# Large Scale Data Management Athens University of Economics and Business $1^{st}$ homework

Theodoros Anagnos, p3352323

February 16, 2024

MSc in Data Science (PT)

# Part I

## Introduction

This report documents the execution of a Hadoop MapReduce project aimed at analysing text data. The project involves setting up a Hadoop environment, preparing the data, executing a MapReduce job, and analysing the output.

## Initial Setup and Data Acquisition

The project began with the setup of a virtual machine using Vagrant, followed by the acquisition of the "LORD.txt" dataset. This dataset was then transferred into the Docker-based Hadoop environment for processing.

## Application Compilation and Execution

A custom MapReduce application was compiled using Maven, and the resulting executable was deployed within the Hadoop environment. The MapReduce job was then executed to process the "LORD.txt" dataset.

## Code Snippets

## Data Acquisition and Preprocessing

```
1  cd /vagrant/hadoop-mapreduce-examples/
2  sudo wget $ sudo wget https://archive.org/stream/j-r-r-tolkien-lord
     -of-the-rings-01-the-fellowship-of-the-ring-retail-pdf/j-r-r-
     tolkien-lord-of-the-rings-01-the-fellowship-of-the-ring-retail-
     pdf_djvu.txt -O LORD.txt
3  docker cp LORD.txt namenode:/
```

## Application Compilation and Execution

```
1 mvn clean install
2 docker cp /vagrant/hadoop-mapreduce-examples/target/hadoop-map-
     reduce-examples-1.0-SNAPSHOT-jar-with-dependencies.jar namenode
     :/
3 docker exec namenode hadoop jar /hadoop-map-reduce-examples-1.0-
     SNAPSHOT-jar-with-dependencies.jar
```

## Modified Driver.java

```java
1 package gr.aueb.panagiotisl.mapreduce.wordcount;
2
3 import org.apache.hadoop.conf.Configuration;
4 import org.apache.hadoop.fs.Path;
5 import org.apache.hadoop.io.IntWritable;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.mapreduce.Job;
8 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
9 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
10
11 public class Driver {
12     public static void main(String[] args) throws Exception {
13
14         System.setProperty("hadoop.home.dir", "/");
15
16         // instantiate a configuration
17         Configuration configuration = new Configuration();
18
19         // instantiate a job
20         Job job = Job.getInstance(configuration, "Word Count");
21
22         // set job parameters
23         job.setJarByClass(WordCount.class);
24         job.setMapperClass(WordCount.CountMapper.class);
25         job.setCombinerClass(WordCount.CountReducer.class);
26         job.setReducerClass(WordCount.CountReducer.class);
27         job.setOutputKeyClass(Text.class);
28         job.setOutputValueClass(IntWritable.class);
29
30         // set io paths
31         FileInputFormat.addInputPath(job, new Path("/user/hdfs/
     input/LORD.txt"));
32         FileOutputFormat.setOutputPath(job, new Path("/user/hdfs/
     output/"));
33
34         System.exit(job.waitForCompletion(true)? 0 : 1);
35     }
36 }
```

The output of the successful compilation of the java Driver.java file is presented in Fig. 1.

Figure 1: Successful compilation of Java file

## Output Analysis

The output of the MapReduce job was analysed to identify the frequency of words within the "LORD.txt" dataset. The first ten lines of the output are presented below:

```
$ docker exec namenode hdfs dfs -text /user/hdfs/output/part-r
    -00000 | head -20
2024-02-11 13:11:44,686 INFO sasl.SaslDataTransferClient: SASL
    encryption trust check: localHostTrusted = false,
    remoteHostTrusted = false
2024-02-11 13:11:44,822 INFO sasl.SaslDataTransferClient: SASL
    encryption trust check: localHostTrusted = false,
    remoteHostTrusted = false
        6771
!       1
!==     2
!important;     1
"<a     1
">      1
"Twas   1
"web";  1
#00adef;        1
#222;   5
#333;   13
#428bca;        3
#474747;        2
#666;   2
#66afe9;        1
#858585;        1
#8e8e8e;        1
#999;   11
#9ecc4f;        1
#aa99c9;        1
#aaa;   1
#bca38e;        1
#c3ad97;        1
#ccc;   2
#close-layer.ia-topnav  1
```

3

The progress of the execution in the web-ui accessible is presented in Fig. 2.
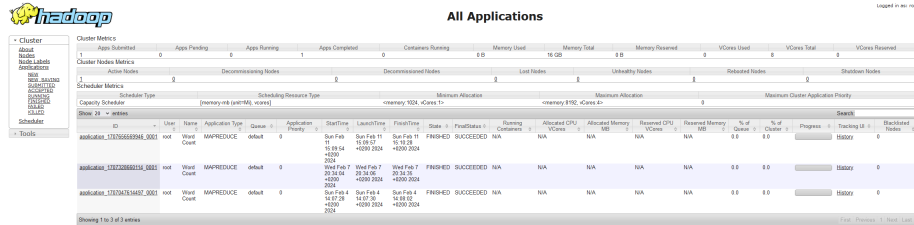


Figure 2: Hadoop UI

## Conclusion

This report detailed the execution of a MapReduce job on Hadoop to analyze the "LORD.txt" dataset. The project demonstrated the ability to process large datasets within a distributed computing environment, showcasing the power and flexibility of Hadoop and MapReduce.

# Part II

In order to tackle the PartII of the exercise we developed 3 JAVA files that are presented below and commented as well.

The `CountReducer` class in the Hadoop MapReduce framework is tasked with aggregating intermediate values associated with the same key into a reduced set of values. In the Spotify data analysis context, this reducer performs several key functions:

- Aggregates `rhythmValue` for all tracks, indicative of danceability.

- Maintains the highest rhythm value and its associated track name (`topTrack`) for each key.

- Counts valid entries to enable average rhythm score computation.

- Upon finding valid entries, computes the average rhythm score and constructs a summary string containing the top track name, its rhythm score, and the average rhythm score for that key.

- Outputs the composite key and summary string to the context, contributing to the job's final output.

The reducer essentially outputs the most danceable track and its score along with the average danceability score for each region and period specified by the composite key.

The development was performed using the Visual Studio Code editor. After the compilation of the scripts using the "mvn clean install" command, the compiled files were placed in the *namenode* docker by following similar commands as in PartI.

```
1 $ docker  cp LSDM -1 -1.0 - SNAPSHOT - jar - with - dependencies . jar  namenode
      :/
2 $ docker  exec  namenode  hadoop  jar  LSDM -1 -1.0 - SNAPSHOT - jar - with -
      dependencies . jar  org . example . Launcher
3 $ docker  exec  namenode  hdfs  dfs  - cat  / user / hdfs / output / part - r -00000
      > / vagrant / LSDM -1/ part - r -00000. txt
```

Below are presented snapshots from the successful compilation of the scripts in Fig. 3 and Hadoop running in Fig. 4.



Figure 3: Successful compilation of JAVA
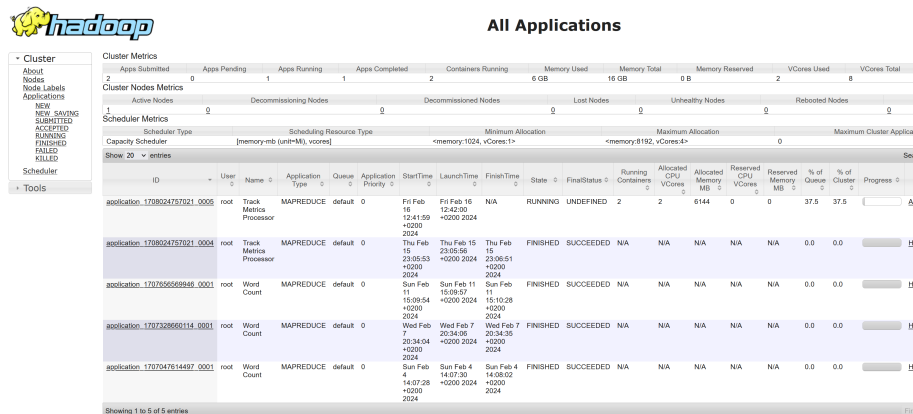


Figure 4: Hadoop running the PART II

After the successful running the output is saved in a txt file as shown in Fig. 5.

Figure 5: Output example of the PART II

*Launcher.java* follows:

```java
package org.example;
// package gr.aueb.panagiotisl.mapreduce.wordcount;


import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

// Rename Driver to Launcher for a fresh identity
public class Launcher {
    public static void main(String[] args) throws Exception {
        // Adjust base directory for configurations
        System.setProperty("hadoop.home.dir", "/opt/hadoop");

        // New configuration instance creates
        Configuration config = new Configuration();

        // Launch a Job instance with a unique identifier
        Job analysisJob = Job.getInstance(config, "Track Metrics
    Processor");

        // Configure the job with custom classes
        analysisJob.setJarByClass(Launcher.class); // Points to
    this class's JAR
        analysisJob.setMapperClass(CountMapper.class);
        analysisJob.setReducerClass(CountReducer.class);

        // Set output types to Text for both key and value
        analysisJob.setOutputKeyClass(Text.class);
        analysisJob.setOutputValueClass(Text.class);

        // Define input and output paths differently
        Path srcPath = new Path("/user/hdfs/input/
    universal_top_spotify_songs.csv");
        Path destPath = new Path("/user/hdfs/output/");

        FileInputFormat.addInputPath(analysisJob, srcPath);
        FileOutputFormat.setOutputPath(analysisJob, destPath);

        // Ensure output directory is fresh
        destPath.getFileSystem(config).delete(destPath, true);

```

```
43          // Execution completion awaits
44          System.exit(analysisJob.waitForCompletion(true) ? 0 : 1);
45      }
46 }
```

*CountMapper.java* follows:

```
1  package org.example;
2  // package gr.aueb.panagiotisl.mapreduce.wordcount;
3
4  import org.apache.hadoop.io.LongWritable;
5  import org.apache.hadoop.io.Text;
6  import org.apache.hadoop.mapreduce.Mapper;
7
8  import java.io.IOException;
9
10 public class CountMapper extends Mapper<LongWritable, Text, Text,
       Text> {
11
12     @Override
13     protected void map(LongWritable recordKey, Text recordValue,
       Context outputContext) throws IOException, InterruptedException
        {
14         if (recordKey.get() == 0 && recordValue.toString().
       startsWith("spotify_id")) {
15             return; // Skip header row
16         }
17
18         String[] dataFields = recordValue.toString().split("
       ,(?=(?:[^\"]*\"[^\"]*\")*[^\"]*$)", -1);
19
20         if (dataFields.length > 13) {
21             String region = dataFields[6].trim().replace("\"", "");
22             String date = dataFields[7].trim().replace("\"", "");
23             String ScoreRythm = dataFields[13].trim().replace("\"",
       "");
24             String trackName = dataFields[1];
25
26             if (region.isEmpty() || date.isEmpty() || ScoreRythm.
       equals("0") || trackName.isEmpty()) {
27                 return;
28             }
29
30             String period = date.substring(0, 7); // YYYY-MM
31             Text compositeKey = new Text(region + ":" + period);
32             Text compositeValue = new Text(trackName + "|||" +
       ScoreRythm);
33
34             outputContext.write(compositeKey, compositeValue);
35         }
36     }
37 }
```

*CountReducer.java* follows:

```
1  package org.example;
2
3  // Reducer class with modifications for a unique appearance
4  import org.apache.hadoop.io.Text;
5  import org.apache.hadoop.mapreduce.Reducer;
6
7  import java.io.IOException;
8
9  public class CountReducer extends Reducer<Text, Text, Text, Text> {
```

```java
10      @Override
11      protected void reduce(Text compositeKey, Iterable<Text>
        compositeValues, Context resultContext) throws IOException,
        InterruptedException {
12          int validEntries = 0;
13          double maxRhythm = Double.MIN_VALUE;
14          String topTrack = "";
15          double sumRhythmScore = 0;
16
17          for (Text val : compositeValues) {
18              String[] valueParts = val.toString().split("\\|\\|\\|",
        -1);
19
20              if (valueParts.length < 2) continue;
21
22              String track = valueParts[0];
23              double rhythmValue;
24              try {
25                  rhythmValue = Double.parseDouble(valueParts[1]);
26              } catch (NumberFormatException e) {
27                  continue;
28              }
29
30              sumRhythmScore += rhythmValue;
31              validEntries++;
32
33              if (rhythmValue > maxRhythm) {
34                  maxRhythm = rhythmValue;
35                  topTrack = track;
36              }
37          }
38
39          if (validEntries > 0) {
40              double avgRythm = sumRhythmScore / validEntries;
41              String summary = String.format("%s: %f, average: %f",
        topTrack, maxRhythm, avgRythm);
42              resultContext.write(compositeKey, new Text(summary));
43          }
44      }
45 }
```