

There are two parts this case study

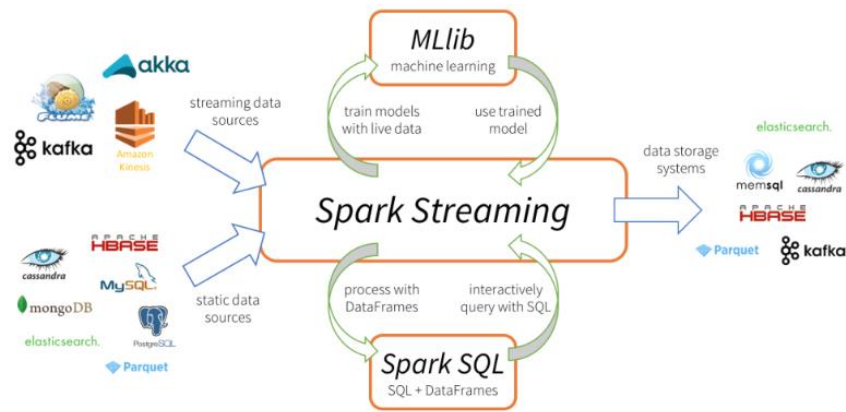
❑ **First Part** - You have to create a Spark Application which streams data from a file on local directory on your machine and does the word count on the fly. The word should be done by the spark application in such a way that as soon as you drop the file in your local directory, your spark application 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

**Solution:**

**Apache Spark Streaming is a scalable fault-tolerant streaming processing system that natively supports both batch and streaming workloads. Spark Streaming is an extension of the core Spark API that allows data engineers and data scientists to process real-time data from various sources including (but not limited to) Kafka, Flume, and Amazon Kinesis. This processed data can be pushed out to file systems, databases, and live dashboards. Its key abstraction is a Discretized Stream or, in short, a DStream, which represents a stream of data divided into small batches. DStreams are built on RDDs, Spark's core data abstraction. This allows Spark Streaming to seamlessly integrate with any other Spark components like MLlib and [Spark SQL](#).**



Below is the screen shot of the code run in eclipse for spark streaming. Separte files are attached in GitHub for the code also.

Hadoop 2.6.1\_1 [Running] - Oracle VM VirtualBox

File Machine View Input Devices Help

You have the **Auto capture keyboard** option turned on. This will cause the Virtual Machine to automatically **capture** the keyboard every time the VM window is activated and make it unavailable to other applications running on your host machine: when the keyboard is captured, all keystrokes (including system ones like **Ctrl**)

Applications Places System Tue Oct 9, 4:26 PM Acadgild

eclipse-workspace - SparkStreamingWordCount/src/pack1/SparkFileStreamingOutput.scala - Eclipse

File Edit Source Refactor Refactor Navigate Search Project Scala Run Window Help

Package Explorer

- Scala\_spark\_project
  - SparkStreamingWordCount
    - src
      - pack1
        - SparkFileStreamingOutput.scala
          - main(args: Array[String]): Unit

Referenced Libraries

```

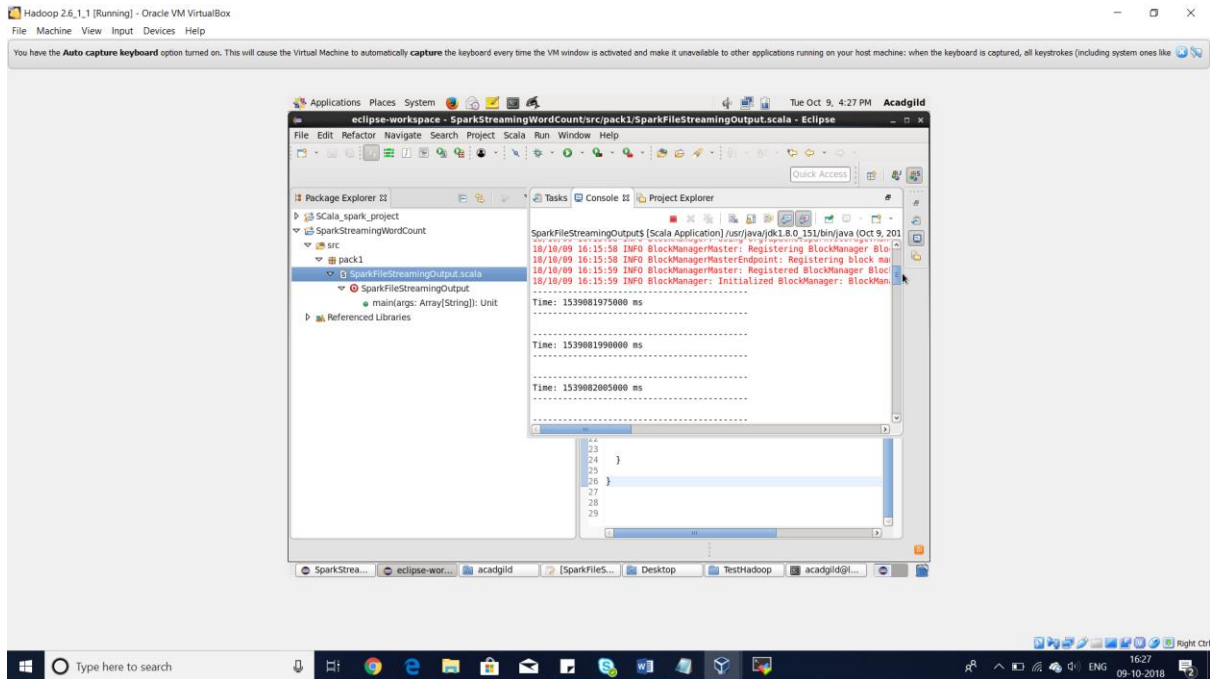
1 package pack1
2 import org.apache.spark.{SparkConf, SparkContext}
3 import org.apache.spark.streaming.{Seconds, StreamingContext}
4 import org.apache.log4j.{Level, Logger}
5
6 object SparkFileStreamingOutput {}
7
8 def main(args: Array[String]): Unit = {
9   println("hey Spark Streaming")
10
11   val conf = new SparkConf().setMaster("local[2]").setAll()
12   val sc = new SparkContext(conf)
13   val rootLogger = Logger.getLogger()
14   rootLogger.setLevel(Level.ERROR)
15   val ssc = new StreamingContext(sc, Seconds(1))
16   val lines = ssc.textFileStream(args(0))
17   val words = lines.flatMap(_.split(" "))
18   val wordCounts = words.map(x => (x, 1)).reduceByKey(_+_).print()
19   ssc.start()
20   ssc.awaitTermination()
21 }
22
23
24
25
26
27
28
29

```

SparkFileStreamingOutput.scala - SparkStreamingWordCount/src

SparkStrea... eclipse-wor... acadgild [SparkFileS... Desktop TestHadoop acadgild@l...

Type here to search 16:26 09-10-2018

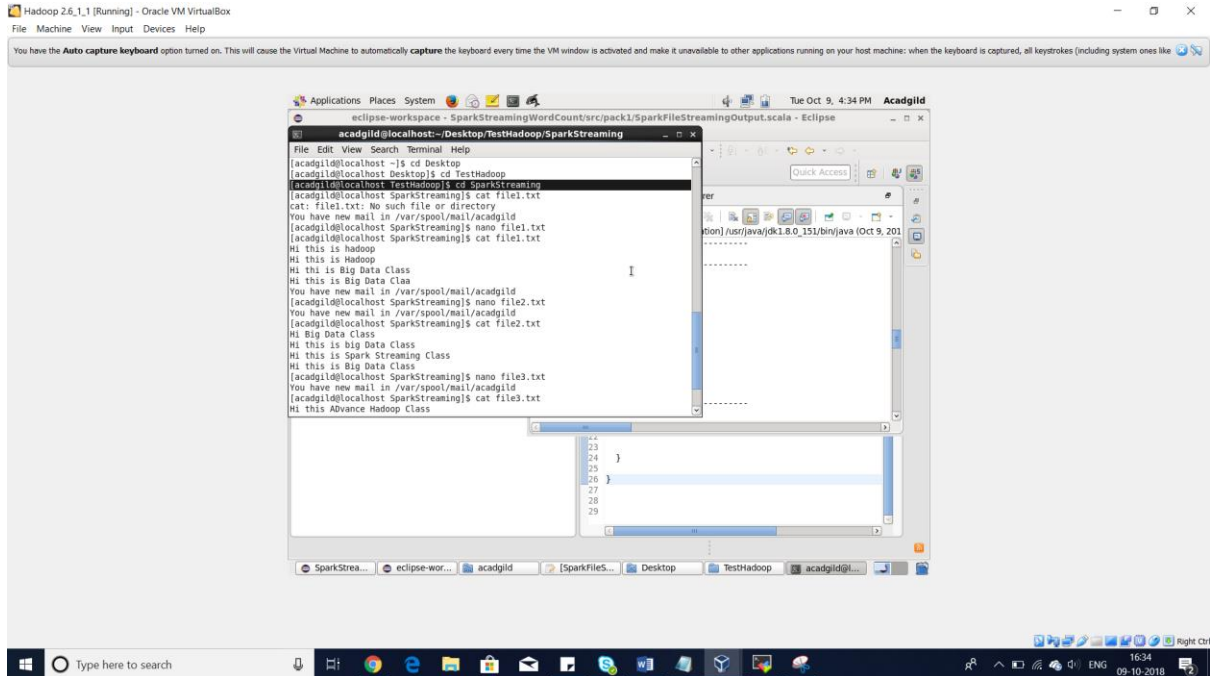


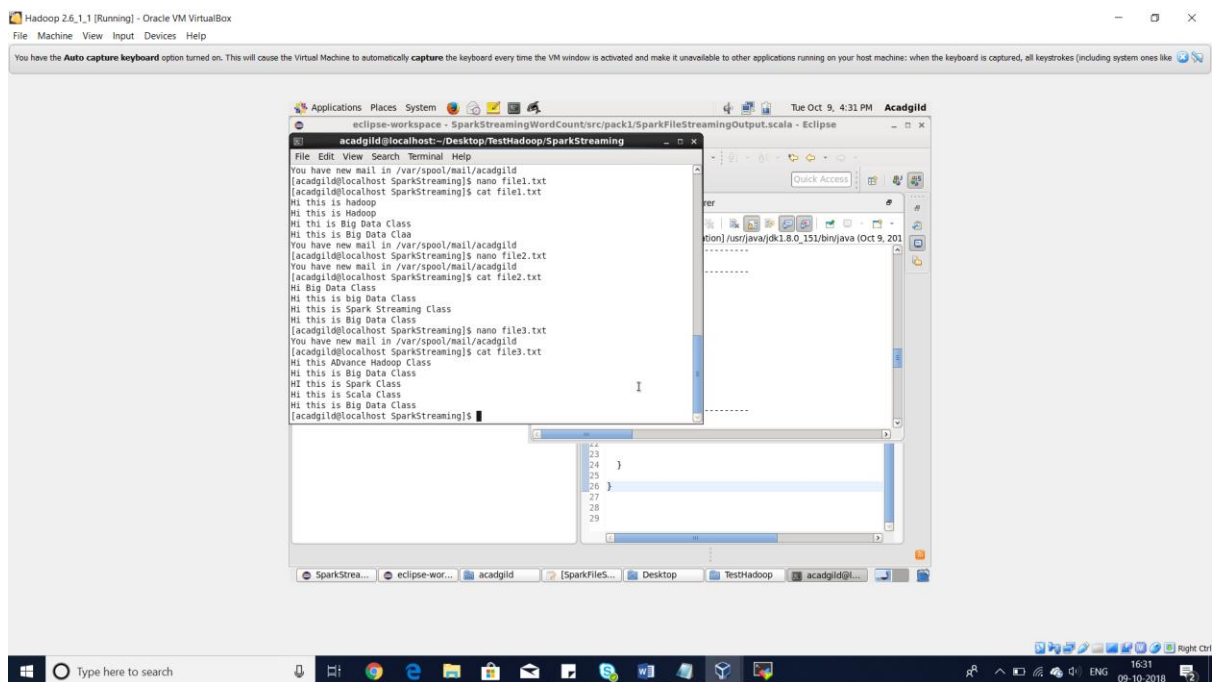
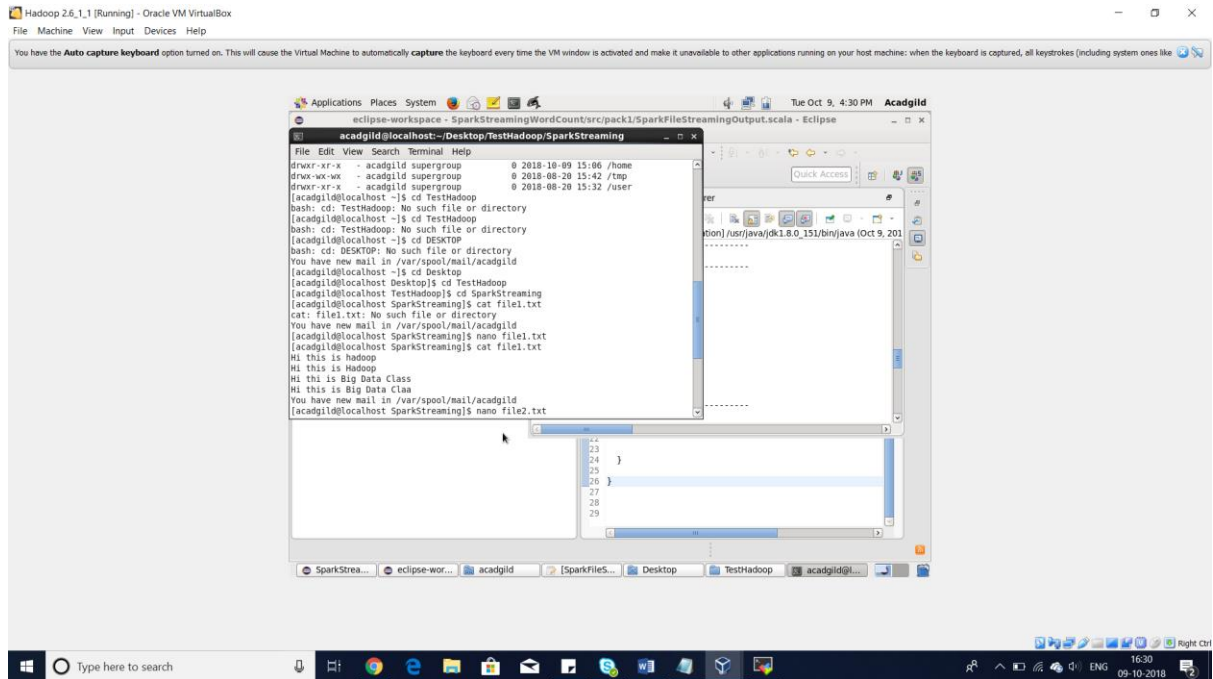
The Directory created for streaming file is: /home/acadgild/Desktop/TestHadoop/SparkStreaming

```

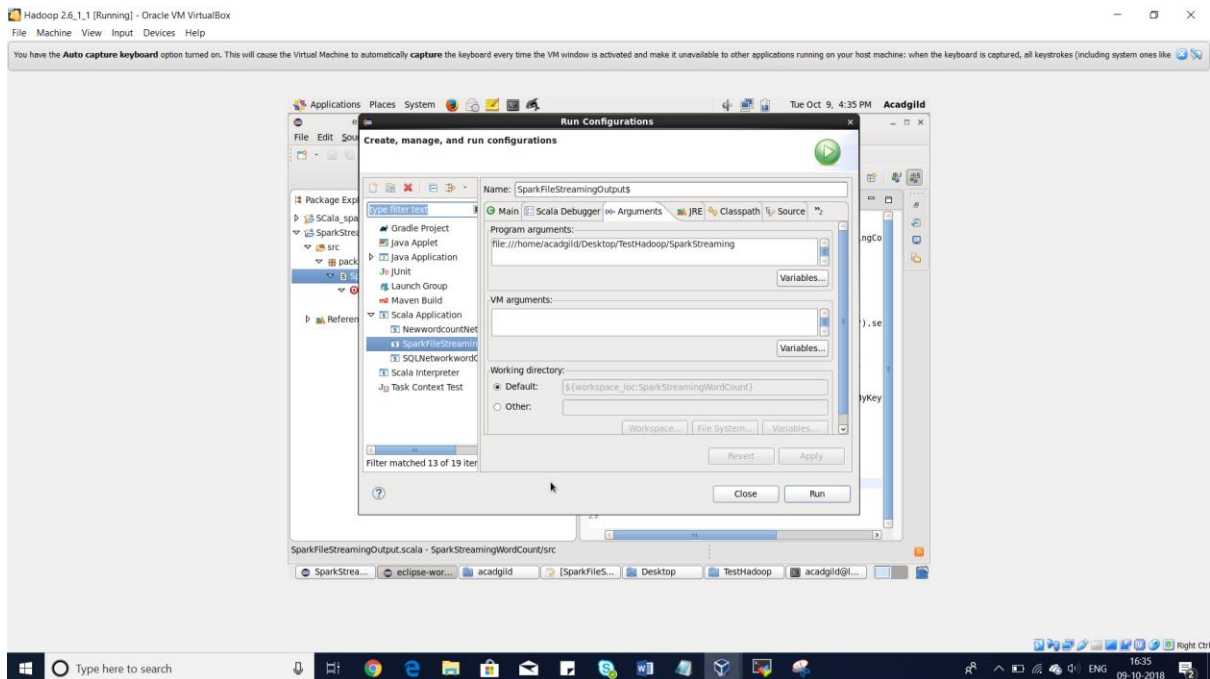
acadgild@localhost: ~/Desktop/TestHadoop/sparkStreaming
File Edit View Search Terminal Help
[acadgild@localhost sparkStreaming]$ pwd
/home/acadgild/Desktop/TestHadoop/sparkStreaming
[acadgild@localhost sparkStreaming]$ ll
total 0
[acadgild@localhost sparkStreaming]$

```

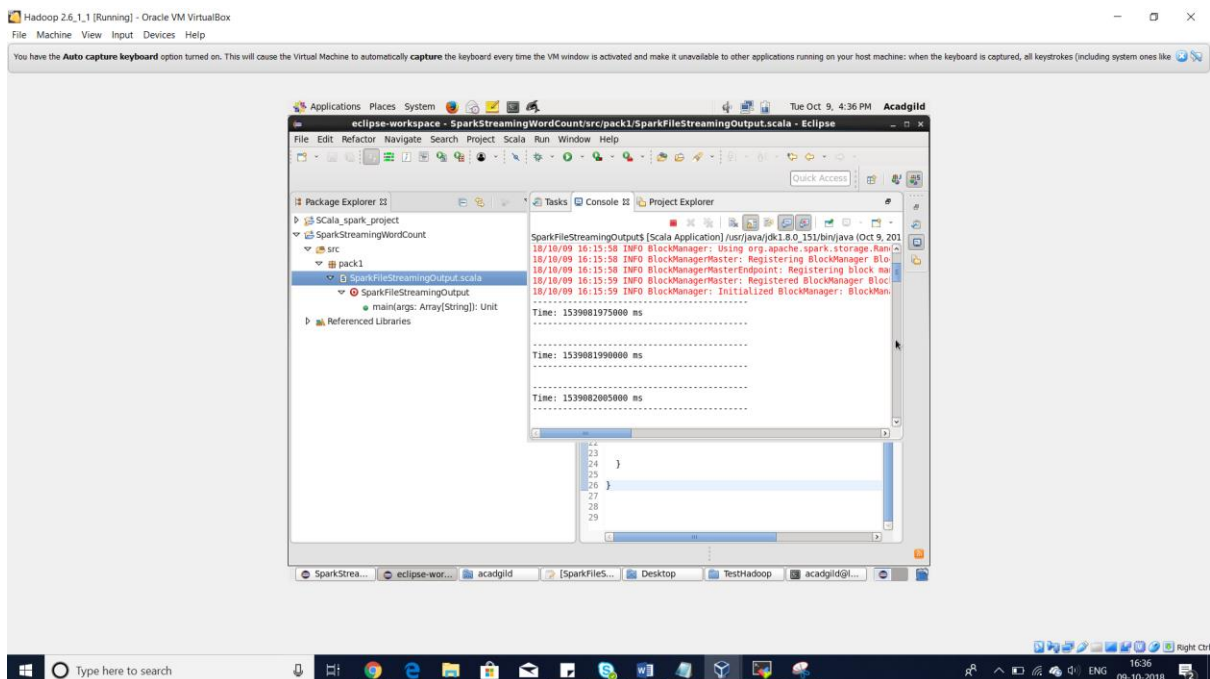




## Specifying the above created directory as input to program arguments

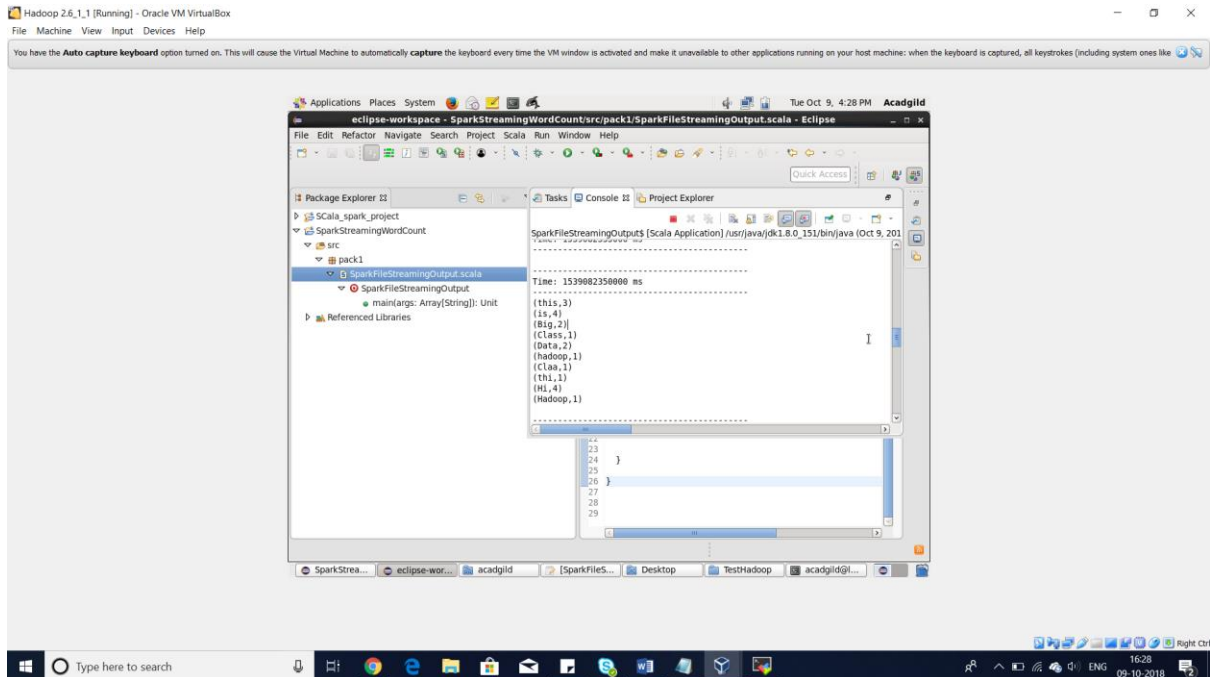


**arg(0) specified within the program reads entire files that are getting added to the above directory & performs the word count operation on the fly. Run the application as Scala Application.**

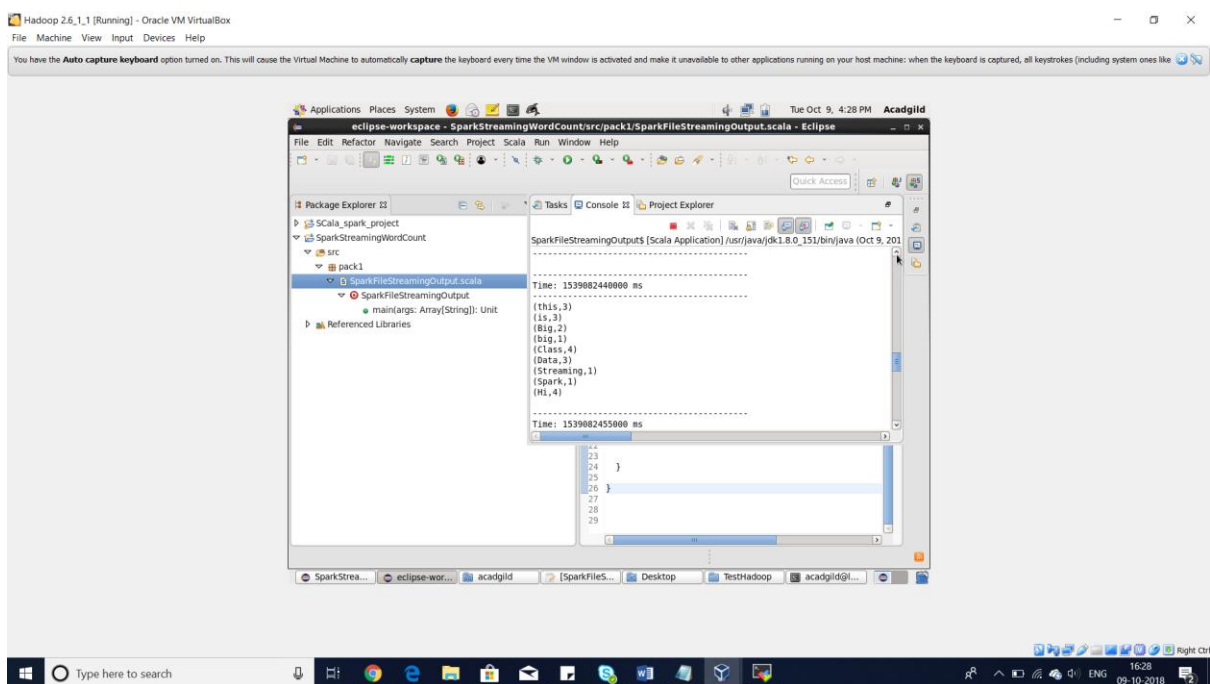


In our program we have created three text file file1, file2 and fly3 in the directory on the fly for spark streaming.

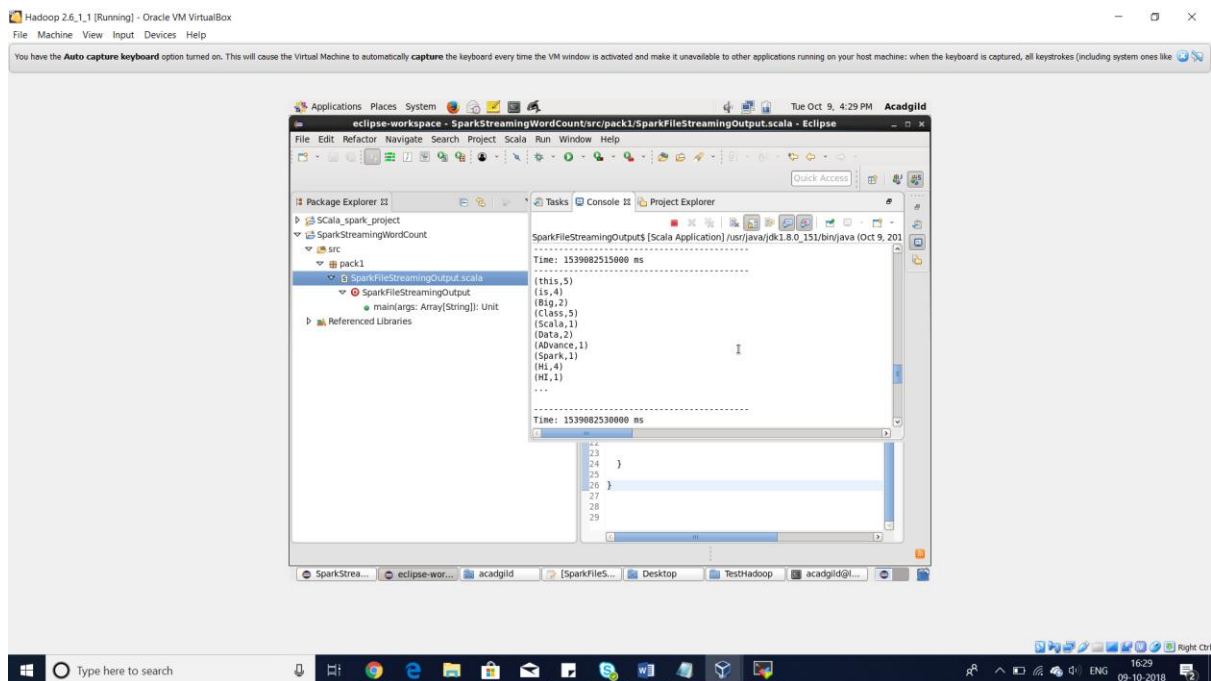
The contents of file1.txt getting read by spark streaming application and computing word count on the fly.



Similarly, creating two other files file2.txt & file3.txt whose words would be counted by spark streaming application on the fly





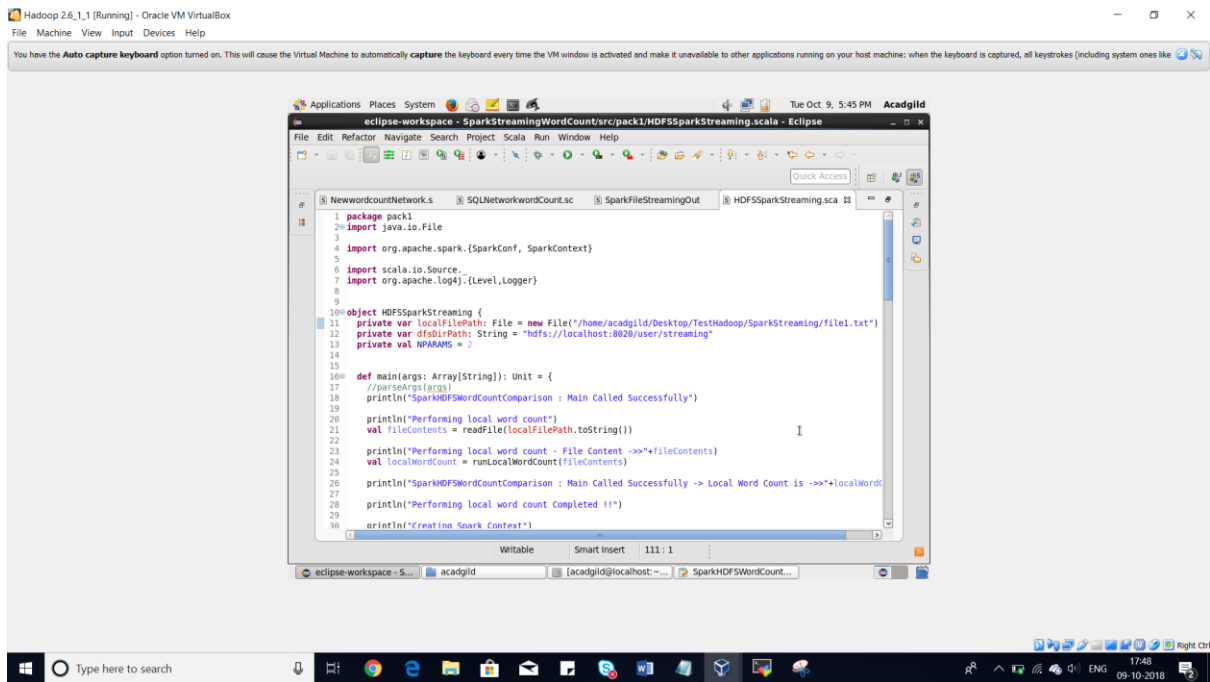
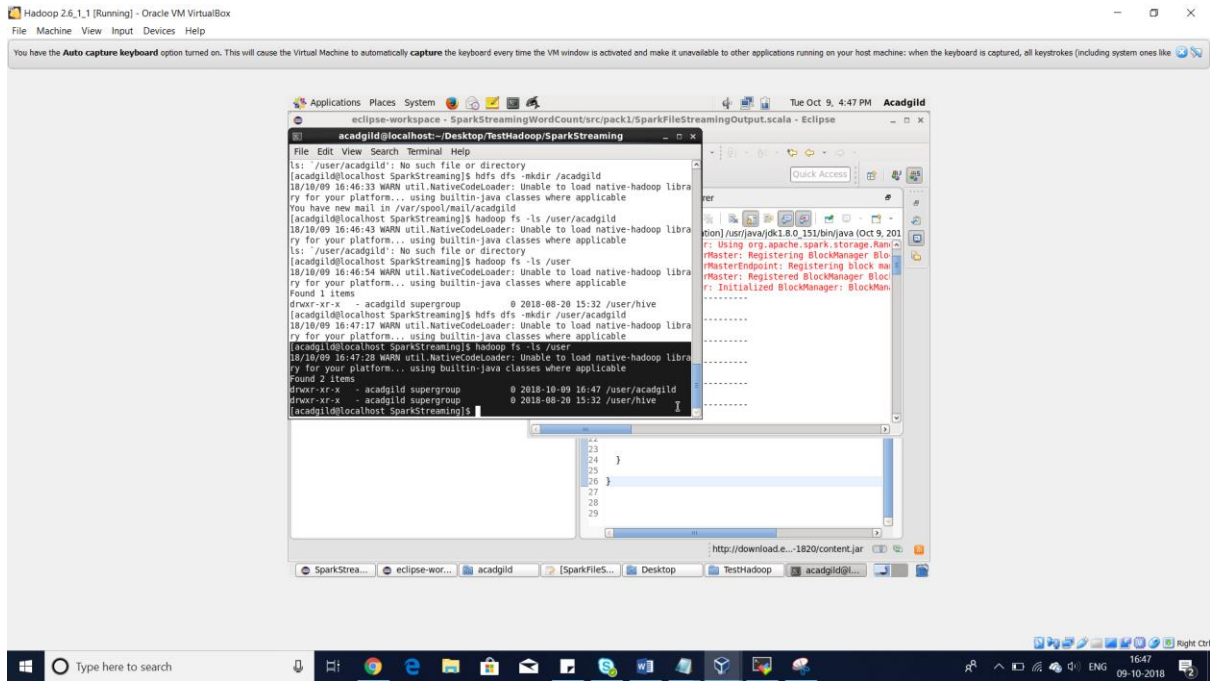


**Objective 2: Second Part** - In this part, you will have to create a Spark Application which should do the following:

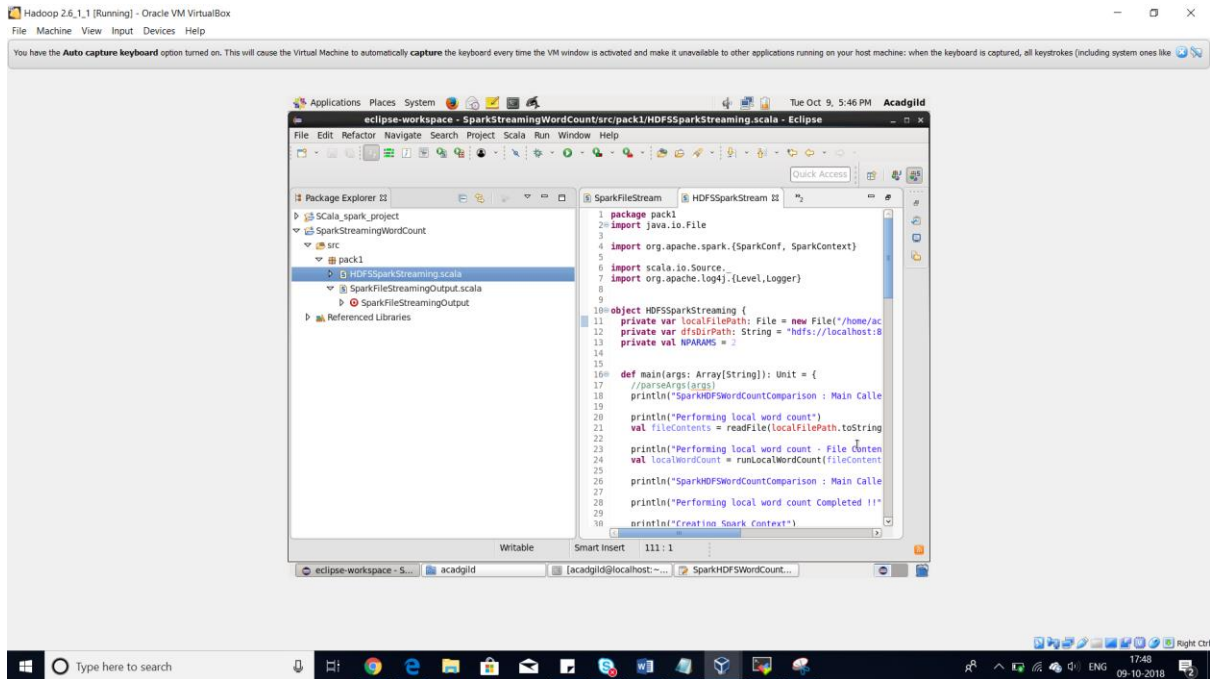
1. Pick up a file from the local directory and do the word count
2. Then in the same Spark Application, write the code to put the same file on HDFS.
3. Then in same Spark Application, do the word count of the file copied on HDFS in step 2
4. Lastly, compare the word count of step 1 and 2. Both should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

**Solution:**

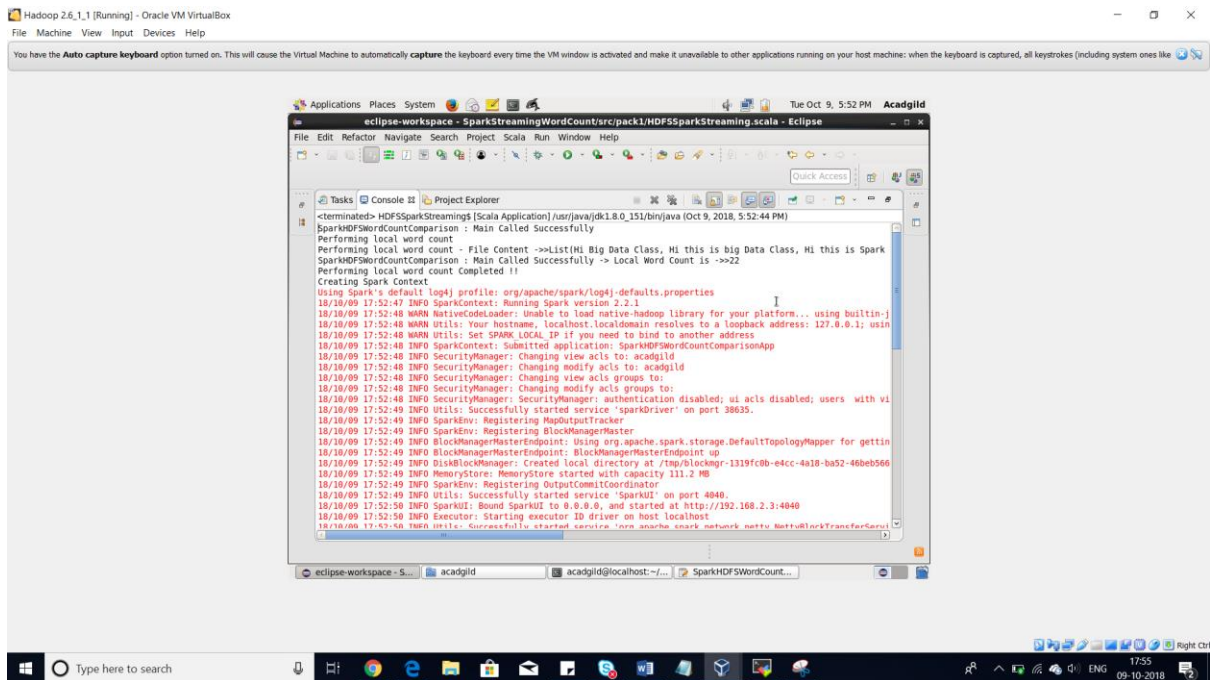
HDFS does not contain streaming directory before the application is run



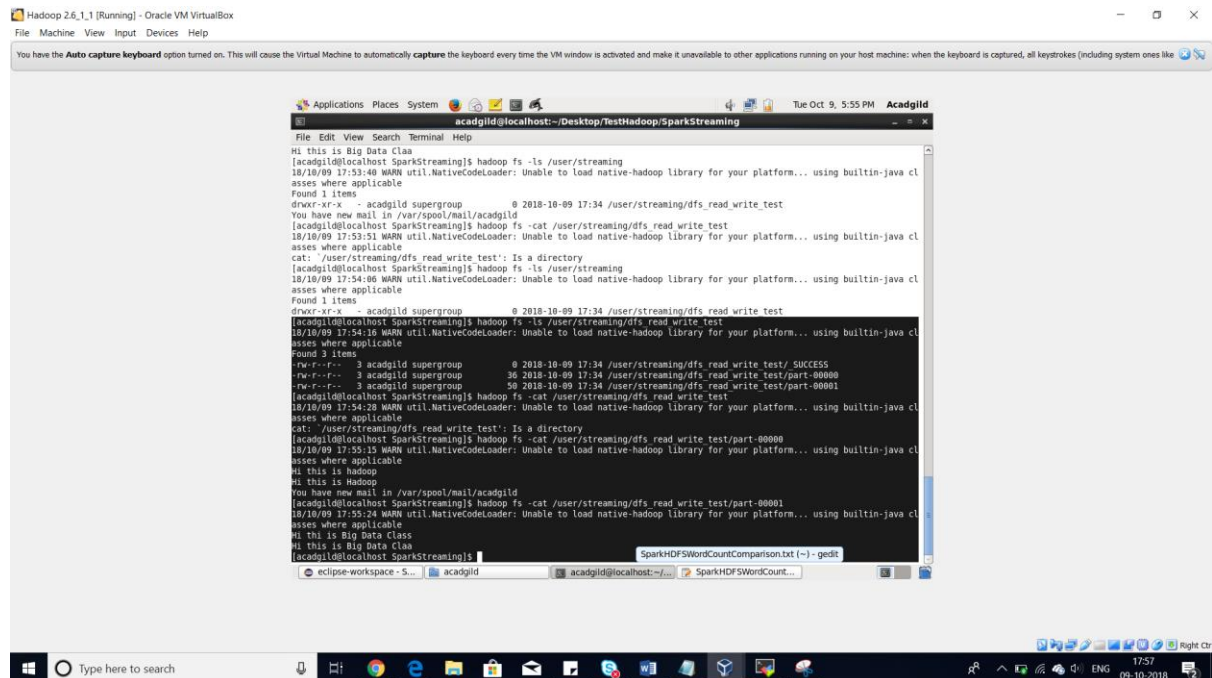




## 1. Pick up a file from the local directory and do the word count



## 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, otherwise throw an error.

Below screen shot shows the word count copied in HDFS and the word count in local directory match.

