

**Exp 11**  
**30/09/24**

## IMPLEMENTATION OF APACHE SPARK

### Aim:

To install the spark in Linux and execute the wordcount program in the spark shell.

### Procedure:

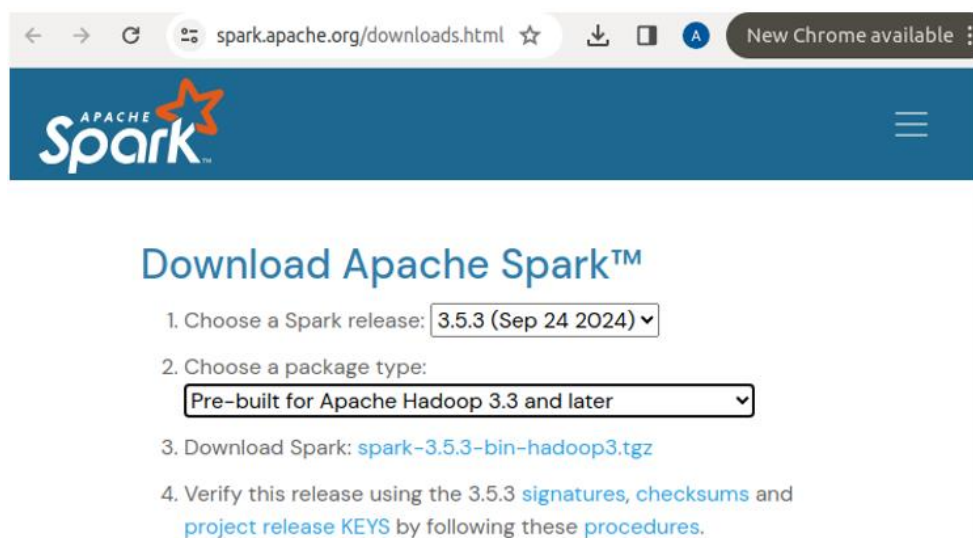
Step 1: Verify if Java is installed

```
tce@tce-VirtualBox:~$ java -version
openjdk version "11.0.19" 2023-04-18
OpenJDK Runtime Environment (build 11.0.19+7-post-Ubuntu-0ubuntu118.04.1)
OpenJDK 64-Bit Server VM (build 11.0.19+7-post-Ubuntu-0ubuntu118.04.1, mixed mode, sharing)
```

Step 2 : Verify if Spark is installed

```
tce@tce-VirtualBox:~$ sudo apt-get install scala
[sudo] password for tce:
Reading package lists... Done
Building dependency tree
Reading state information... Done
The following packages were automatically installed and are no longer required:
  gir1.2-goa-1.0 gir1.2-snapd-1
Use 'sudo apt autoremove' to remove them.
The following additional packages will be installed:
  libhawtjni-runtime-java libjansi-java libjansi-native-java libjline2-java
  scala-library scala-parser-combinators scala-xml
Suggested packages:
  scala-doc
The following NEW packages will be installed:
  libhawtjni-runtime-java libjansi-java libjansi-native-java libjline2-java
```

Step 3: Download and Install Apache Spark : <https://spark.apache.org/downloads.html>



Step 4: Check spark file using `ls` in Downloads and extract the spark file.

```
tce@tce-VirtualBox:~$ cd Downloads
tce@tce-VirtualBox:~/Downloads$ ls
'a (1).php'      'JPS Virus Maker1-20240405T180900Z-001.zip'
a.php           'JPS Virus Maker.tar.gz'
hadoop          spark-3.5.3-bin-hadoop3.tgz
Hadoop.docx     Spark_Manual.docx
'JPS Virus Maker'

tce@tce-VirtualBox:~/Downloads$ tar xvf spark-3.5.3-bin-hadoop3.tgz
spark-3.5.3-bin-hadoop3/
spark-3.5.3-bin-hadoop3/data/
spark-3.5.3-bin-hadoop3/data/graphx/
spark-3.5.3-bin-hadoop3/data/graphx/users.txt
spark-3.5.3-bin-hadoop3/data/graphx/followers.txt
spark-3.5.3-bin-hadoop3/data/mllib/
spark-3.5.3-bin-hadoop3/data/mllib/sample_linear_regression_data.txt
spark-3.5.3-bin-hadoop3/data/mllib/sample_fpgrowth.txt
spark-3.5.3-bin-hadoop3/data/mllib/sample_libsvm_data.txt
spark-3.5.3-bin-hadoop3/data/mllib/gmm_data.txt
spark-3.5.3-bin-hadoop3/python/pyspark/
tce@tce-VirtualBox:~/Downloads$ ls
'a (1).php'      'JPS Virus Maker1-20240405T180900Z-001.zip'
a.php           'JPS Virus Maker.tar.gz'
hadoop          spark-3.5.3-bin-hadoop3
Hadoop.docx     spark-3.5.3-bin-hadoop3.tgz
'JPS Virus Maker'  Spark_Manual.docx
tce@tce-VirtualBox:~/Downloads$
```

Step 5: Move the spark-3.5.3-bin-hadoop3 to spark directory.

```
tce@tce-VirtualBox:~/Downloads$ cd
tce@tce-VirtualBox:~$ sudo su
[sudo] password for tce:
root@tce-VirtualBox:/home/tce# cd /home/tce/Downloads/
root@tce-VirtualBox:/home/tce/Downloads# mv spark-3.5.3-bin-hadoop3 /usr/local/
spark
root@tce-VirtualBox:/home/tce/Downloads# ls
'a (1).php'      'JPS Virus Maker1-20240405T180900Z-001.zip'
a.php           'JPS Virus Maker.tar.gz'
hadoop          spark-3.5.3-bin-hadoop3.tgz
Hadoop.docx     Spark_Manual.docx
'JPS Virus Maker'

root@tce-VirtualBox:/home/tce/Downloads# cd /usr/local
root@tce-VirtualBox:/usr/local# ls
bin  etc  games  include  lib  man  sbin  share  spark  src

root@tce-VirtualBox:~# cd /usr/local/spark
root@tce-VirtualBox:/usr/local/spark# ls
bin  data  jars  LICENSE  NOTICE  R  RELEASE  yarn
conf  examples  kubernetes  licenses  python  README.md  sbin
root@tce-VirtualBox:/usr/local/spark#
```

STEP 6: Open the bashrc

```
tce@tce-VirtualBox:~$ nano ~/.bashrc
tce@tce-VirtualBox:~$
```



STEP 9: Open the file map it and split and save to sparkoutput.txt

```
scala> val textFile = sc.textFile("file:///home/tce/input.txt")
textFile: org.apache.spark.rdd.RDD[String] = file:///home/tce/input.txt MapPartitionsRDD[32] at textFile at <console>:23

scala> val words = textFile.flatMap(line => line.split(" "))
words: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[33] at flatMap at <console>:23

scala> val wordPairs = words.map(word => (word, 1))
wordPairs: org.apache.spark.rdd.RDD[(String, Int)] = MapPartitionsRDD[34] at map at <console>:23

scala> val wordCounts = wordPairs.reduceByKey((a, b) => a + b)
wordCounts: org.apache.spark.rdd.RDD[(String, Int)] = ShuffledRDD[35] at reduceByKey at <console>:23

scala> wordCounts.saveAsTextFile("file:///home/tce/sparkoutput")

scala> val output = sc.textFile("file:///home/tce/sparkoutput")
output: org.apache.spark.rdd.RDD[String] = file:///home/tce/sparkoutput MapPartitionsRDD[38] at textFile at <console>:23
```

STEP 10: Print the output in console by below command

```
scala> output.collect().foreach(println)
(football,1)
(plays,1)
(Rohankumar,2)
(with,1)
(Tejeshwar,2)
(Anbarasan,2)
(and,1)
```

### Result:

Thus the installation of Spark in linux and execution of sample program has been executed successfully and output has been verified.