop Components: MapReduce

Example: Word Count

Input Data

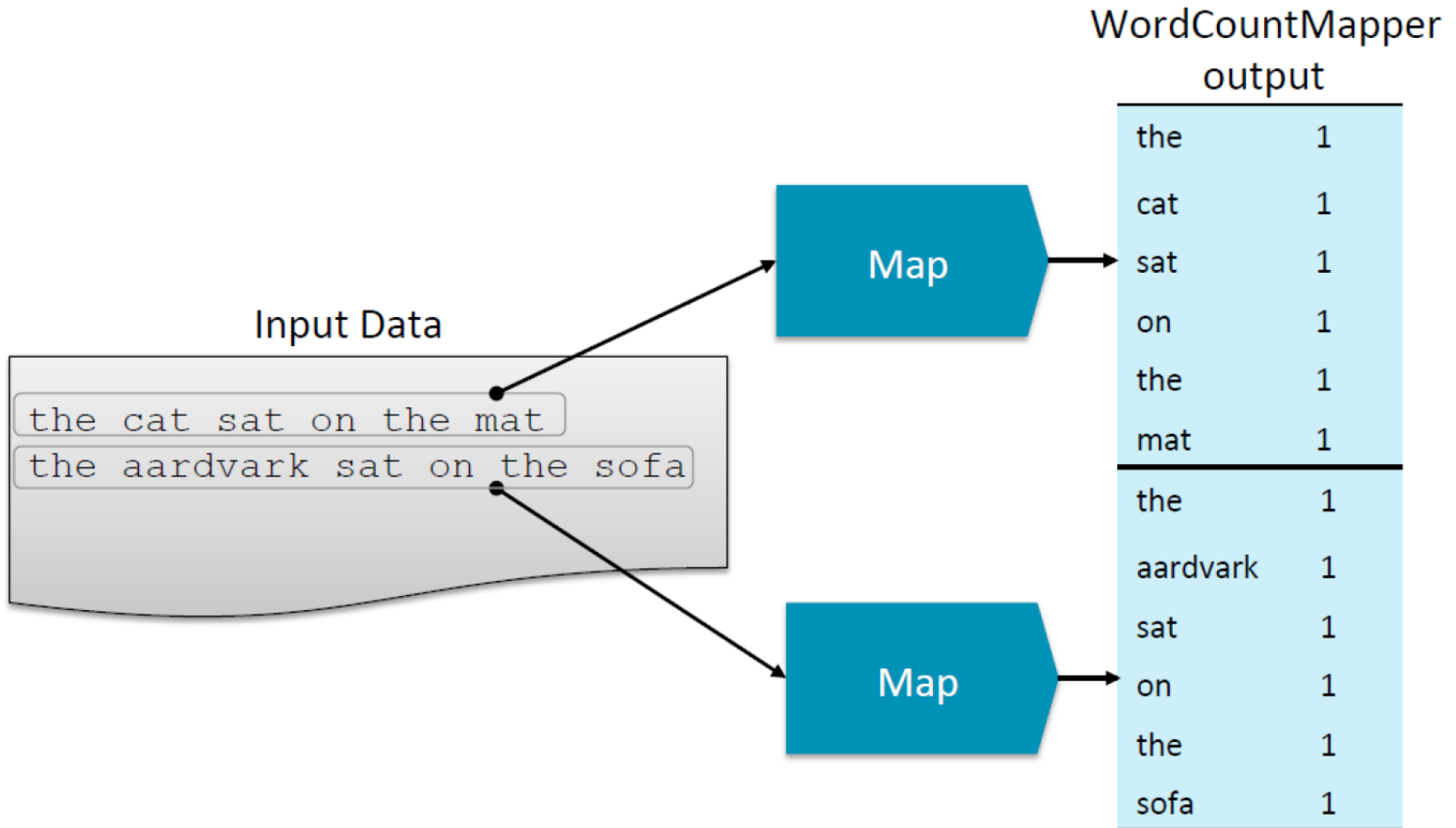
```
the cat sat on the mat  
the aardvark sat on the sofa
```



Result

aardvark	1
cat	1
mat	1
on	2
sat	2
sofa	1
the	4

Example: Word Count Mapper



Example: Shuffle & Sort

Mapper Output

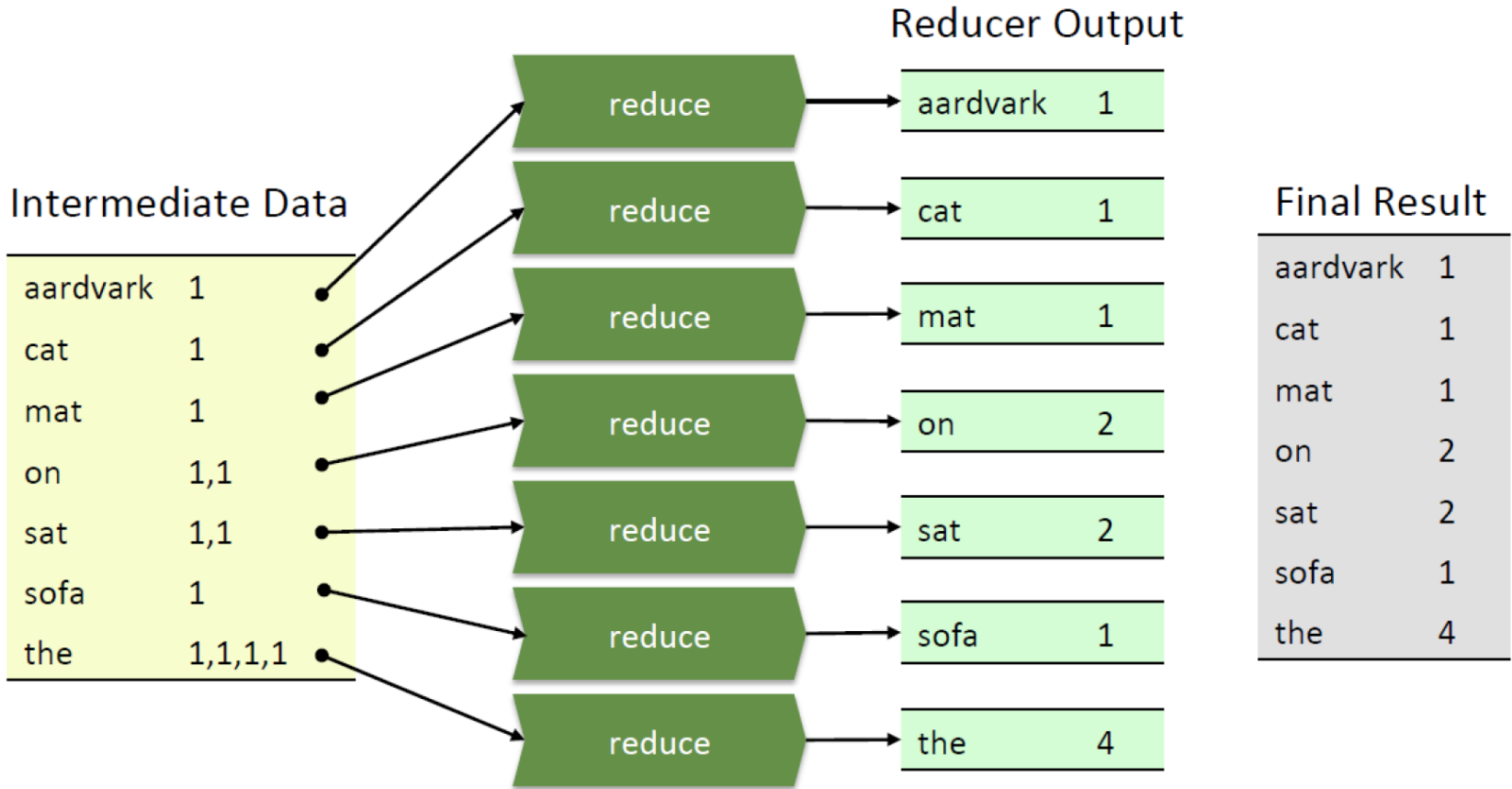
the	1
cat	1
sat	1
on	1
the	1
mat	1
the	1
aardvark	1
sat	1
on	1
the	1
sofa	1



Intermediate Data

aardvark	1
cat	1
mat	1
on	1,1
sat	1,1
sofa	1
the	1,1,1,1

Example: SumReducer



MapReduce in Scala

- Click on the green Play button and the results will be displayed on the right side of the window.

The screenshot shows an IDE window with a Scala file named Spark.sc. The code defines a SparkContext, reads a file, splits it into words, and performs a MapReduce operation. The output on the right shows the results of the reduceByKey operation, which are pairs of words and their counts.

```
1 import org.apache.log4j.{Level, Logger}
2 import org.apache.spark.SparkContext
3
4 Logger.getRootLogger.setLevel(Level.INFO)
5 val sc = new SparkContext( master = "local[*]", appName = "SparkTest")
6 val lines = sc.textFile( path = "C:\\Program Files\\Java\\jdk1.8.0_241\\COPYRIGHT.txt")
7 val words = lines.flatMap(line => line.split(' '))
8 val wordsKV RDD = words.map(x => (x,1))
9 val count = wordsKV RDD.reduceByKey((x,y) => x + y).map(x => (x._1, x._2))
10 count.foreach(println)
11
```

```
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
lines: org.apache.spark.rdd.RDD[String] = C:\Program Files\Java\jdk1.8.0_241\COPYRIGHT.txt
words: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[2] at flatMap at <console>
wordsKV RDD: org.apache.spark.rdd.RDD[(String, Int)] = MapPartitionsRDD[3] at map at <console>
count: Array[(String, Int)] = Array((and,19), (or,18), (to,15), (of,14), (any,13), (the,11), (are,10), (trademarks,9), (is,9), (,8))
```

References

- Google File System - <https://research.google/pubs/pub51/>
- MapReduce - <https://research.google/pubs/pub62/>
- Apache Hadoop Ecosystem - <https://www.cloudera.com/products/open-source/apache-hadoop.html>
- Apache Spark - <https://spark.apache.org/>
- Scala - <https://www.scala-lang.org/>

Bibliography

- Bernard Marr, “Data Strategy: How to Profit from a World of Big Data, Analytics and Artificial Intelligence”, Kogan Page, 2nd edition, 2021.
- Florin Pop, Gabriel Neagu, “Big Data Platforms and Applications: Case Studies, Methods, Techniques, and Performance Evaluation (Computer Communications and Networks)”, Springer; 1st edition, 2021.