# CSCI-5408

# DATA MANAGEMENT, WAREHOUSING, & ANALYTICS

# LAB ASSIGNMENT - 5

Banner ID: B00977669

GitLab Assignment Link:

https://git.cs.dal.ca/saji/csci5408_w24_b00977669_christin_saji

# Table of Contents

# Task 1: Screenshots of the step-by-step process followed to create the Apache Spark (GCP Dataproc) cluster and execute the job (WordCounter.jar) file on it.

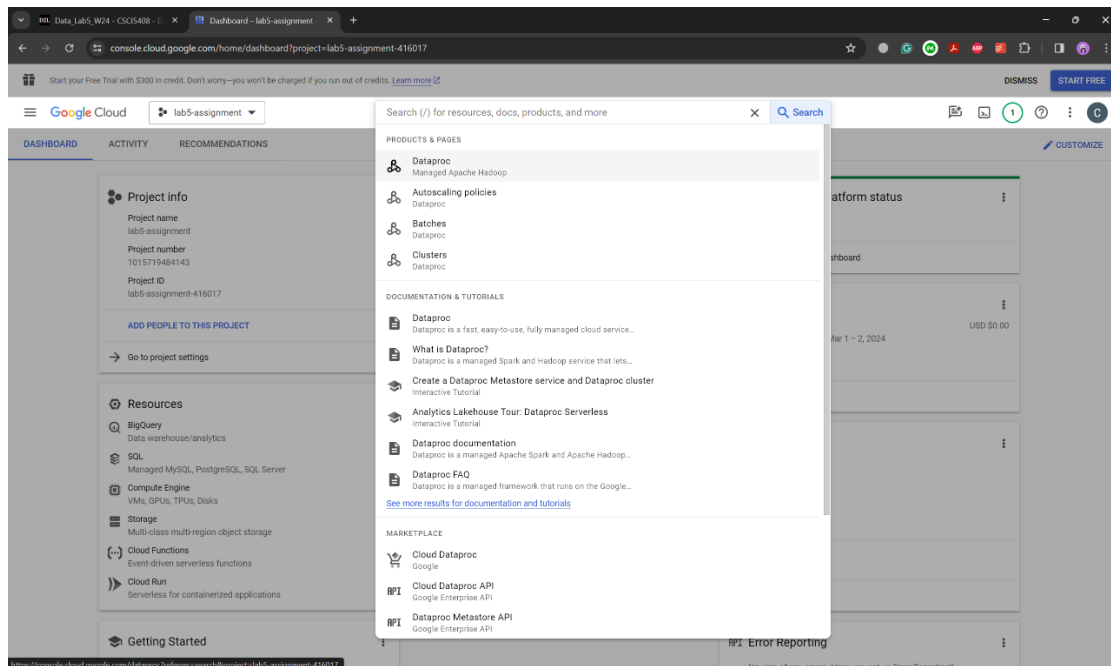First, I created a project named lab5-assignment then selected Dataproc.
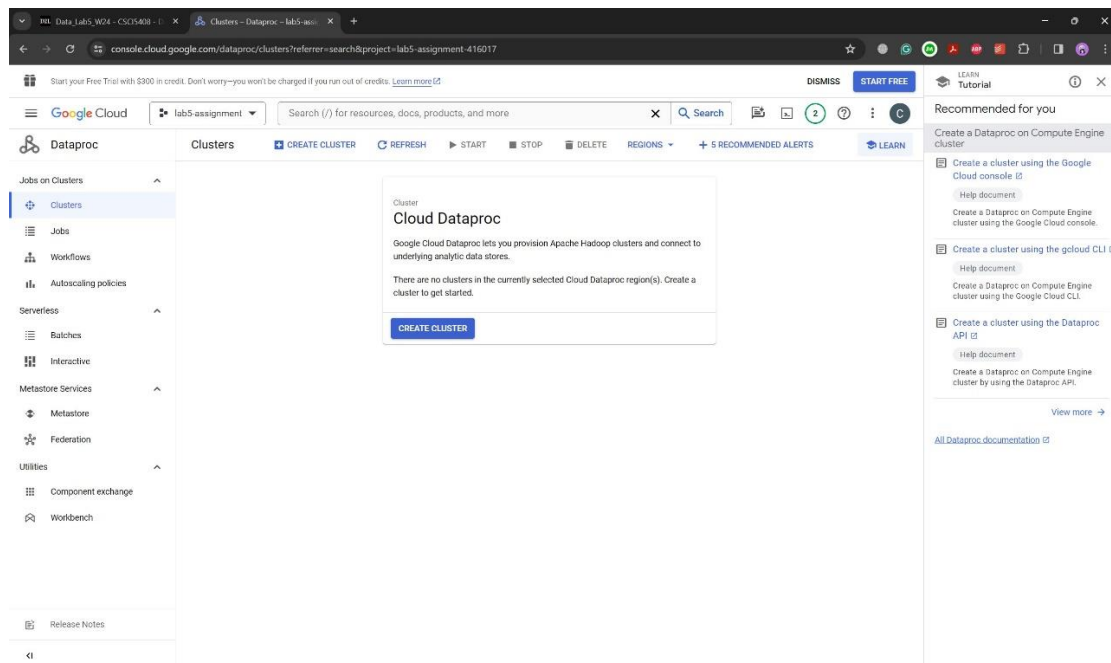


*Figure 1 Selected Dataproc*



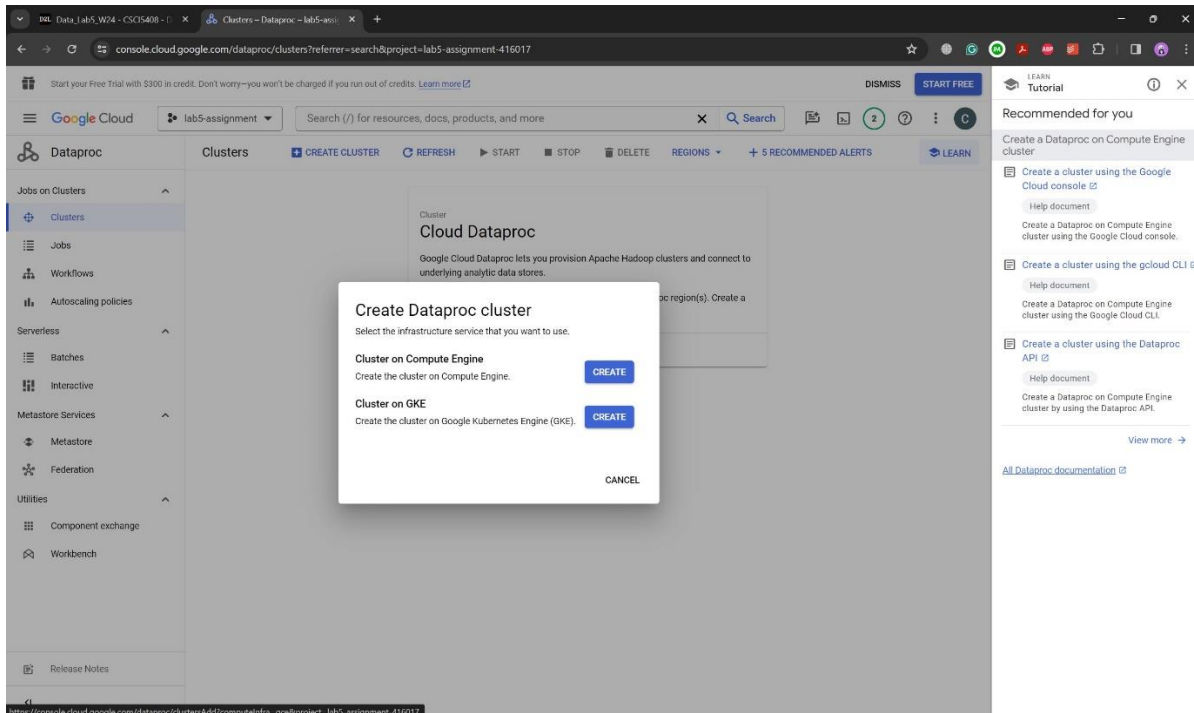*Figure 2 Create Cluster*

*Figure 3 Cluster on Compute Engine*

Renamed the cluster name to wordcounter-cluster and selected single node option for cluster type.
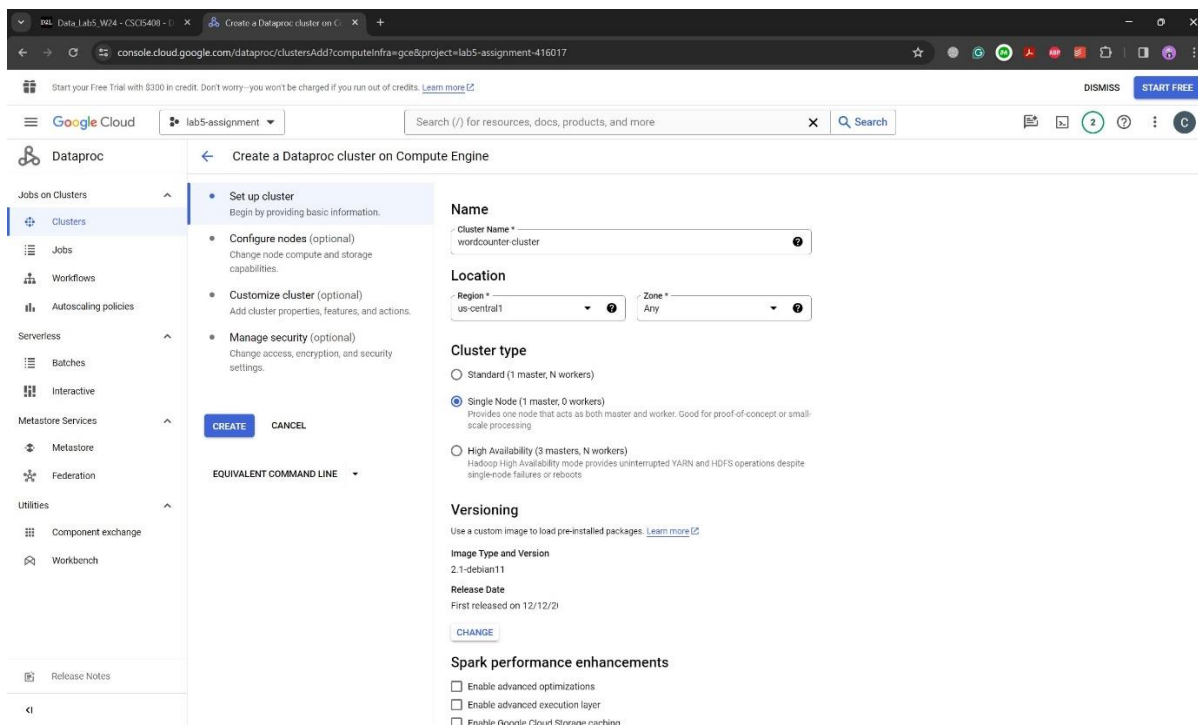


*Figure 4 Settings for set up cluster*

4

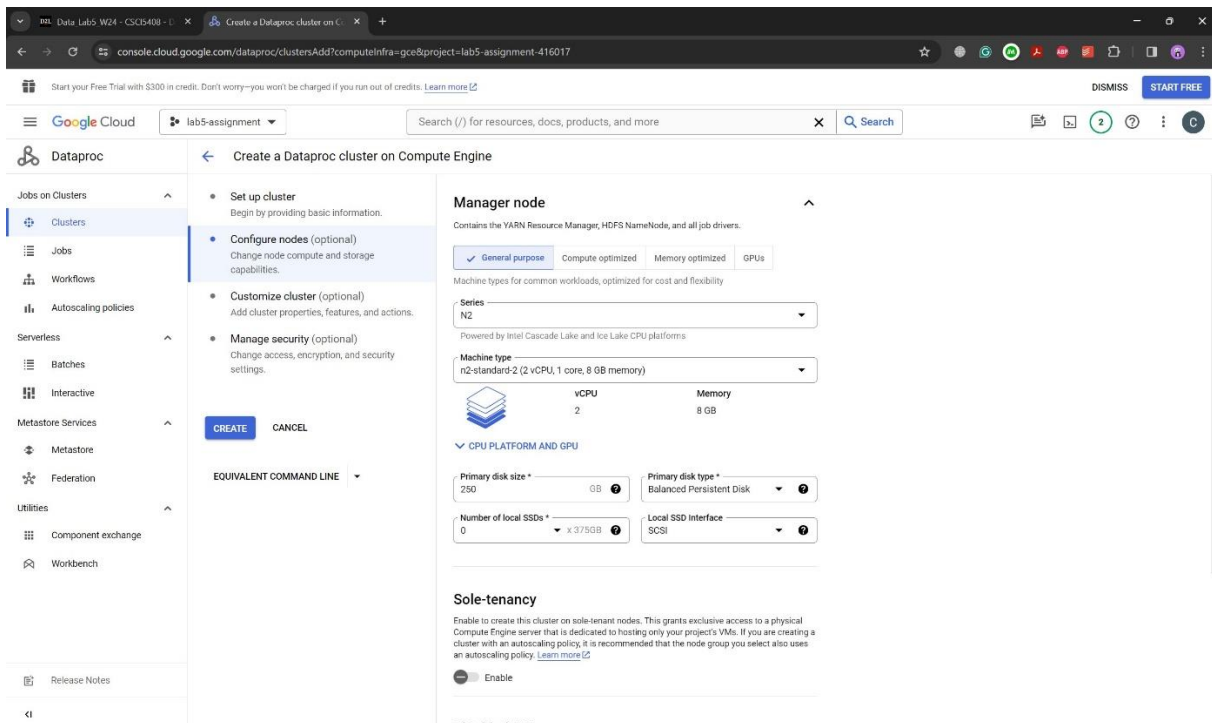In configure nodes, I changed machine type to 2 vCpu with 8 GB memory and primary disk size to 250 GB.



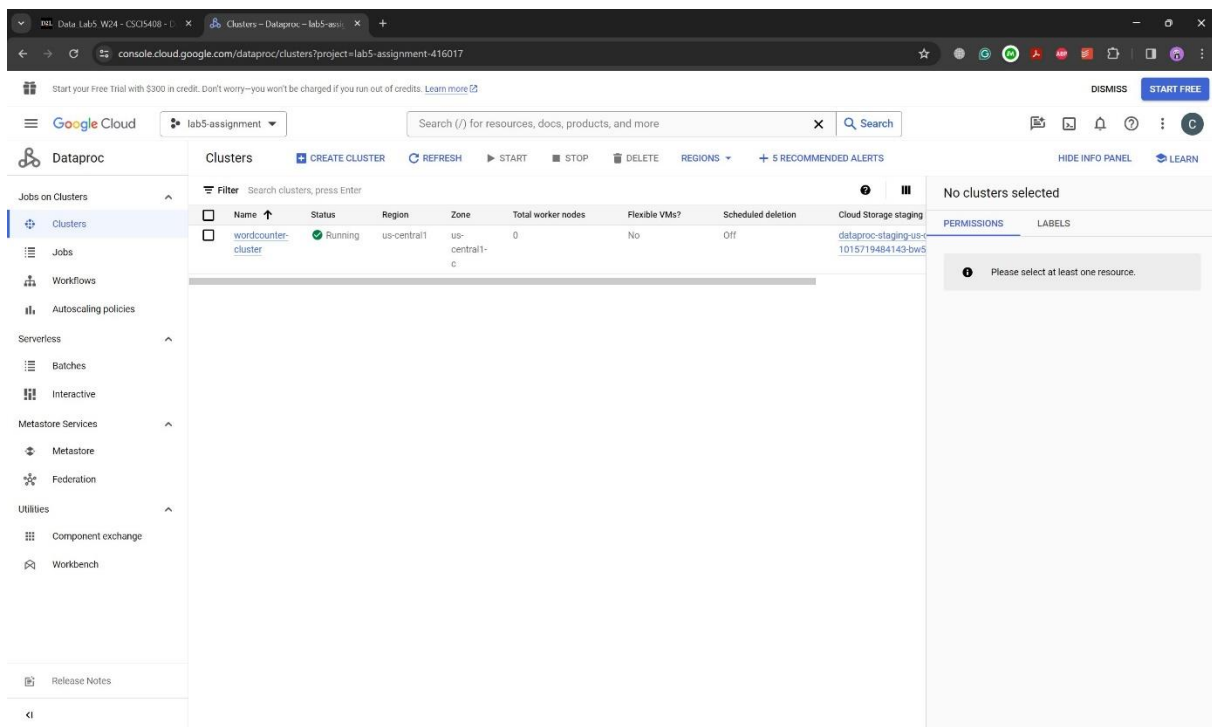*Figure 5 Settings for configure nodes*
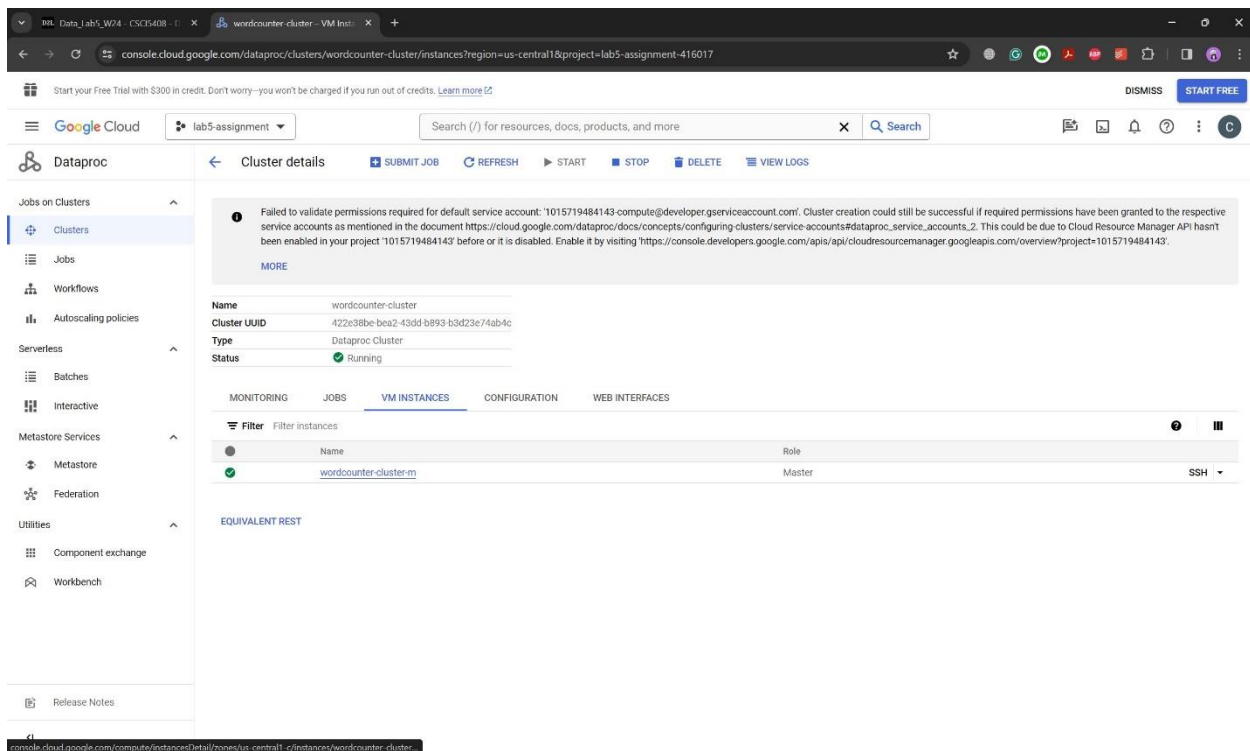


*Figure 6 Clusters menu*

*Figure 7 VM instances for wordcounter-cluster*
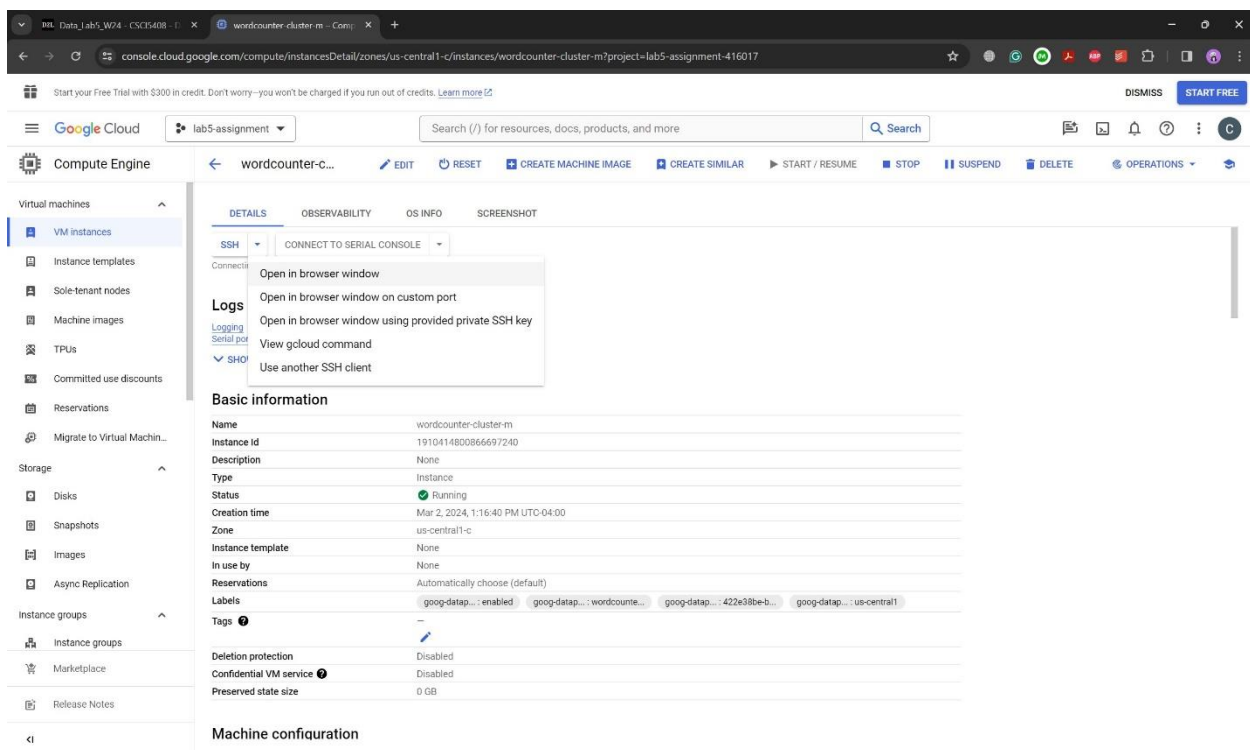


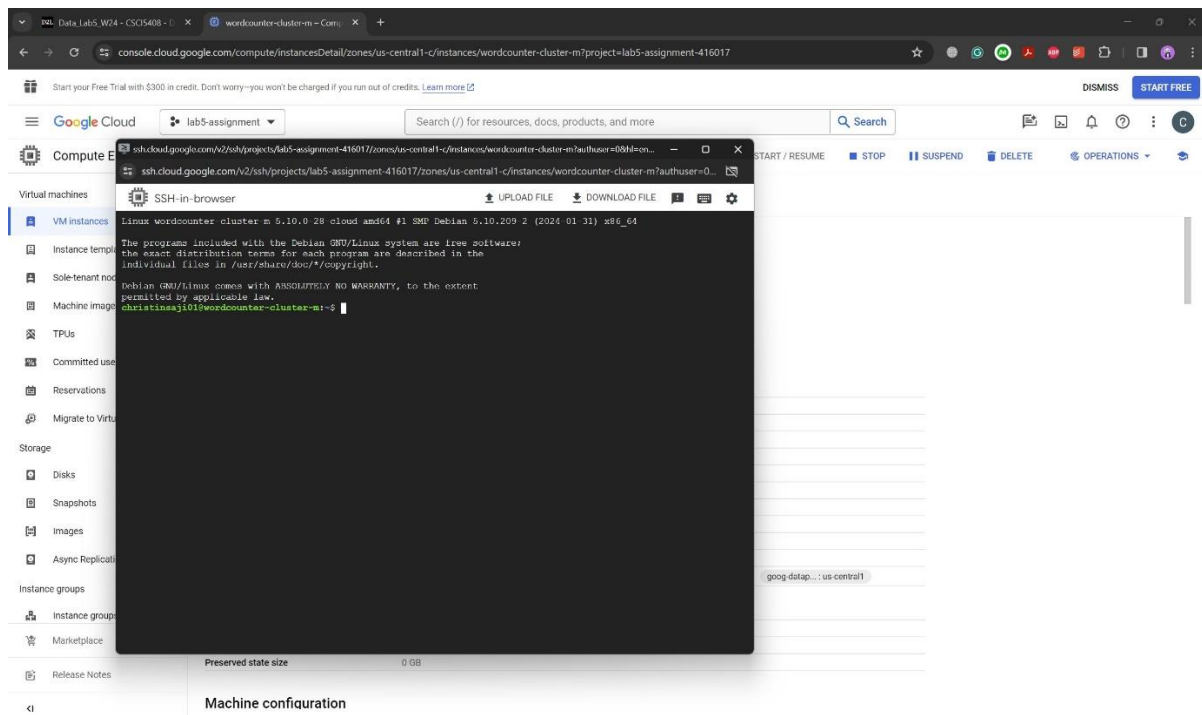*Figure 8 SSH in browser window*

*Figure 9 SSH-in browser terminal*

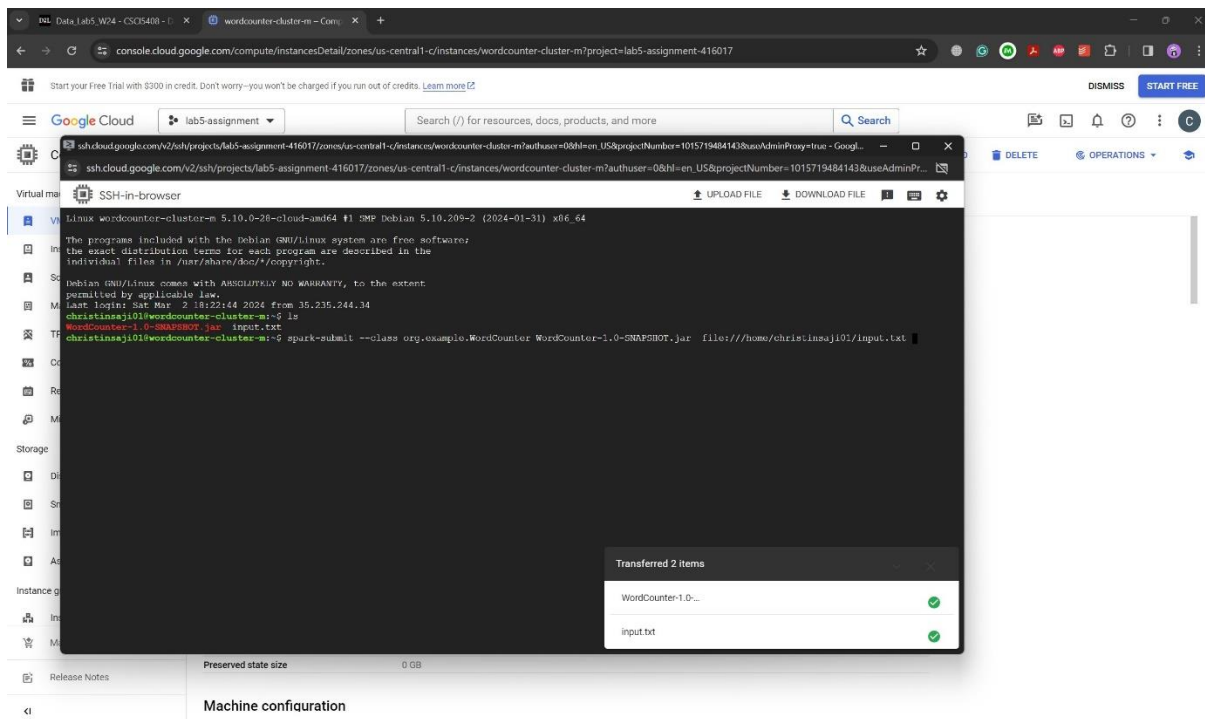I uploaded WordCounter-1.0-SNAPSHOT.jar and input.txt



*Figure 10 Uploaded files in the terminal*

Executed the WordCounter program using the command "spark-submit –class org.example.WordCounter WordCounter-1.0-SNAPSHOT.jar file:///home/christinsaji01/input.txt".



*Figure 11 Output of the WordCounter program*

# Task 2: Explanation of the Java Spark program with the screenshots of the code.
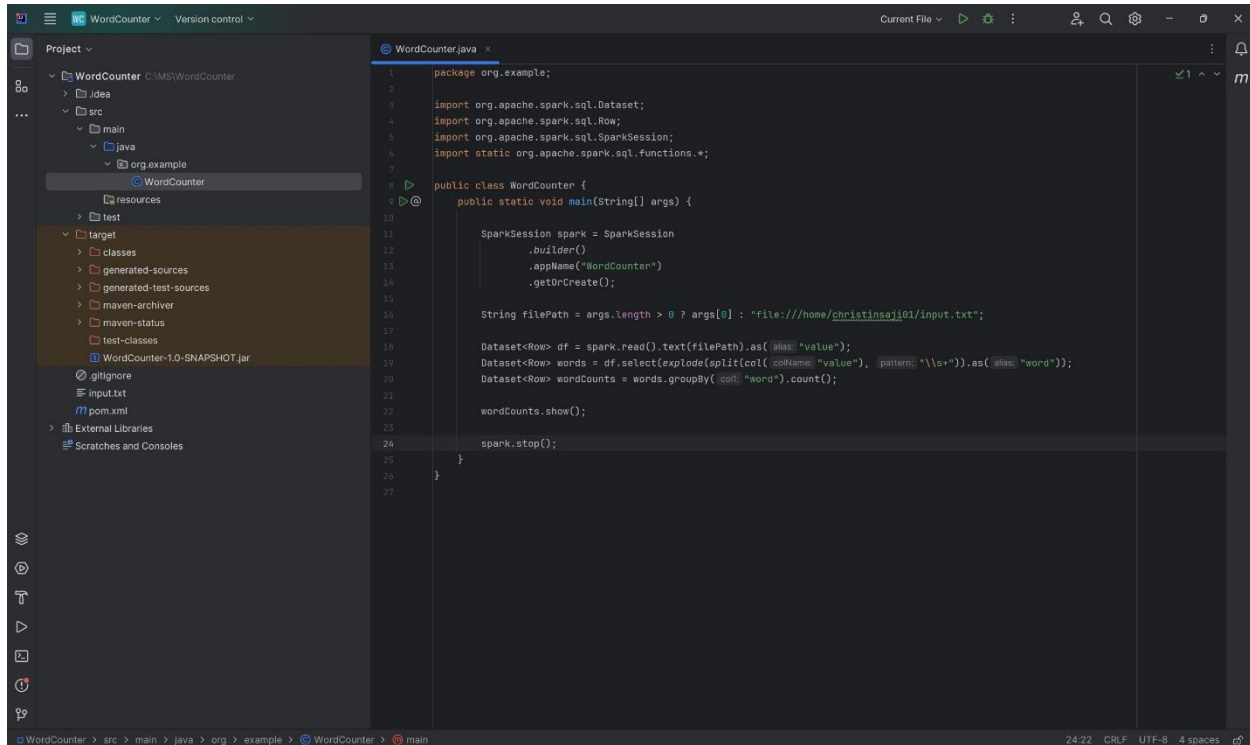


*Figure 12 Java Spark program to count words*

Step 1: First, I created an instance for "SparkSession" using the Dataset and DataFrame API.

Step 2: I created a filePath to fetch the file path from the command-line argument and used it; otherwise, I set a default path.

Step 3: The input text was read into the DataFrame, with each row containing a single column called "value."

Step 4: I split the single line using regex by whitespaces and used the "explode" function to separate each word into a separate row with a column named "words."

Step 5: I grouped the resulting DataFrame by the word and counted the occurrence of that word.

Step 6: I displayed the result using the "show" method.

Step 7: Finally, I used "stop" to terminate the Spark session.

# References

[1]    Naveen (NNK), "Spark Read Text File: RDD: DataFrame," *Spark By {Examples}*, [Online], Feb 8, 2023. Available: https://sparkbyexamples.com/spark/spark-read-text-file-rdd-dataframe/ [Accessed: March 2, 2024].

[2]    Singh, Chandan, "An Introduction to Apache Spark with Java," *Stack Abuse*, [Online] Aug 3, 2023. Available: https://stackabuse.com/an-introduction-to-apache-spark-with-java/ [Accessed: March 2, 2024].