# CSCI 5408

# DATA MANAGEMENT AND WAREHOUSING

# LAB-5: BIG DATA: HADOOP AND APACHE SPARK

**Gitlab Link**: https://git.cs.dal.ca/jspatel/csci5408_s24_b00982253_jay_patel.git

# Table of Contents

B00982253

# 1: Create Apache Spark Cluster on GCP
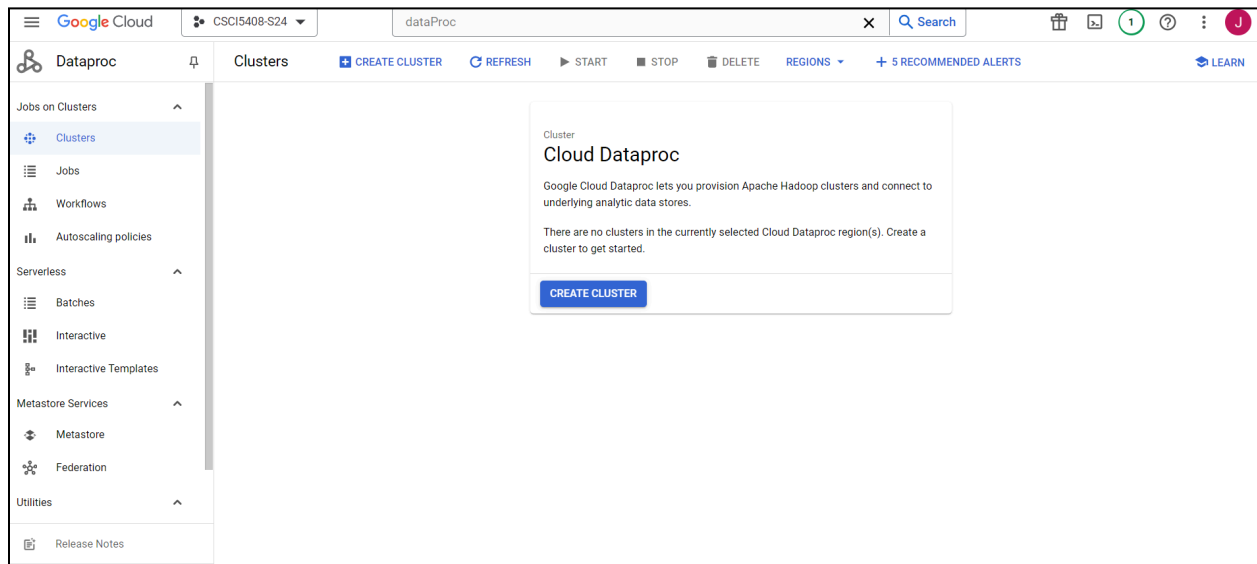


*Figure 1: Enable API for Dataproc*



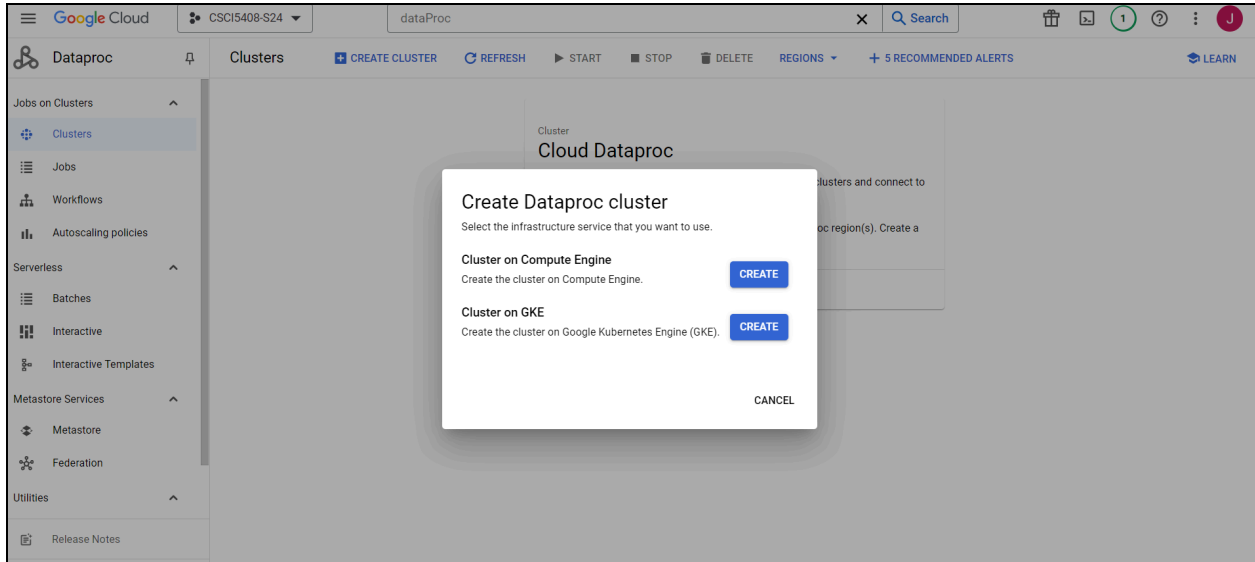*Figure 2: Create Cluster*

B00982253

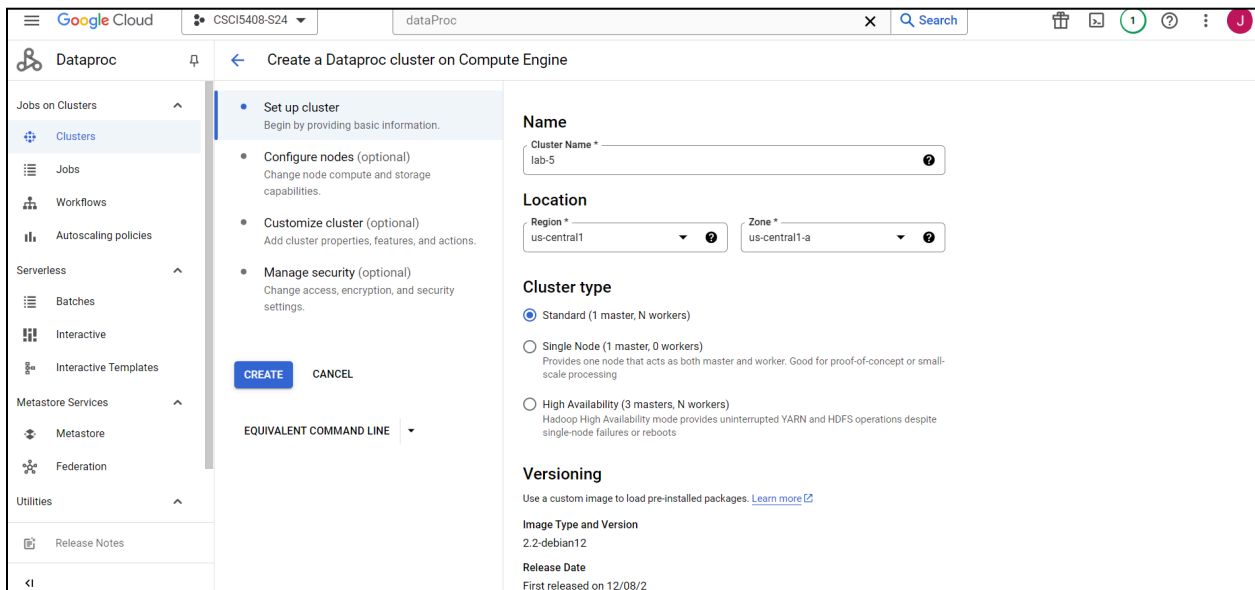*Figure 3: Select Compute Engine*



*Figure 4: Give a name, assign a location, and select cluster type.*

*Figure 5: Error due to incorrect use of machine type for worker node.*



*Figure 6: Change machine type to n2-standard-2 for worker node.*

*Figure 7: Error due to incorrect use of machine type, and primary disk size for manager node.*



*Figure 8: Change machine type to n2-standard-2, and set disk size to 50 GB for the manager node.*

*Figure 9: Error due to incorrect network configuration.*



*Figure 10: Set primary network, and sub-network to default.*

B00982253

*Figure 11: Cluster Creation in Progress.*



*Figure 12: Apache Spark cluster is successfully created and running.*

## 2: Problems faced while running the .jar file on a spark cluster



*Figure 13: Enter into the master node and check the connection using ssh.*



*Figure 14: Upload .jar and input.txt files*

**Problem 1**: Unable to find the main class - cluster was not able to find the driver or main class to run the Java program

**Solution**: Specify that main classpath with the package name using the following command.

```
jaypatel41120@lab5-m:~$ spark-submit --class Spark.DriverClass lab-5-1.0-SNAPSHOT.jar
```
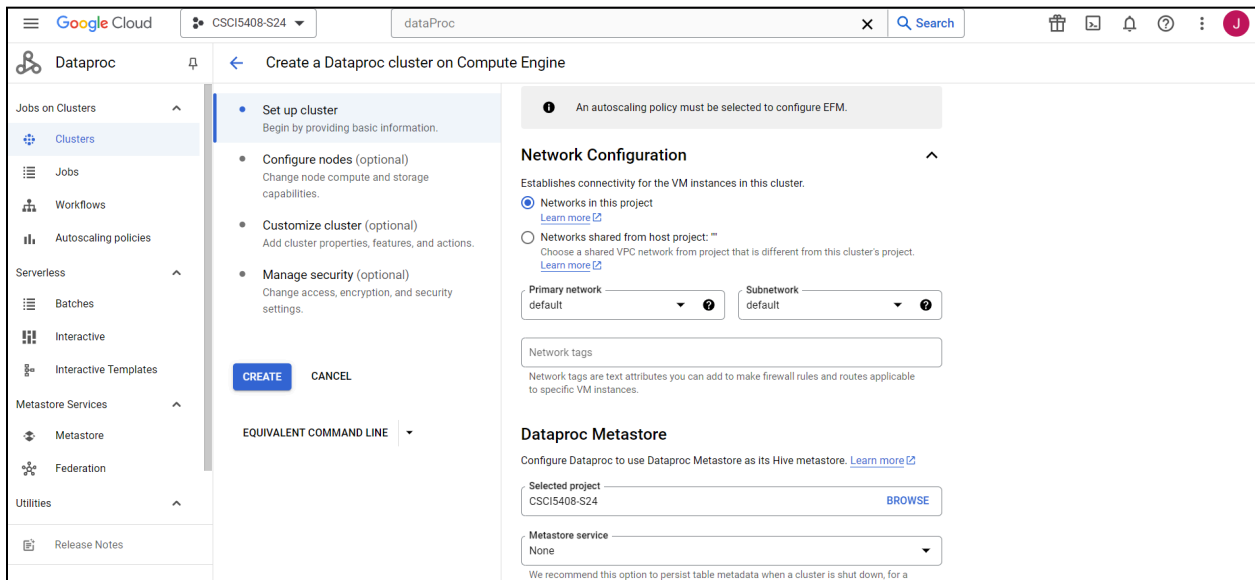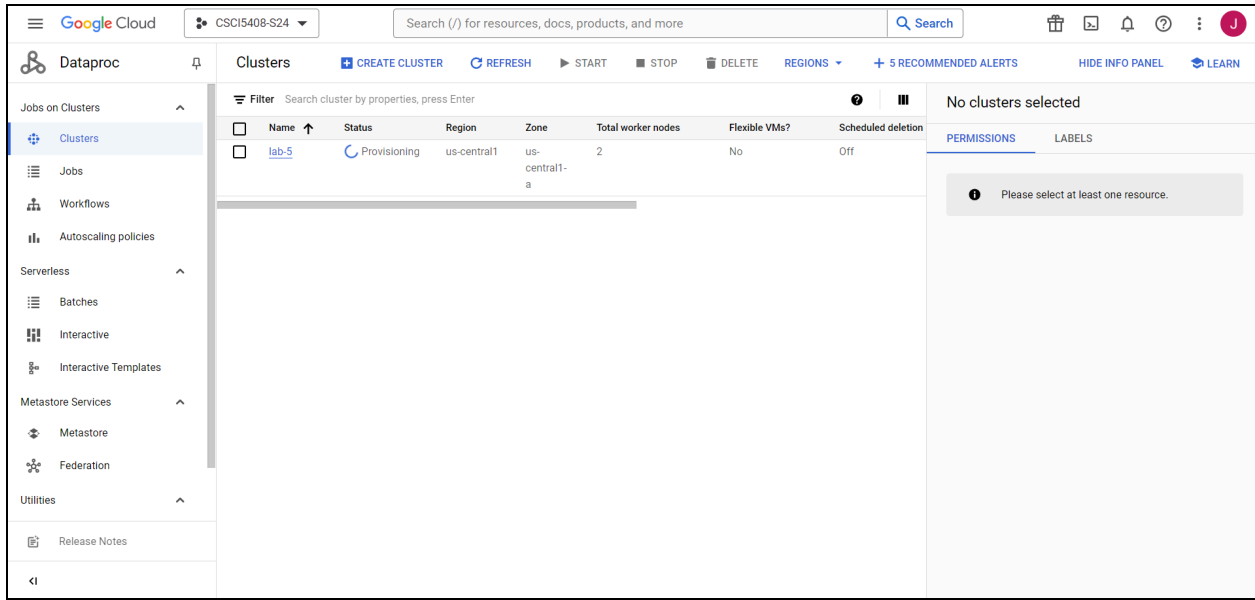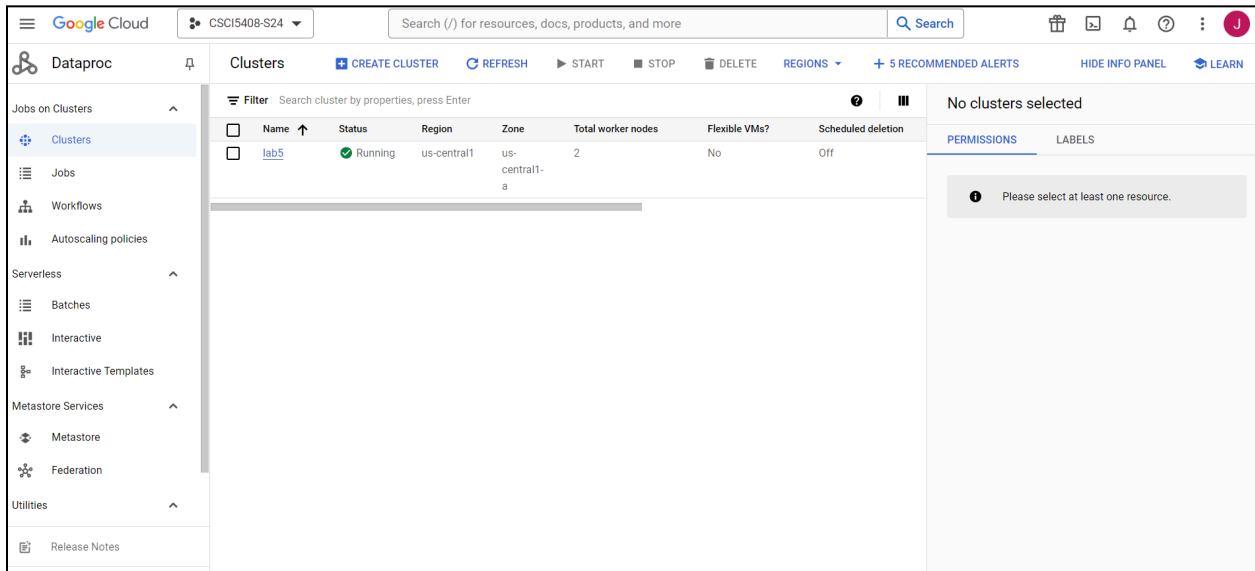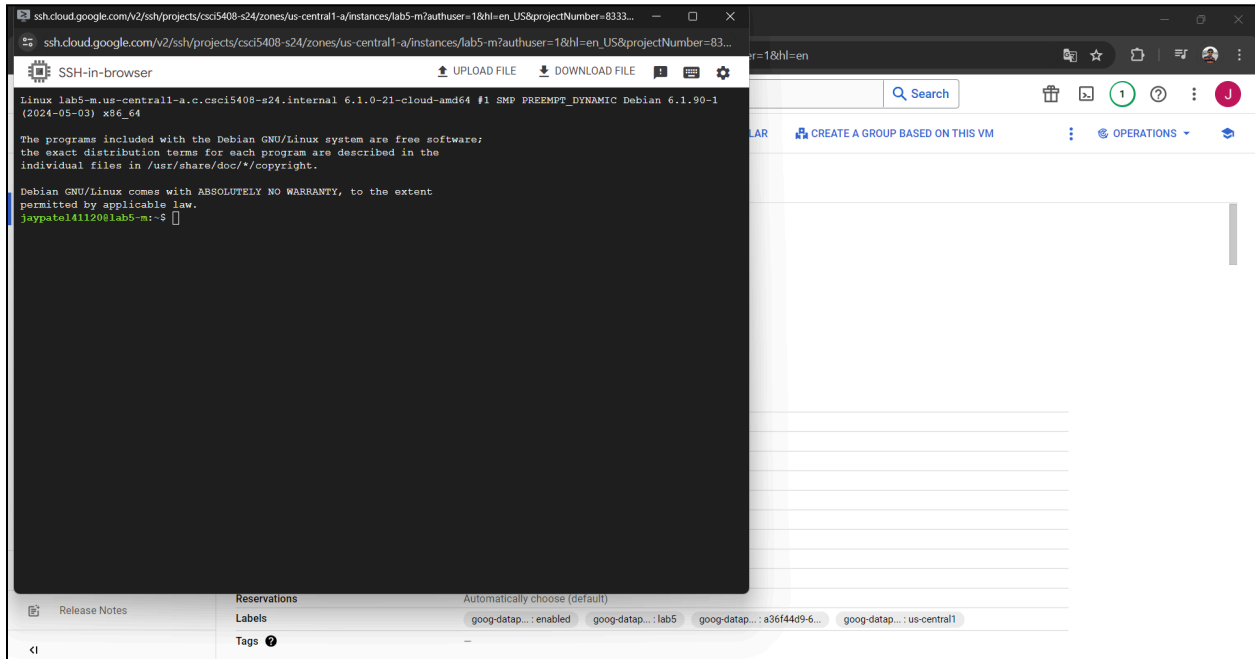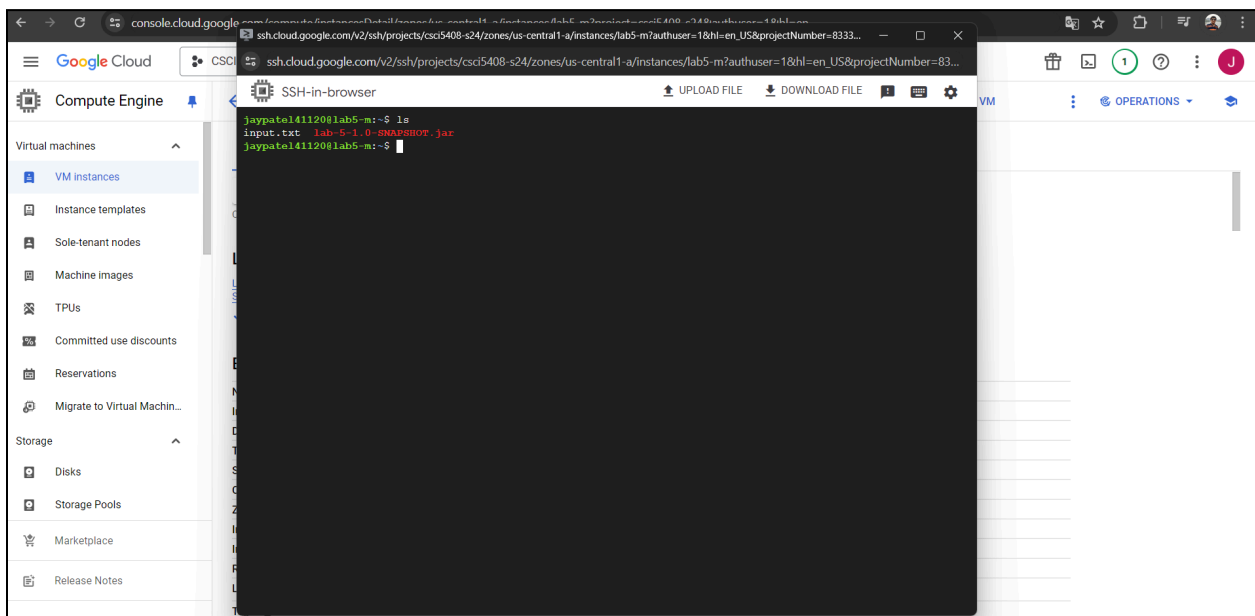
**Problem 2**: Path not found - program didn't find the input.txt file.

```
jaypatel41120@lab5-m:~$ ls
input.txt  lab-5-1.0-SNAPSHOT.jar
jaypatel41120@lab5-m:~$ spark-submit --class Spark.DriverClass lab-5-1.0-SNAPSHOT.jar
24/06/16 01:42:36 INFO SparkEnv: Registering MapOutputTracker
24/06/16 01:42:36 INFO SparkEnv: Registering BlockManagerMaster
24/06/16 01:42:36 INFO SparkEnv: Registering BlockManagerMasterHeartbeat
24/06/16 01:42:37 INFO SparkEnv: Registering OutputCommitCoordinator
24/06/16 01:42:38 INFO DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at lab5-m.us-central1-a.c.csci5408-s24.internal./10.128.0.5:8032
24/06/16 01:42:38 INFO AHSProxy: Connecting to Application History server at lab5-m.us-central1-a.c.csci5408-s24.internal./10.128.0.5:10200
24/06/16 01:42:40 INFO Configuration: resource-types.xml not found
24/06/16 01:42:40 INFO ResourceUtils: Unable to find 'resource-types.xml'.
24/06/16 01:42:41 INFO YarnClientImpl: Submitted application application_1718499166482_0003
24/06/16 01:42:42 INFO DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at lab5-m.us-central1-a.c.csci5408-s24.internal./10.128.0.5:8030
24/06/16 01:42:44 INFO MetricsConfig: Loaded properties from hadoop-metrics2.properties
24/06/16 01:42:44 INFO MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
24/06/16 01:42:44 INFO MetricsSystemImpl: google-hadoop-file-system metrics system started
24/06/16 01:42:45 INFO GoogleCloudStorageImpl: Ignoring exception of type GoogleJsonResponseException; verified object already exists with desired state.
24/06/16 01:42:46 INFO GoogleHadoopOutputStream: hflush(): No-op due to rate limit (RateLimiter[stableRate=0.2qps]): readers will *not* yet see flushed data for gs://dataproc-temp-us-central1-
833315076776-hxcmzzxk/a36f44d9-6db0-4023-9b22-3c2b4f824b73/spark-job-history/application_1718499166482_0003.inprogress [CONTEXT ratelimit_period="1 MINUTES" ]
Exception in thread "main" org.apache.spark.sql.AnalysisException: [PATH_NOT_FOUND] Path does not exist: hdfs://lab5-m/user/jaypatel41120/input.txt.
        at org.apache.spark.sql.errors.QueryCompilationErrors$.dataPathNotExistError(QueryCompilationErrors.scala:1500)
        at org.apache.spark.sql.execution.datasources.DataSource$.$anonfun$checkAndGlobPathIfNecessary$4(DataSource.scala:757)
        at org.apache.spark.sql.execution.datasources.DataSource$.$anonfun$checkAndGlobPathIfNecessary$4$adapted(DataSource.scala:754)
        at org.apache.spark.util.ThreadUtils$.$anonfun$parmap$2(ThreadUtils.scala:380)
        at scala.concurrent.Future$.$anonfun$apply$1(Future.scala:659)
        at scala.util.Success.$anonfun$map$1(Try.scala:255)
        at scala.util.Success.map(Try.scala:213)
        at scala.concurrent.Future.$anonfun$map$1(Future.scala:292)
        at scala.concurrent.impl.Promise.liftedTree1$1(Promise.scala:33)
        at scala.concurrent.impl.Promise.$anonfun$transform$1(Promise.scala:33)
        at scala.concurrent.impl.CallbackRunnable.run(Promise.scala:64)
        at java.base/java.util.concurrent.ForkJoinTask$RunnableExecuteAction.exec(ForkJoinTask.java:1426)
        at java.base/java.util.concurrent.ForkJoinTask.doExec(ForkJoinTask.java:290)
        at java.base/java.util.concurrent.ForkJoinPool$WorkQueue.topLevelExec(ForkJoinPool.java:1020)
        at java.base/java.util.concurrent.ForkJoinPool.scan(ForkJoinPool.java:1656)
        at java.base/java.util.concurrent.ForkJoinPool.runWorker(ForkJoinPool.java:1594)
        at java.base/java.util.concurrent.ForkJoinWorkerThread.run(ForkJoinWorkerThread.java:183)
```

*Figure 15: Error input.txt not found on the given path.*

**Solution**: Copy the input.txt file to the Hadoop file system using the following command.

```
jaypatel41120@lab5-m:~$ hadoop fs -copyFromLocal ./input.txt hdfs://lab5-m/user/jaypatel41120/input.txt
```
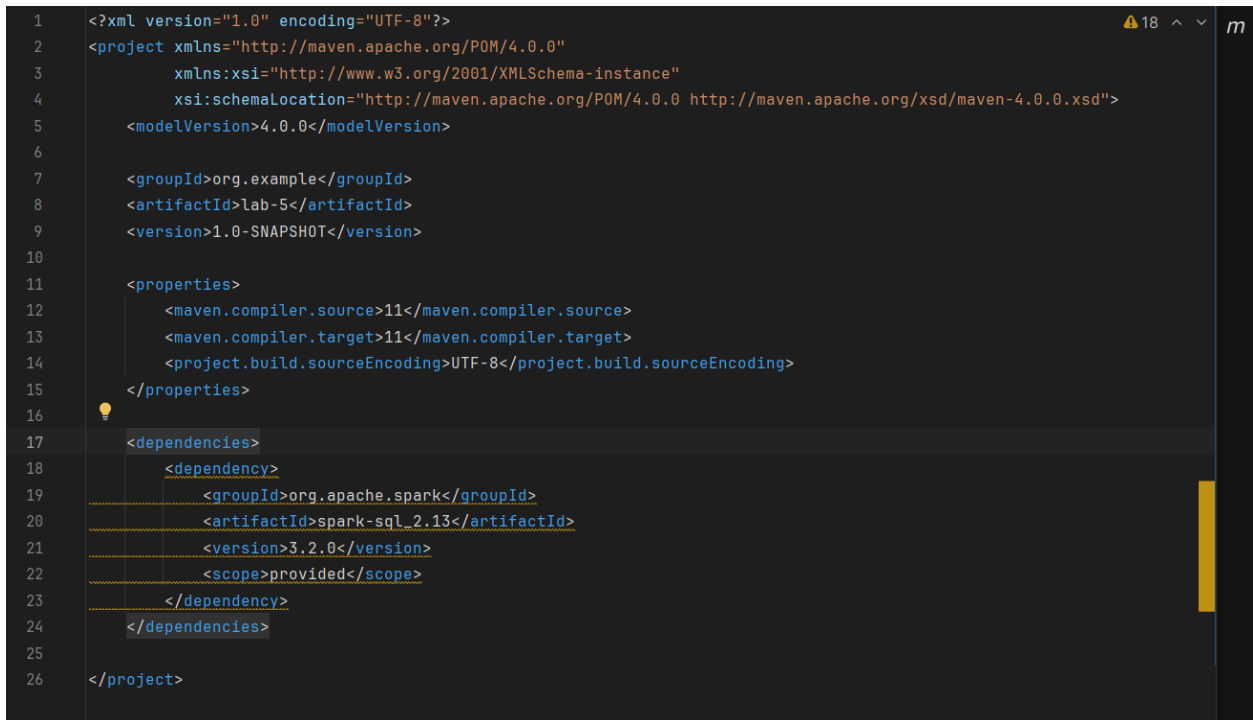
# Final Output



```
jaypatel41120@lab5-m:~$ hadoop fs -copyFromLocal ./input.txt hdfs://lab5-m/user/jaypatel41120/input.txt
jaypatel41120@lab5-m:~$ spark-submit --class Spark.DriverClass lab-5-1.0-SNAPSHOT.jar
24/06/16 01:46:03 INFO SparkEnv: Registering MapOutputTracker
24/06/16 01:46:03 INFO SparkEnv: Registering BlockManagerMaster
24/06/16 01:46:03 INFO SparkEnv: Registering BlockManagerMasterHeartbeat
24/06/16 01:46:04 INFO SparkEnv: Registering OutputCommitCoordinator
24/06/16 01:46:05 INFO DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at lab5-m.us-central1-a.c.csci5408-s24.internal./10.128.0.5:8032
24/06/16 01:46:05 INFO AHSProxy: Connecting to Application History server at lab5-m.us-central1-a.c.csci5408-s24.internal./10.128.0.5:10200
24/06/16 01:46:06 INFO Configuration: resource-types.xml not found
24/06/16 01:46:06 INFO ResourceUtils: Unable to find 'resource-types.xml'.
24/06/16 01:46:07 INFO YarnClientImpl: Submitted application application_1718499166482_0004
24/06/16 01:46:08 INFO DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at lab5-m.us-central1-a.c.csci5408-s24.internal./10.128.0.5:8030
24/06/16 01:46:10 INFO MetricsConfig: Loaded properties from hadoop-metrics2.properties
24/06/16 01:46:10 INFO MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
24/06/16 01:46:10 INFO MetricsSystemImpl: google-hadoop-file-system metrics system started
24/06/16 01:46:12 INFO GoogleCloudStorageImpl: Ignoring exception of type GoogleJsonResponseException; verified object already exists with desired state.
24/06/16 01:46:13 INFO GoogleHadoopOutputStream: hflush(): No-op due to rate limit (RateLimiter[stableRate=0.2qps]): readers will *not* yet see flushed data for gs://dataproc-temp-us-central1-
833315076776-hxcmzzxk/a36f44d9-6db0-4023-9b22-3c2b4f824b73/spark-job-history/application_1718499166482_0004.inprogress [CONTEXT ratelimit_period="1 MINUTES" ]
+----------------+
|           value|
+----------------+
|12,4,23,56,88,100|
+----------------+

Sum is : 283
jaypatel41120@lab5-m:~$
```

*Figure 16: Final Output*

# 2: Java Code Explanation

pom.xml (Insert Apache Spark dependency)

```xml
1    <?xml version="1.0" encoding="UTF-8"?>
2    <project xmlns="http://maven.apache.org/POM/4.0.0"
3             xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
4             xsi:schemaLocation="http://maven.apache.org/POM/4.0.0 http://maven.apache.org/xsd/maven-4.0.0.xsd">
5        <modelVersion>4.0.0</modelVersion>
6
7        <groupId>org.example</groupId>
8        <artifactId>lab-5</artifactId>
9        <version>1.0-SNAPSHOT</version>
10
11       <properties>
12           <maven.compiler.source>11</maven.compiler.source>
13           <maven.compiler.target>11</maven.compiler.target>
14           <project.build.sourceEncoding>UTF-8</project.build.sourceEncoding>
15       </properties>
16
17       <dependencies>
18           <dependency>
19               <groupId>org.apache.spark</groupId>
20               <artifactId>spark-sql_2.13</artifactId>
21               <version>3.2.0</version>
22               <scope>provided</scope>
23           </dependency>
24       </dependencies>
25
26   </project>
```

*Figure 17: pom.xml file of Java program*

- Here I use Apache Spark version 3.2.0 and I also have to change the compiler version from 21 to 11 because when I initially run using the 21 version it gives me an error on the GCP cluster, so after that, I just change to 11 and it works for me.

```java
8  ▷  public class DriverClass {
9  ▷      public static void main(String[] args) {
10
11          String appName = "lab-5";
12
13          SparkSession sparkSession = SparkSession.builder().appName(appName).getOrCreate();
14
15          String filePath = "input.txt";
16
17          Dataset<String> fileData = sparkSession.read().option("multiline", false).textFile(filePath);
18
19          fileData.show();
20
21          List<String> lines = fileData.collectAsList();
22
23          int sum = 0;
24          for (String line : lines){
25              String[] numbers= line.split( regex: ",");
26              for (String currNum : numbers){
27                  sum += Integer.parseInt(currNum);
28              }
29          }
30
31          System.out.println("Sum is : " + sum);
32
33          sparkSession.stop();
34      }
35  }
36
```

*Figure 18: Java Main Class*

- Here I first start the spark session by providing any random app name, and after that set the file path to read the input.txt file.
- Show the contents of the file to the console.
- Convert the file data into a list of strings.
- Now here we have a list so we have to iterate through each object of that list and here we have only one object as we have only one line in our input.txt file.
- Now in the loop first we have to split the line using the "," operator and we get the list in which all the numbers are in string format.
- Now just change numbers from string to integer and add them into the global variable named sum.
- Show the sum of all the numbers to the console.
- Stop the spark session.