

CSCI 5408

DATA MANAGEMENT AND
WAREHOUSING

LAB - 5

Banner ID: B00958098

GitLab Assignment Link:

https://git.cs.dal.ca/abhishekk/csci5408_w24_b00958098_abhishek_kapoor

Table of Contents

Deliverable #1	3
Direction.....	3
Deliverable #2	6
Direction.....	6
Deliverable #3	7
Direction.....	7
Deliverable #4	9
Direction.....	9
Deliverable #5	9
Direction.....	9

Deliverable #1

Direction

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.

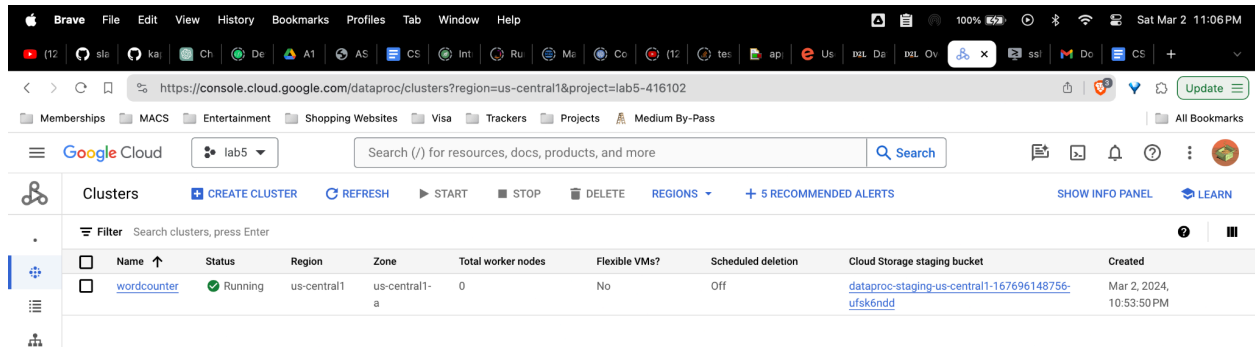


Fig. 1 WordCounter cluster being created

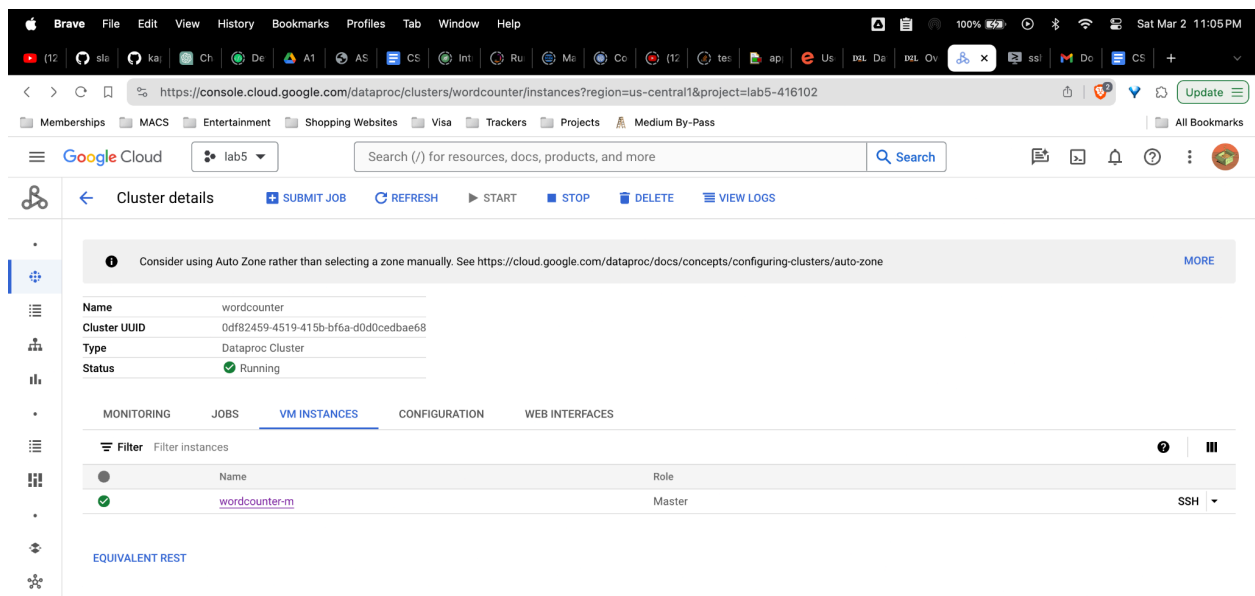


Fig. 2 Cluster Dashboard

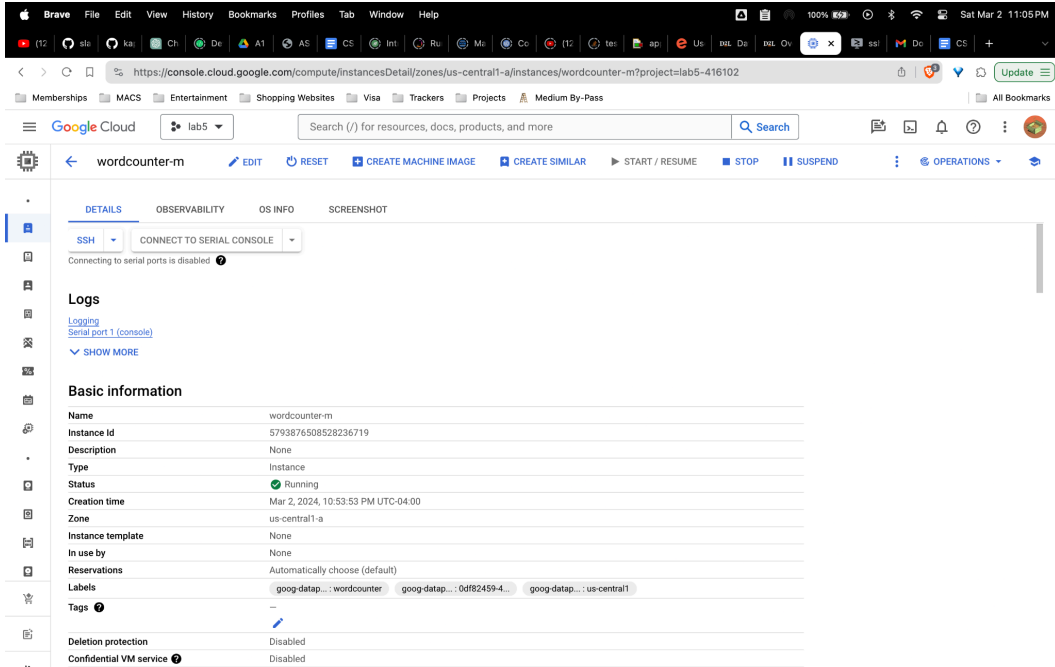


Fig. 3 Wordcounter VM instances

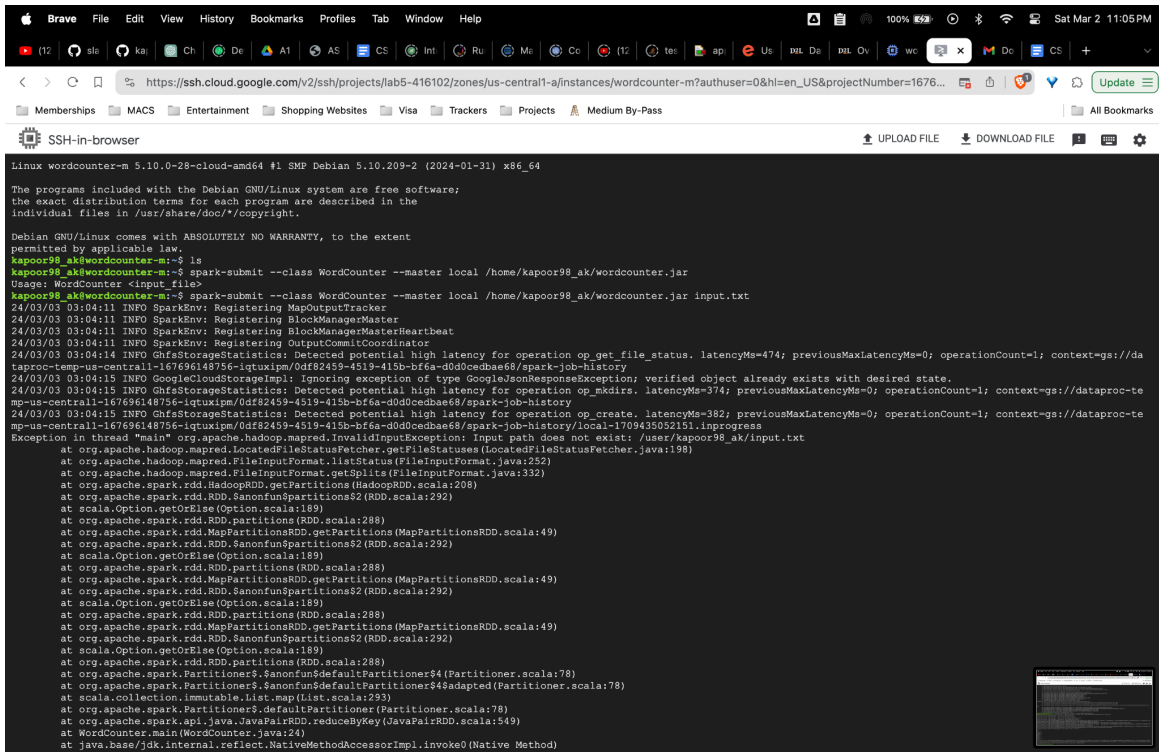


Fig. 4 Spark-submit command being run with appropriate flags to run the jar

```
at WordCounter.main(WordCounter.java:24)
at java.base/jdk.internal.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
at java.base/jdk.internal.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:62)
at java.base/jdk.internal.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
at java.base/java.lang.reflect.Method.invoke(Method.java:566)
at org.apache.spark.deploy.JavaMainApplication.start(SparkApplication.scala:52)
at org.apache.spark.deploy.SparkSubmit.org$apache$spark$deploy$SparkSubmit$$runMain(SparkSubmit.scala:973)
at org.apache.spark.deploy.SparkSubmit.doRunMain$1(SparkSubmit.scala:180)
at org.apache.spark.deploy.SparkSubmit.submit(SparkSubmit.scala:203)
at org.apache.spark.deploy.SparkSubmit.doSubmit(SparkSubmit.scala:90)
at org.apache.spark.deploy.SparkSubmit$anon$2.doSubmit(SparkSubmit.scala:1061)
at org.apache.spark.deploy.SparkSubmit$.main(SparkSubmit.scala:1070)
at org.apache.spark.deploy.SparkSubmit.main(SparkSubmit.scala)
Caused by: java.io.IOException: Input path does not exist: /user/kapoor98_ak/input.txt
at org.apache.hadoop.mapred.LocalizedFileStatusFetcher$ProcessInitialInputPathCallable.call(LocalizedFileStatusFetcher.java:402)
at org.apache.hadoop.mapred.LocalizedFileStatusFetcher$ProcessInitialInputPathCallable.call(LocalizedFileStatusFetcher.java:380)
at org.apache.hadoop.thirdparty.com.google.common.util.concurrent.TrustedListenableFutureTask$TrustedFutureInterruptibleTask.runInterruptibly(TrustedListenableFutureTask.java:1
25)
at org.apache.hadoop.thirdparty.com.google.common.util.concurrent.InterruptibleTask.run(InterruptibleTask.java:69)
at org.apache.hadoop.thirdparty.com.google.common.util.concurrent.TrustedListenableFutureTask.run(TrustedListenableFutureTask.java:78)
at java.base/java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1128)
at java.base/java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:628)
at java.base/java.lang.Thread.run(Thread.java:829)
kapoor98_ak@wordcounter-m:~$ hdfs dfs -copyFromLocal ./input.txt /user/kapoor98_ak/input.txt
24/03/03 03:05:07 INFO SparkEnv: Registering MapOutputTracker
24/03/03 03:05:07 INFO SparkEnv: Registering BlockManagerMaster
24/03/03 03:05:07 INFO SparkEnv: Registering BlockManagerMasterHeartbeat
24/03/03 03:05:07 INFO SparkEnv: Registering OutputCommitCoordinator
24/03/03 03:05:10 INFO GhsfsStorageStatistics: Detected potential high latency for operation op_get_file_status. latencyMs=427; previousMaxLatencyMs=0; operationCount=1; context-gs://da
dataproc-temp-us-central1-167696148756-igtuxipm/0df82459-4519-415b-bf6a-d0d0cebdae68/spark-job-history
24/03/03 03:05:10 INFO GoogleCloudStorageImpl: Ignoring exception of type GoogleJsonResponseException; verified object already exists with desired state.
24/03/03 03:05:10 INFO GhsfsStorageStatistics: Detected potential high latency for operation op_mkdirs. latencyMs=323; previousMaxLatencyMs=0; operationCount=1; context-gs://dataproc-te
mp-us-central1-167696148756-igtuxipm/0df82459-4519-415b-bf6a-d0d0cebdae68/spark-job-history
24/03/03 03:05:10 INFO GhsfsStorageStatistics: Detected potential high latency for operation op_create. latencyMs=416; previousMaxLatencyMs=0; operationCount=1; context-gs://dataproc-te
mp-us-central1-167696148756-igtuxipm/0df82459-4519-415b-bf6a-d0d0cebdae68/spark-job-history/local-1709435107653.inprogress
24/03/03 03:05:12 INFO FileInputFormat: Total input files to process : 1
fox: 1
the: 1
jumps: 1
quick: 1
lazy: 1
dog: 1
over: 1
brown: 1
tea: 1
24/03/03 03:05:15 INFO GhsfsStorageStatistics: Detected potential high latency for operation op_rename. latencyMs=340; previousMaxLatencyMs=0; operationCount=1; context-rename(gs://data
proc-temp-us-central1-167696148756-igtuxipm/0df82459-4519-415b-bf6a-d0d0cebdae68/spark-job-history/local-1709435107653.inprogress -> gs://dataproc-temp-us-central1-167696148756-igtuxip
m/0df82459-4519-415b-bf6a-d0d0cebdae68/spark-job-history/local-1709435107653)
kapoor98_ak@wordcounter-m:~$
```

Fig. 5 Expected Results shown in the output

Deliverable #2

Direction

Report any challenges faced while executing the .jar file on the Apache spark cluster.

During the lab, I faced these challenges:

1. Configuration and Setup:

- Setting up the Spark cluster environment correctly was a significant challenge as I could not follow properly in the Lab and then could not follow the Lab 5 PDF.
- It involved taking help from other classmates, and ensuring that all the servers in the cluster have the required dependencies and configurations.
- To not repeat this in future, I would test the setup on a small scale before deploying to a larger cluster.

2. Resource Management and Optimization:

- Efficiently managing and optimizing resources on a Spark cluster is crucial for performance and cost management.
- Issues can be resolved by fine-tuning the Spark configurations for memory, CPU, and other resource optimization.
- I utilized Spark's built-in monitoring to monitor resource utilization.

3. Debugging and Logging:

- Debugging issues on a distributed Spark cluster was challenging due to the distributed nature of the computation. Identifying and resolving errors, especially in a production environment, could be very time-consuming without proper logging and debugging mechanisms.
- Enabling logging can help us to monitor job progress and identify errors. Incorporate proper error handling and logging within your application code to aid in troubleshooting.

Deliverable #3

Direction

Explanation of the Java Spark program with the screenshots of the code.

Main Class - WordCounter:

1. This class serves as the entry point for the application.
2. It checks if the input file path is provided as a command-line argument. If not, it prints an error message and exits.
3. If the input file path is provided, it creates an object of WordCounterProcessor and calls the process method to handle the word counting.

Processor Class - WordCounterProcessor:

1. This class has the processing logic for word counting using Apache Spark.
2. Methods:
 - a. **process():**
 - Takes the input file path as a parameter.
 - Initializes Spark configuration and context using helper methods.
 - Read the input file, split lines into words, map words to tuples, reduce word counts, collect results, print results, and finally stop the Spark context.
3. Helper Methods:
 - a. **createSparkConf():**
 - Creates a Spark configuration with the application name set to "WordCounter" and the master set to "local."
 - b. **readInputFile():**
 - Reads the input file using the provided Spark context and returns a JavaRDD<String> containing lines from the file.
 - c. **splitLinesIntoWords():**
 - Takes a JavaRDD<String> of lines, splits each line into words, and returns a new JavaRDD<String> containing individual words.
 - d. **mapWordsToTuples():**
 - Takes a JavaRDD<String> of words and maps each word to a tuple (word, 1) using mapToPair, resulting in a JavaPairRDD<String, Integer>.
 - e. **reduceWordCountsByKey():**
 - Takes a JavaPairRDD<String, Integer> and reduces the tuples by key (word) by summing up the values, resulting in a new JavaPairRDD<String, Integer>.
 - f. **collectResults():**
 - Takes the reduced word counts as a JavaPairRDD<String, Integer> and collects the results as a list of tuples.
 - g. **printResults():**
 - Takes the list of tuples and prints each word and its count to the console.

This approach follows the SOLID principles by separating concerns and ensuring that each method has a single responsibility. It improves code readability by using meaningful names. The Spark-specific logic is encapsulated within the WordCounterProcessor class.

Deliverable #4

Direction

Provide references of any sources used for this lab.

References provided in the report.

Deliverable #5

Direction

Submit code on Gitlab.

GitLab Repository:

https://git.cs.dal.ca/abhishekk/csci5408_w24_b00958098_abhishek_kapoor