## Spark Streaming Case Study

There are two parts this case study
- **First Part** - You have to create a Spark Application which streams data from a file on local directory on your machine and does the word count on the fly.

 The word should be done by the spark application in such a way that as soon as you drop the file in your local directory, your spark application should immediately do the word count for you.

- **Second Part** – In this part, you will have to create a Spark Application which should do the following

1. Pick up a file from the local directory and do  the word count
2. Then in the same Spark Application, write the code to put the same file on HDFS.
3. Then in same Spark Application, do the word count of the file copied on HDFS in step 2
4. Lastly, compare the word count of step 1 and 2. Both should match, other throw an error


### Below is the spark program for first case

import org.apache.spark.{SparkConf, SparkContext}

import org.apache.spark.streaming.{Seconds, StreamingContext}

 import org.apache.log4j.{Level,Logger}

object SparkFileStreamingWordCount

 {

 def main(args: Array[String]): Unit = { println("hey Spark Streaming")

 val conf = new SparkConf().setMaster("local[2]").setAppName("SparkSteamingExample")

val sc = new SparkContext(conf)

val rootLogger =Logger.getRootLogger()

rootLogger.setLevel(Level.ERROR)

 // Create Streaming context to set batch duration 5 seconds

 val ssc = new StreamingContext(sc, Seconds(5))

//Create RDD for text file streaming by

val lines = ssc.textFileStream("/home/acadgild/Desktop/Spark_Streaming")

//Split each line into words val words = lines.flatMap(_.split(" "))

//Count each word in each batch

val wordCounts = words.map(x => (x, 1)).reduceByKey(_ + _)

wordCounts.print()

//Start the computation

ssc.start()

//wait for the computation to terminate
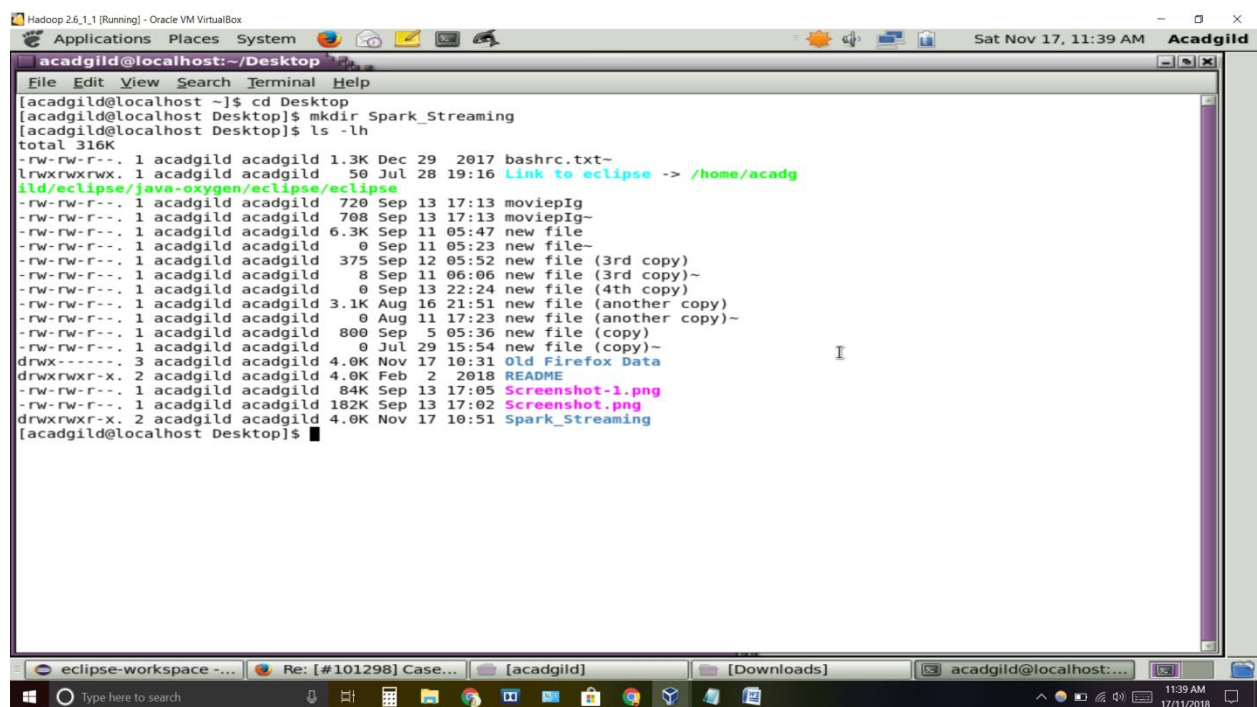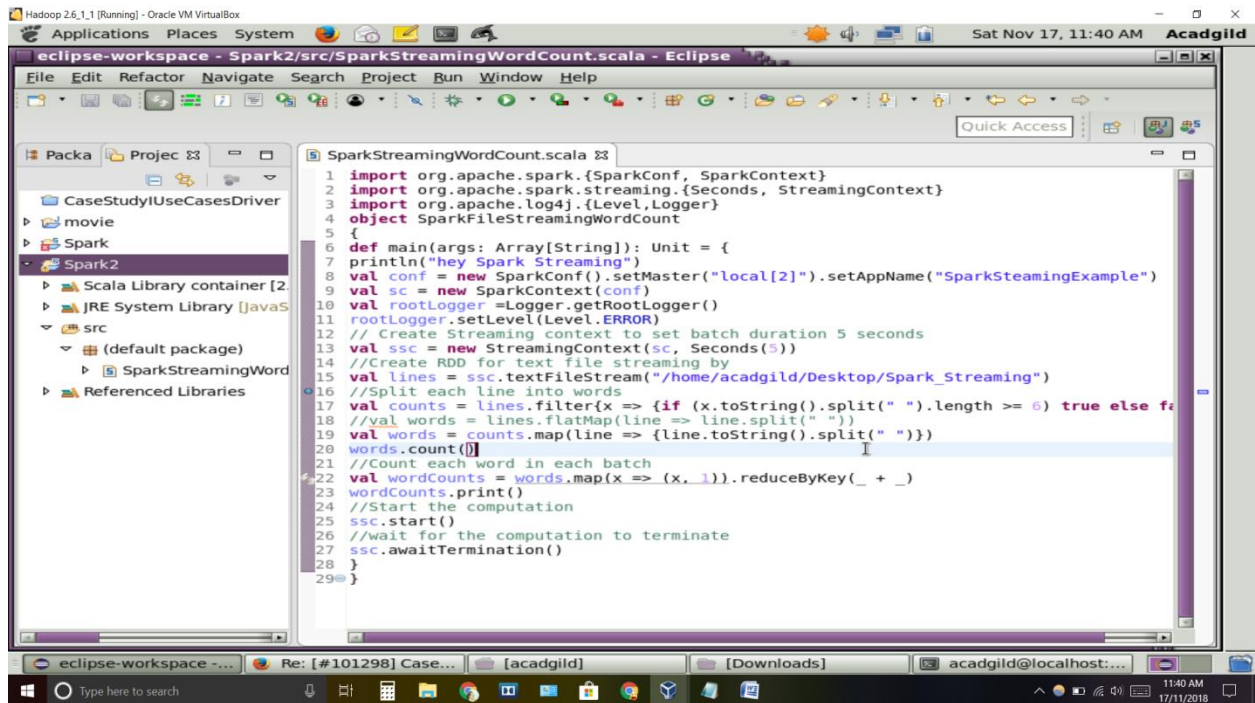
ssc.awaitTermination()

 } }

Let us create a directory Spark_Streaming in the Desktop folder of home acadgild directory by using following command:
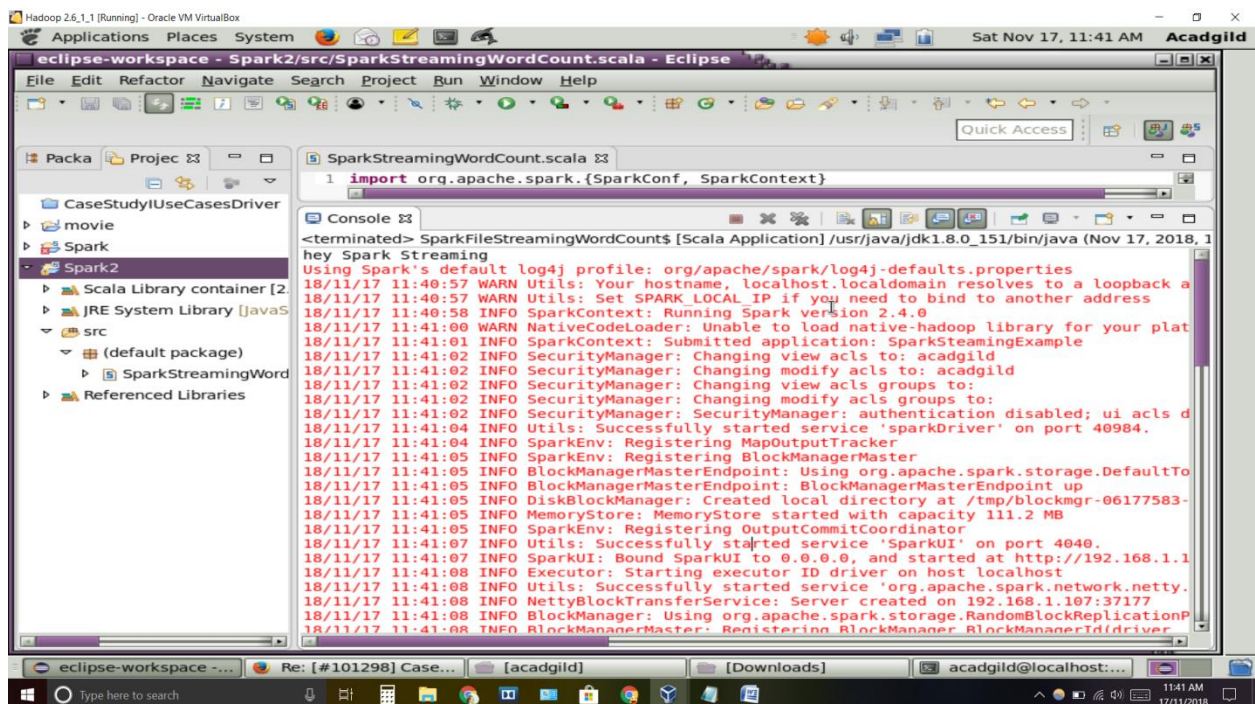
**cd Desktop**

**mkdir Spark_Streaming**

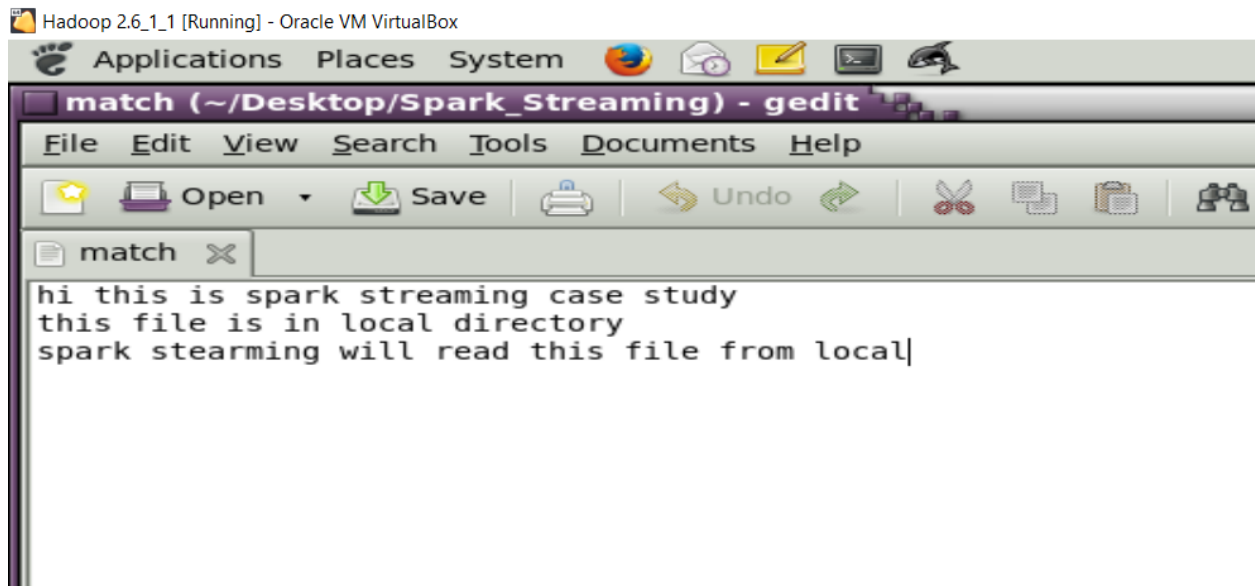 Below screenshot shows the required directory in Desktop

Run the spark application:



In below screenshot we are able to see that spark streaming is running every 5 seconds

```
Time: 1527230490000 ms
-------------------------------------------------

-------------------------------------------------
Time: 1527230495000 ms
-------------------------------------------------

-------------------------------------------------
Time: 1527230500000 ms
-------------------------------------------------

-------------------------------------------------
Time: 1527230505000 ms
-------------------------------------------------

-------------------------------------------------
Time: 1527230510000 ms
-------------------------------------------------
```

Now we put a blank text file 'match' and then putting some words in this file

Hadoop 2.6_1_1 [Running] - Oracle VM VirtualBox

Applications   Places   System

**match (~/Desktop/Spark_Streaming) - gedit**

File   Edit   View   Search   Tools   Documents   Help

Open   ▾   Save   Undo

match ✖

```
hi this is spark streaming case study
this file is in local directory
spark stearming will read this file from local
```

Output of word count in console:

```
(this,3)
(is,2)
(will,1)
(read,1)
(streaming,2)
(case,1)
(file,2)
(spark,2)
(directory,1)
(a,1)
```

Above Screenshot shows the word count of each words in text file from local file directory

## Below is the spark programming for second case

```scala
import java.io.File

import org.apache.spark.{SparkConf, SparkContext}

import scala.io.Source._

import org.apache.log4j.{Level,Logger}

object SparkHDFSWordCountComparison

{

// defining the local file directory

private var localFilePath: File = new File("/home/acadgild/Desktop/Spark_Streaming/text")

//defining the directory in hdfs path

private var dfsDirPath: String = "hdfs://localhost:8020/user"

private val NPARAMS = 2

def main(args: Array[String]): Unit = {

//parseArgs(args)

println("SparkHDFSWordCountComparison : Main Called Successfully")

println("Performing local word count")

//read the file which is present in local directory and convert into string

val fileContents = readFile(localFilePath.toString())

println("Performing local word count - File Content ->>"+fileContents)

val localWordCount = runLocalWordCount(fileContents)

println("SparkHDFSWordCountComparison : Main Called Successfully -> Local Word Count
is - >>" +localWordCount)

println("Performing local word count Completed !!")

println("Creating Spark Context")

//Create spark context

val conf = new
SparkConf().setMaster("local[2]").setAppName("SparkHDFSWordCountComparisonApp")

val sc = new SparkContext(conf     )
```

```scala
// Setting log level to [WARN] for streaming executions and to override add a custom
log4j.properties to the //classpath

val rootLogger =Logger.getRootLogger()

rootLogger.setLevel(Level.ERROR)

println("Spark Context Created")

println("Writing local file to DFS")

val dfsFilename = dfsDirPath + "/dfs_read_write_test"

 val fileRDD = sc.parallelize(fileContents) fileRDD.saveAsTextFile(dfsFilename)

println("Writing local file to DFS Completed")

println("Reading file from DFS and running Word Count")

val readFileRDD = sc.textFile(dfsFilename)

 val dfsWordCount = readFileRDD .flatMap(_.split(" ")) .flatMap(_.split("\t"))
.filter(_.nonEmpty) .map(w => (w, 1)) .countByKey() .values .sum


sc.stop()

//apply if condition to check word count result from both the directories

if (localWordCount == dfsWordCount)

{

 println(s"Success! Local Word Count ($localWordCount) " + s"

and DFS Word Count ($dfsWordCount) agree.") }

 else {

println(s"Failure! Local Word Count ($localWordCount) " + s"

and DFS Word Count ($dfsWordCount) disagree.")

} }

 /***private def parseArgs(args: Array[String]): Unit =

 { if (args.length != NPARAMS) {

 printUsage()

 System.exit(1) } }***/
```

```scala
private def printUsage(): Unit = {

val usage: String = "DFS Read-Write Test\n"

 + "\n"

+ "Usage: localFile dfsDir\n"

+ "\n" + "localFile - (string) local file to use in test\n"

 + "dfsDir - (string) DFS directory for read/write tests\n"

println(usage) }

private def readFile(filename: String): List[String] = {

 val lineIter: Iterator[String] = fromFile(filename).getLines()

val lineList: List[String] = lineIter.toList lineList }

 def runLocalWordCount(fileContents: List[String]): Int =

{ fileContents.flatMap(_.split(" "))

.flatMap(_.split("\t"))

.filter(_.nonEmpty)

.groupBy(w => w)

.mapValues(_.size)

.values

.sum }

}
```

**Screenshot 1 — eclipse-workspace - Spark2/src/SparkHDFSWord.scala - Eclipse**

```scala
1  import java.io.File
2  import org.apache.spark.{SparkConf, SparkContext}
3  import scala.io.Source._
4  import org.apache.log4j.{Level,Logger}
5
6  object SparkHDFSWordCountComparison
7  {
8    // defining the local file directory
9    private var localFilePath: File = new File("/home/acadgild/Desktop/Spark_Streaming/text")
10   //defining the directory in hdfs path
11   private var dfsDirPath: String = "hdfs://localhost:8020/user"
12   private val NPARAMS = 2
13   def main(args: Array[String]): Unit = {
14   //parseArgs(args)
15   println("SparkHDFSWordCountComparison : Main Called Successfully")
16   println("Performing local word count")
17   //read the file which is present in local directory and convert into string
18   val fileContents = readFile(localFilePath.toString())
19   println("Performing local word count - File Content ->>"+fileContents)
20   val localWordCount = runLocalWordCount(fileContents)
21   println("SparkHDFSWordCountComparison : Main Called Successfully -> Local Word Count is -"+localWordCount)
22   println("Performing local word count Completed !!")
23   println("Creating Spark Context")
24   //Create spark context
25   val conf = new SparkConf().setMaster("local[2]").setAppName("SparkHDFSWordCountComparisonApp")
26   val sc = new SparkContext(conf)
27   // Setting log level to [WARN] for streaming executions and to override add a custom log4j.properties
28   //classpath
```

Console
`<terminated> SparkFileStreamingWordCount$ [Scala Application] /usr/java/jdk1.8.0_151/bin/java (Nov 17, 2018, 11:40:47 AM)`

**Screenshot 2 — eclipse-workspace - Spark2/src/SparkHDFSWord.scala - Eclipse**

```scala
27   // Setting log level to [WARN] for streaming executions and to override add a custom log4j.properties
28   //classpath
29   val rootLogger =Logger.getRootLogger()
30   rootLogger.setLevel(Level.ERROR)
31   println("Spark Context Created")
32   println("Writing local file to DFS")
33   val dfsFilename = dfsDirPath + "/dfs_read_write_test"
34   val fileRDD = sc.parallelize(fileContents)
35   fileRDD.saveAsTextFile(dfsFilename)
36   println("Writing local file to DFS Completed")
37   println("Reading file from DFS and running Word Count")
38   val readFileRDD = sc.textFile(dfsFilename)
39   val dfsWordCount = readFileRDD.flatMap(_.split(" ")).flatMap(_.split("\t")).filter(_.nonEmpty).
40   map(w => (w, 1)).countByKey().values .sum
41   sc.stop()
42   //apply if condition to check word count result from both the directories
43   if (localWordCount == dfsWordCount)
44   {
45   println(s"Success! Local Word Count ($localWordCount) " +
46   s"and DFS Word Count ($dfsWordCount) agree.")
47   }
48   else {
49   println(s"Failure! Local Word Count ($localWordCount) " +
50   s"and DFS Word Count ($dfsWordCount) disagree.")
51   }
52   }
53   /***private def parseArgs(args: Array[String]): Unit = {
54   if (args.length != NPARAMS) {
```
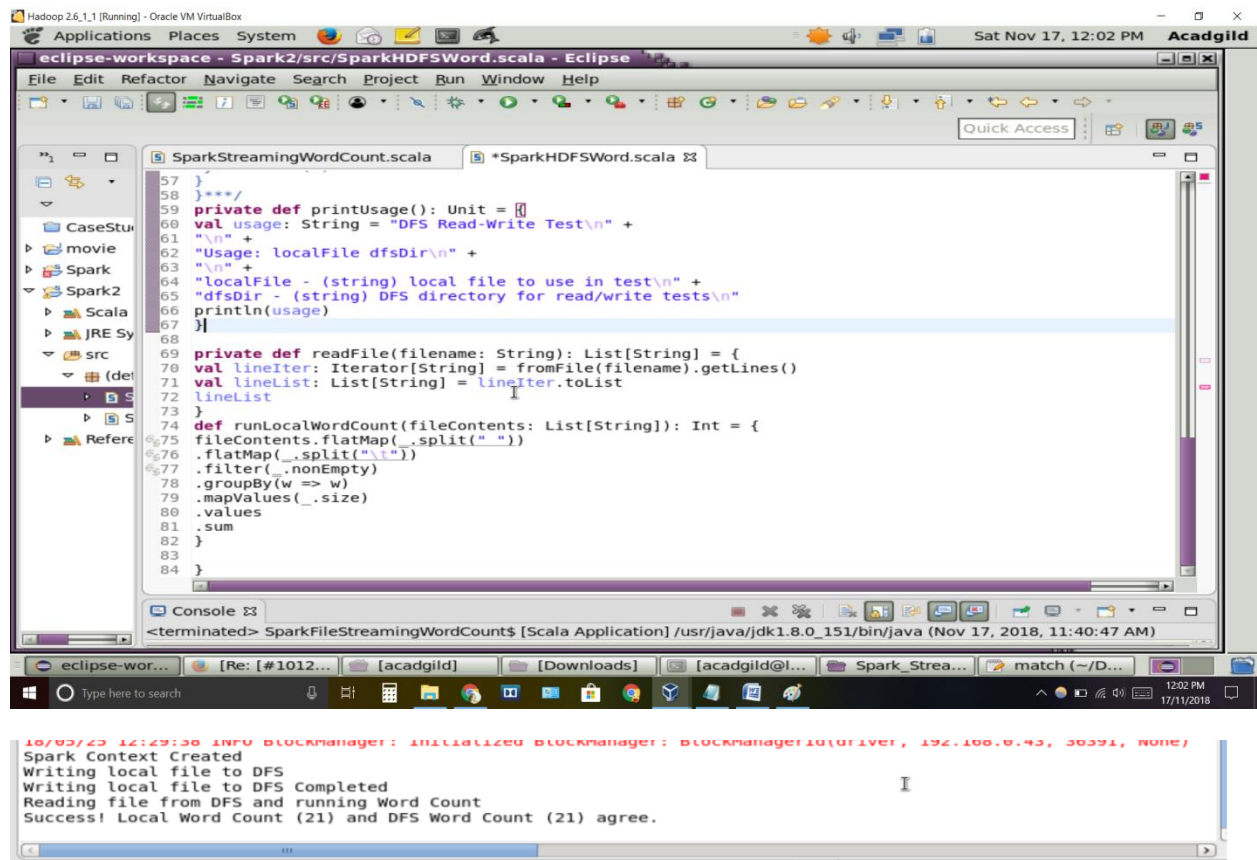
Console
`<terminated> SparkFileStreamingWordCount$ [Scala Application] /usr/java/jdk1.8.0_151/bin/java (Nov 17, 2018, 11:40:47 AM)`

Above screen shot shows that, word count of file from local file directory and HDFS directory file is same. Output file in HDFS directory shown below in screen shot.



Above screen, shot shows the content of the output file **part-00000** which same as the file present in the local directory