
Programming for Big Data – Assignment 1

Hadoop/Map-Reduce

Ciaran Finnegan – TU060 (Part Time) – Second
Year

PROG9813: 2022/2023

MSc in Computer Science (Data Science)

Student No : D21124026

26/2/2023

Table of Contents

1	Introduction	3
1.1	Purpose of Report	3
2	Loading Data Into HDFS	4
2.1	Setting Up Shared Folders	4
2.2	Language Files Source (Books).....	4
2.3	Transferring and Loading Data Files.....	5
3	Map-Reduce Process Design	7
3.1	High Level Overview.....	7
3.2	Explanation of Map-Reduce Workflow	8
3.3	Key Design Decisions and Assumptions.....	15
4	Map-Reduce: Java Code.....	17
4.1	Driver.....	17
4.2	Job 1	20
4.3	Job 2	23
5	Python: Graph Analysis of MR Outputs	26
5.1	Set Up Excel From HSFs Output	26
5.2	Python Jupyter Notebook	27
5.3	Analysis Commentary.....	33



1 Introduction

1.1 Purpose of Report

This report collates all the delivery requirements for Assignment 1 in the Programming for Big Data module (PROG9813: 2022-23).

The primary content of the report covers;

- How language data (book files) were sourced and then loaded into a HDFS data store on a Linux VM running Hadoop (the VM provided during lecture 1).
- How the Map-Reduce workflow is constructed to ingest the language files stored in HDFS and ultimately compile a breakdown of the average character frequency in each language.
- Design assumptions and decisions made during the implementation of this PBD assignment.
- The Java source written to build the Map-Reduce processes, and how the features offered by Map-Reduce were employed to solve the assignment challenge.
- Python source code and graphs (via a Jupyter Notebook) of the analysis on the output from the Map-Reduce process.

Supporting source code and HDFS outputs accompany this report as part of the overall Brightspace submission.

2 Loading Data Into HDFS

2.1 Setting Up Shared Folders

The VM provided during lectures was used as the Hadoop environment for this assignment. Following lecturer advice, a shared folder on the host machine was set up.



Figure – VM Shared Folder Set Up

2.2 Language Files Source (Books)

Six books were selected from the <http://www.gutenberg.org> website;

- Two in Spanish.
- Two in German.
- Two in Italian.

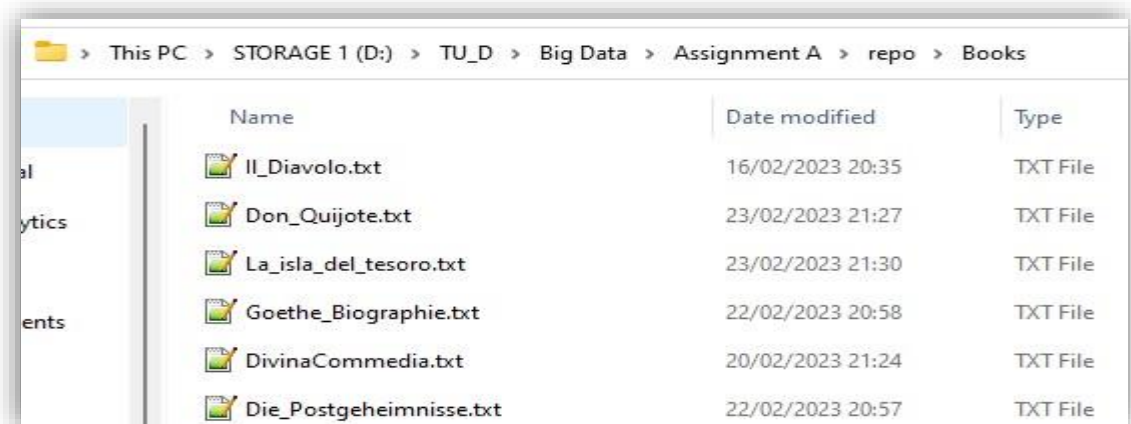
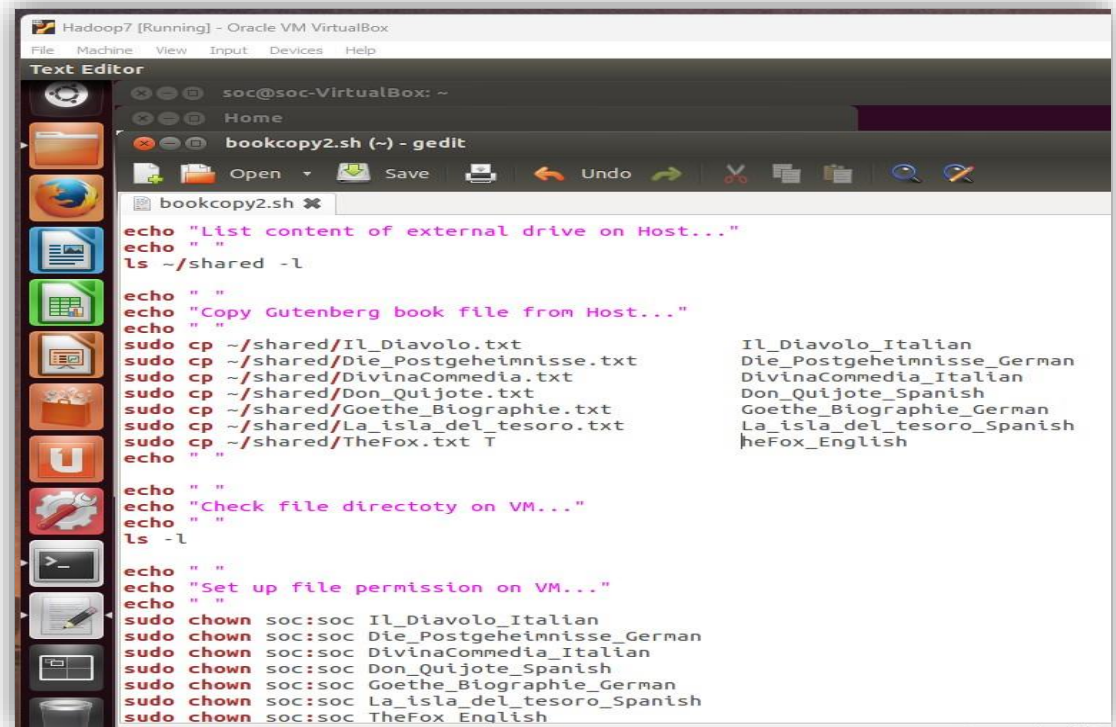


Figure – Book Files Extracted From Gutenberg Site

2.3 Transferring and Loading Data Files

The filenames of each book were altered in the transfer process from host machine to VM to add a suffix indicating language.



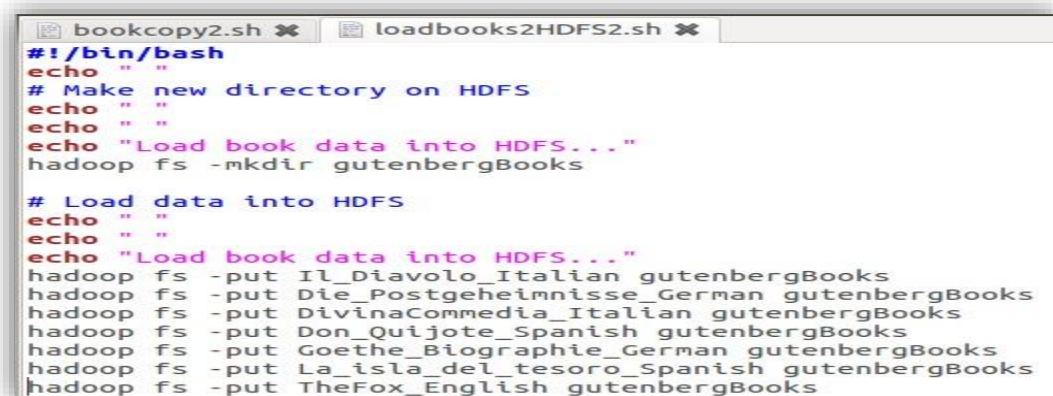
```

Hadoop7 [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
Text Editor
soc@soc-VirtualBox: ~
bookcopy2.sh (~) - gedit
bookcopy2.sh
echo "List content of external drive on Host..."
echo " "
ls ~/shared -l
echo " "
echo "Copy Gutenberg book file from Host..."
echo " "
sudo cp ~/shared/Il_Diavolo.txt Il_Diavolo_Italian
sudo cp ~/shared/Die_Postgeheimnisse.txt Die_Postgeheimnisse_German
sudo cp ~/shared/DivinaCommedia.txt DivinaCommedia_Italian
sudo cp ~/shared/Don_Quijote.txt Don_Quijote_Spanish
sudo cp ~/shared/Goethe_Biographie.txt Goethe_Biographie_German
sudo cp ~/shared/La_isla_del_tesoro.txt La_isla_del_tesoro_Spanish
sudo cp ~/shared/TheFox.txt TheFox_English
echo " "
echo "Check file directoty on VM..."
echo " "
ls -l
echo " "
echo "Set up file permission on VM..."
echo " "
sudo chown soc:soc Il_Diavolo_Italian
sudo chown soc:soc Die_Postgeheimnisse_German
sudo chown soc:soc DivinaCommedia_Italian
sudo chown soc:soc Don_Quijote_Spanish
sudo chown soc:soc Goethe_Biographie_German
sudo chown soc:soc La_isla_del_tesoro_Spanish
sudo chown soc:soc TheFox_English

```

Figure – Shell Script with Linux Commands to Copy/Rename Files in VM

The renaming was done as part of the design assumptions covered in section 3.3 of this report. Setting up a book directory and loading the files into HDFS was done through the following Linux commands;



```

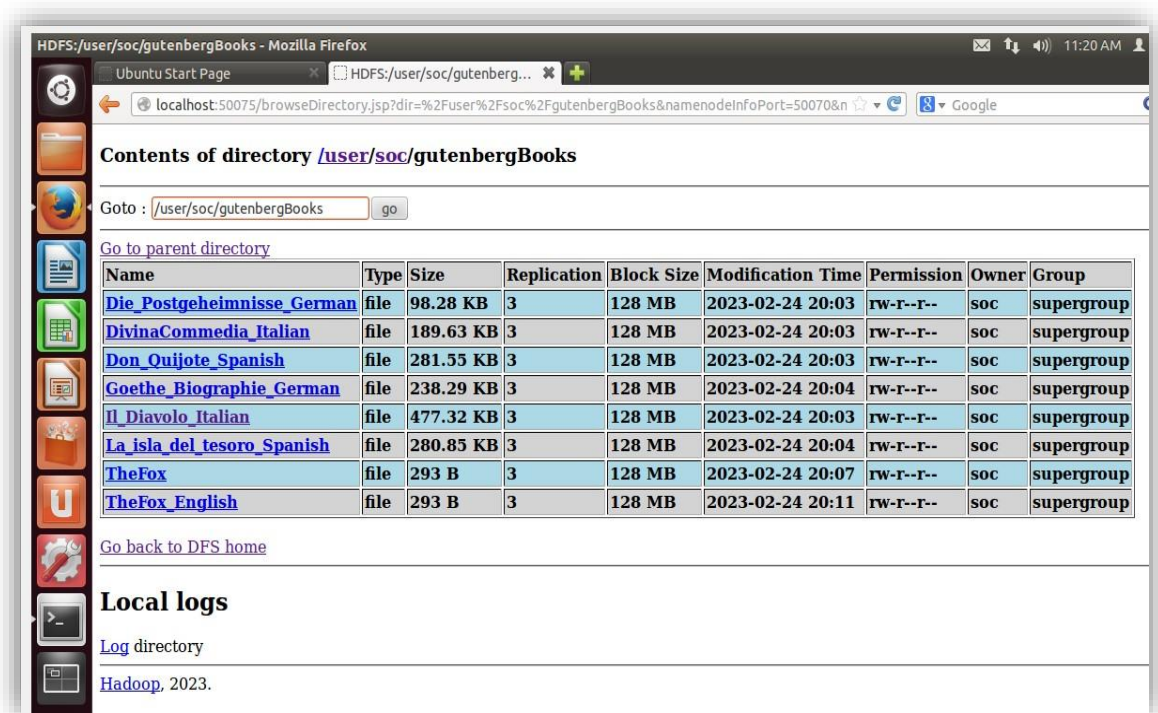
bookcopy2.sh loadbooks2HDFS2.sh
#!/bin/bash
echo " "
# Make new directory on HDFS
echo " "
echo "Load book data into HDFS..."
hadoop fs -mkdir gutenbergsBooks

# Load data into HDFS
echo " "
echo " "
echo "Load book data into HDFS..."
hadoop fs -put Il_Diavolo_Italian gutenbergsBooks
hadoop fs -put Die_Postgeheimnisse_German gutenbergsBooks
hadoop fs -put DivinaCommedia_Italian gutenbergsBooks
hadoop fs -put Don_Quijote_Spanish gutenbergsBooks
hadoop fs -put Goethe_Biographie_German gutenbergsBooks
hadoop fs -put La_isla_del_tesoro_Spanish gutenbergsBooks
hadoop fs -put TheFox_English gutenbergsBooks

```

Figure – Shell Script with Linux Commands to Copy Files to HDFS

The Hadoop interface confirmed that the language book files were successfully loaded.



Contents of directory **/user/soc/gutenbergBooks**

Goto : go

[Go to parent directory](#)

Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
Die_Postgeheimnisse_German	file	98.28 KB	3	128 MB	2023-02-24 20:03	rw-r--r--	soc	supergroup
DivinaCommedia_Italian	file	189.63 KB	3	128 MB	2023-02-24 20:03	rw-r--r--	soc	supergroup
Don_Quijote_Spanish	file	281.55 KB	3	128 MB	2023-02-24 20:03	rw-r--r--	soc	supergroup
Goethe_Biographie_German	file	238.29 KB	3	128 MB	2023-02-24 20:04	rw-r--r--	soc	supergroup
Il_Diavolo_Italian	file	477.32 KB	3	128 MB	2023-02-24 20:03	rw-r--r--	soc	supergroup
La_isla_del_tesoro_Spanish	file	280.85 KB	3	128 MB	2023-02-24 20:04	rw-r--r--	soc	supergroup
TheFox	file	293 B	3	128 MB	2023-02-24 20:07	rw-r--r--	soc	supergroup
TheFox_English	file	293 B	3	128 MB	2023-02-24 20:11	rw-r--r--	soc	supergroup

[Go back to DFS home](#)

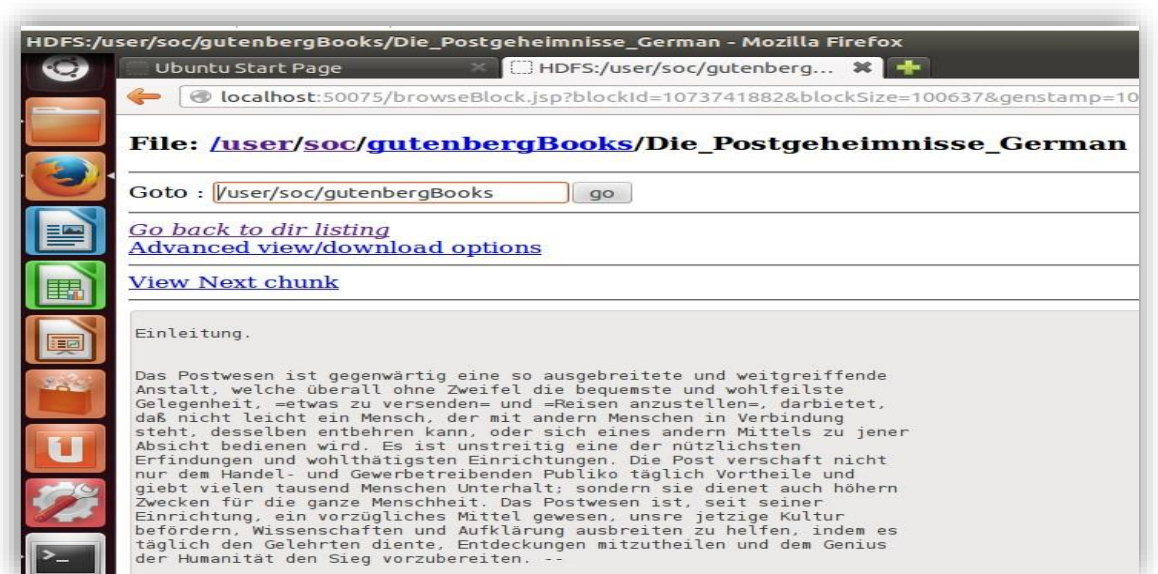
Local logs

[Log](#) directory

[Hadoop](#), 2023.

Figure –HDFS Directory Populated with Book Files

A quick inspection (example below) showed that the data was intact and ready for processing.



File: **/user/soc/gutenbergBooks/Die_Postgeheimnisse_German**

Goto : go

[Go back to dir listing](#)

[Advanced view/download options](#)

[View Next chunk](#)

Einleitung.

Das Postwesen ist gegenwärtig eine so ausgebreitete und weitgreifende Anstalt, welche überall ohne Zweifel die bequemste und wohlfeilste Gelegenheit, =etwas zu versenden= und =Reisen anzustellen=, darbietet, daß nicht leicht ein Mensch, der mit andern Menschen in Verbindung steht, desselben entbehren kann, oder sich eines andern Mittels zu jener Absicht bedienen wird. Es ist unstreitig eine der nützlichsten Erfindungen und wohlthätigsten Einrichtungen. Die Post verschafft nicht nur dem Handel- und Gewerbetreibenden Publikum täglich Vortheile und giebt vielen tausend Menschen Unterhalt; sondern sie dienet auch höhern Zwecken für die ganze Menschheit. Das Postwesen ist, seit seiner Einrichtung, ein vorzügliches Mittel gewesen, unsre jetzige Kultur befördern, Wissenschaften und Aufklärung ausbreiten zu helfen, indem es täglich den Gelehrten diene, Entdeckungen mitzutheilen und dem Genius der Humanität den Sieg vorzubereiten. --

Figure – Example of a German Book File 'Chunk' Stored in HDFS

3 Map-Reduce Process Design

3.1 High Level Overview

This diagram displays the high-level Map-Reduce process built and deployed for the assignment.

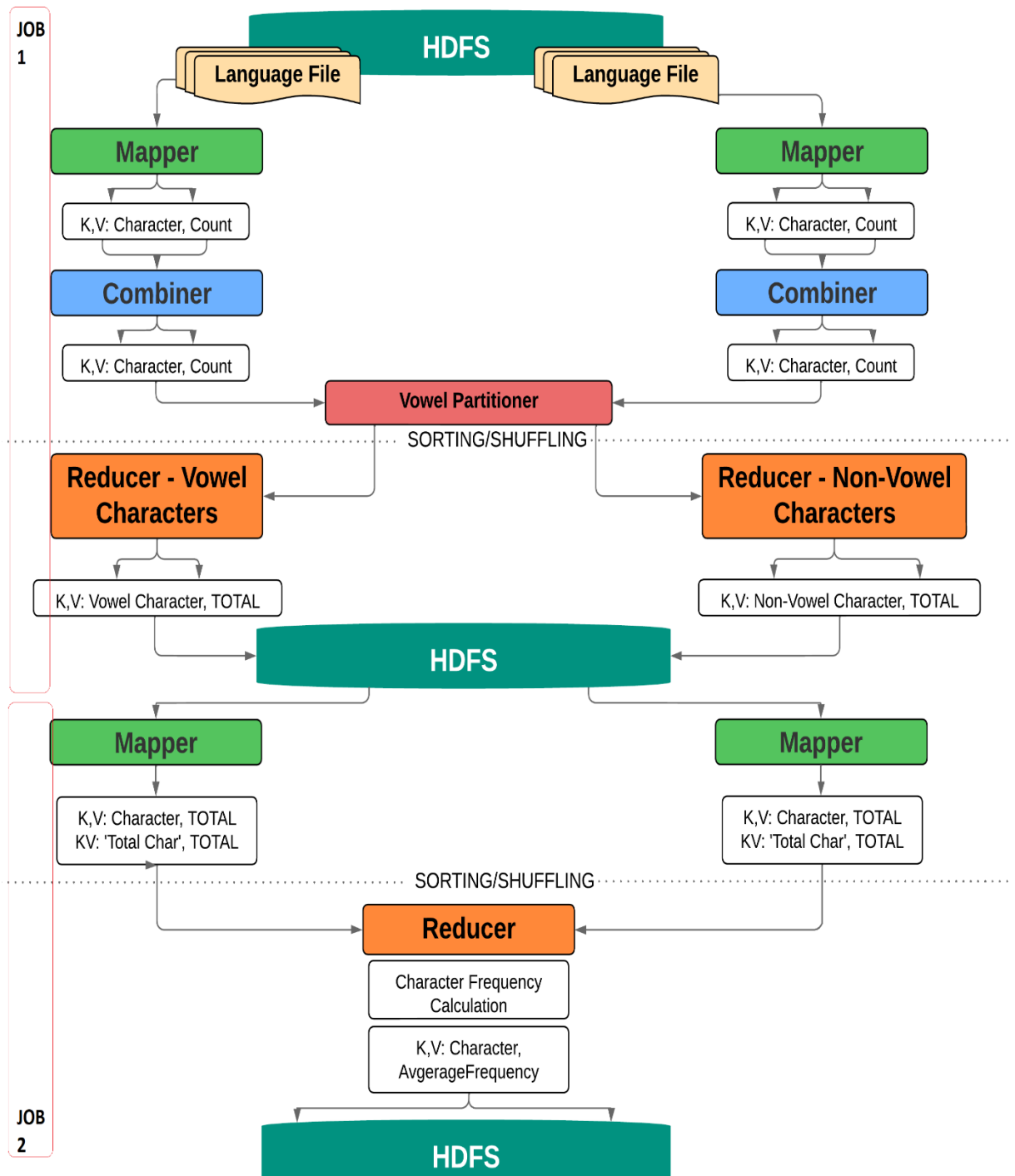


Figure – Map-Reduce Process Overview to Calculate Language Letter Frequency

The following sections of this report elaborate on the individual Map/Reduce processes.

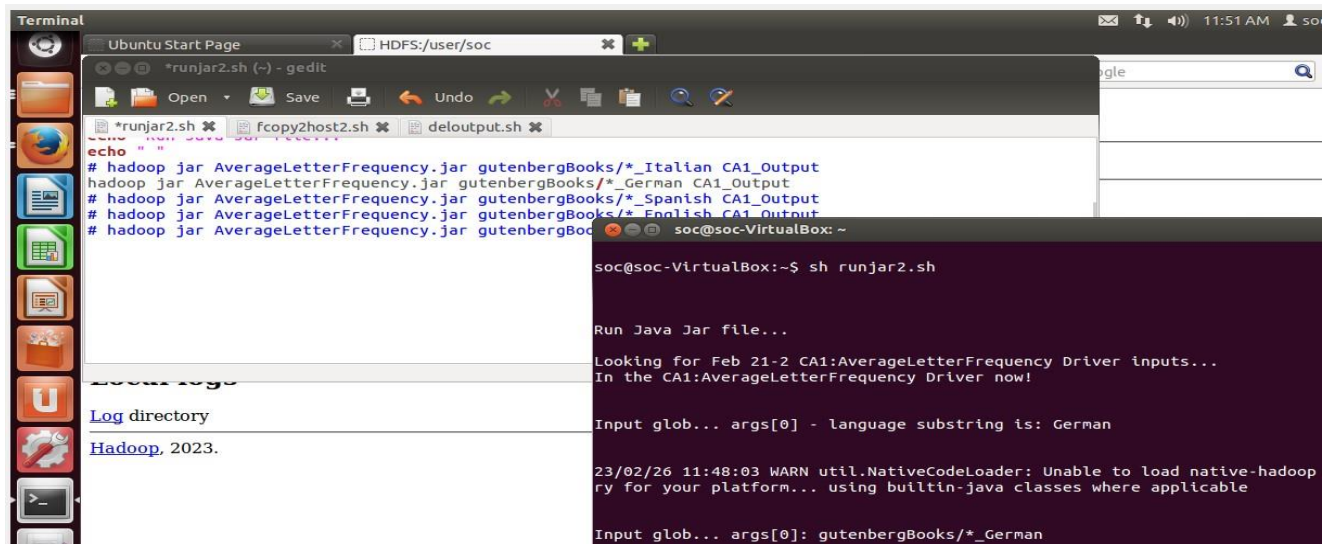
3.2 Explanation of Map-Reduce Workflow

The Java code written for this assignment, compiled in Eclipse under a jar file called `AverageLetterFrequency`, is provided in Section 4 of this report. This section (Section 3.2) provides a generalised overview of the steps in the assignment Map-Reduce process.

3.2.1 Driver Process

There is one Driver process for the entire assignment. Key steps in the process are;

- Invoke jar with command line parameters. This is the main Java project class programme. The `AverageLetterFrequency` jar file is invoked with two parameters.
 - The first is the location of the language files in the HDFS repository.
 - The second is the output directory to be created in HDFS to store the output from the (job1) Map-Reduce process.



```

Terminal
Ubuntu Start Page
HDFS:/user/soc
*runjar2.sh (~) - gedit
Open Save Undo Redo
*runjar2.sh *fcopy2host2.sh *deloutput.sh
echo
# hadoop jar AverageLetterFrequency.jar gutenbergsBooks/*_Italian CA1_Output
hadoop jar AverageLetterFrequency.jar gutenbergsBooks/*_German CA1_Output
# hadoop jar AverageLetterFrequency.jar gutenbergsBooks/*_Spanish CA1_Output
# hadoop jar AverageLetterFrequency.jar gutenbergsBooks/*_English CA1_Output
soc@soc-VirtualBox: ~
soc@soc-VirtualBox:~$ sh runjar2.sh

Run Java Jar file...

Looking for Feb 21-2 CA1:AveragelLetterFrequency Driver inputs...
In the CA1:AveragelLetterFrequency Driver now!

Input glob... args[0] - language substring is: German

23/02/26 11:48:03 WARN util.NativeCodeLoader: Unable to load native-hadoop
ry for your platform... using builtin-java classes where applicable

Input glob... args[0]: gutenbergsBooks/*_German
  
```

Figure – Cmd Line Parameters to Invoke Jar file for Map-Reduce Assignment

- Use Wildcards to select a subset of language specific book files stored in HDFS (see figure above). This glob is used to set the input path of the first job.
- File name/format is used to indicate language. All files processed by the Map-Reduce process must be in the file name format: **<bookname>_<language>**. The MR process will read in all files with the given language name used in the command line parameter invoking the jar file (see figure above).
- Mapper and Reducer input and output formats are set up as a Key Value pair of lists of <Text, Integer> datatypes.
- Set up two jobs in the driver class:



- Job 1 reads all the selected language files in the HDFS repository. The Job 1 output to HDFS is a count for the total number of appearances for each individual letter character in the books for the selected languages.
 - Job 2 takes the output from Job 1 as an input. Character totals are recreated. The total count of ALL characters in the books is also created in the Job2 MR process. This total of all characters is used to calculate the frequency of each individual language character. The Job 2 Reducer output contains the frequency values for all characters, which is stored in a file in the newly created **CA1_Output_out** folder.
- Job1 sets up a Combiner and Partitioner Hadoop process, which are explained in the following sections.
 - Job 2 only uses a Mapper and a Reducer process.
 - A job sequence is set up so that Job 1 must complete first, before Job 2 is permitted to commence.

3.2.2 Job 1: Mapper/Combiner Process

Key steps in the first (Job1) Mapper process are;

- The function `AverageLetterMapper` is set up to process HDFS input line by line from the specific language files.
- Each line is split into individual characters.
- Non-alpha characters are excluded.
- A regex expression is used to ensure local language characters are included, along with English language characters .
- Counts of characters per line are stored in a Java HashMap and then extracted to write out individual `<Char, Count>` Key:Value pairs from the Mapper.
- A Combiner class is set up in the Driver class using the same code as the Mapper (`job1.setCombinerClass(AverageLetterReducer.class);`). This acts as a mini-Reducer running locally on the Mapper output and generates the same output format as the Mapper. It is used here to improve Job1 performance.

..

3.2.3 Job 1: Partitioner Process

Key steps in the Partitioner process are;

- A Partitioner class is set up in the Driver code (`job1.setPartitionerClass(AverageLetterPartitioner.class);`).
- The number of Reducer tasks for Job 1 is explicitly set to '2' (`job1.setNumReduceTasks(2);`).

- The choice of Reducer to be selected is determined in the Partitioner process. If a character is a vowel the Reducer processing is directed to Reducer 1. Non vowels are directed to Reducer 2.
- The reason for this partition logic is that there are more occurrences of vowels than non-vowels in a language text. However, the number of vowels is much less (5), therefore this split is used in this MR process distribute effort across the Reducers.

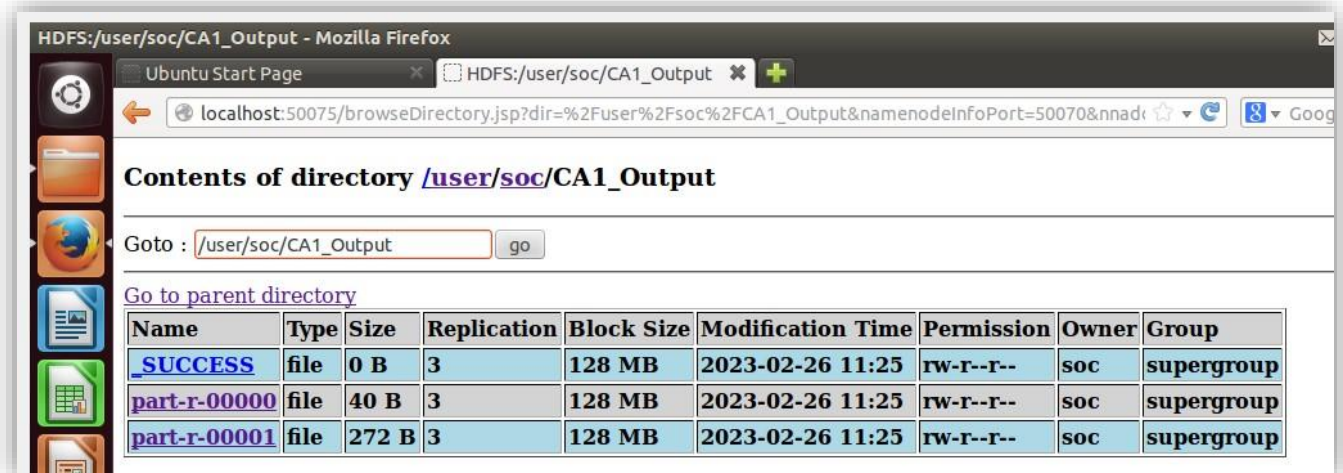
3.2.4 Job 1: Reducer Process

Key steps in the first (Job1) Reducer processes are;

- The function `AverageLetterReducer` is set to read the inputs from the Mapper process, after Hadoop sorting and grouping has taken place.
- Key values, which are individual language characters, are ingested by the Job1 Reducer processes with a list of integer counts for character occurrences in each line of the input.
- The Reducers sum all the character integer counts for each character (Key).
- The output of `<character, Total Count>` is written to HDFS. The two Reducer jobs write to separate files in the **CA_Output** directory, which was set in the command line to execute the assignment jar file.

3.2.5 Job 1: MR Process Outputs

Job 1 output is stored under in the **CA_Output** directory, created based on the configuration set up in the Driver class.



Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
_SUCCESS	file	0 B	3	128 MB	2023-02-26 11:25	rw-r--r--	soc	supergroup
part-r-00000	file	40 B	3	128 MB	2023-02-26 11:25	rw-r--r--	soc	supergroup
part-r-00001	file	272 B	3	128 MB	2023-02-26 11:25	rw-r--r--	soc	supergroup

Figure – CA1_Output directory – Job 1 Output

The '**part-r-0000**' file contains the total count of occurrences for vowel characters read in the language file inputs.



Figure – CA1_Output directory – Job 1 Output – Reducer 1 – Vowels

The '**part-r-0001**' file contains the total count of occurrences for non-vowel characters read in the language file inputs (sample represented below).

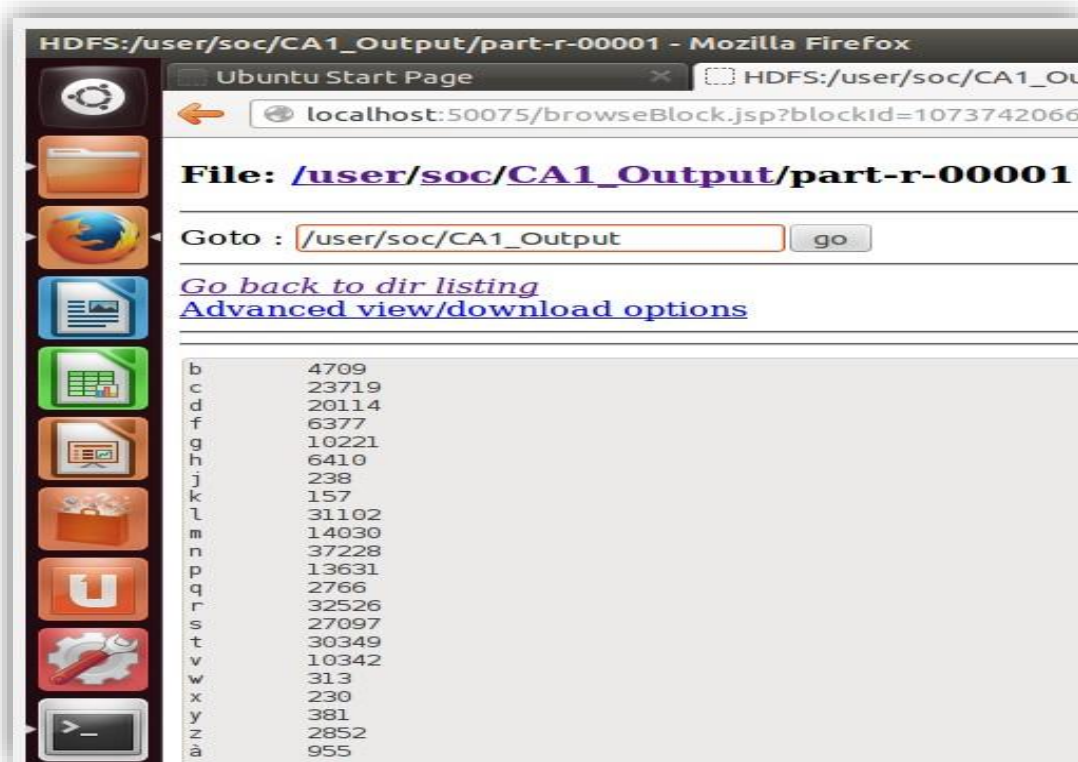


Figure – CA1_Output directory – Job 1 Output – Reducer 2 – Non-Vowels



3.2.6 Job 2: Mapper Process

Key steps in the second (Job2) Mapper process are;

- The `FreqDistribMapper` class reads the output from Job 1, which is stored in the **CA1_Output** directory in HDFS.
- The input for this Mapper process is much smaller than in Job 1. The Key values in the **CA1_Output** directory are unique, so the Mapper effectively repeats the output process of Job1.
- However, a static key of '**Total_Chars**' is also written out by the Mapper as each individual character is re-output with the number of total occurrences. The use of this static key means that it generates a master total output that grows as each separate character is read/output.
- In addition, a '**Local_Char**' static key is set up to group all non-English language characters. This is done to collate characters such as *ä*, *ö*, and *ü* (for example) into a single Key category. The purpose of this grouping is to simplify the end result comparison graphs in the Python output for this project.

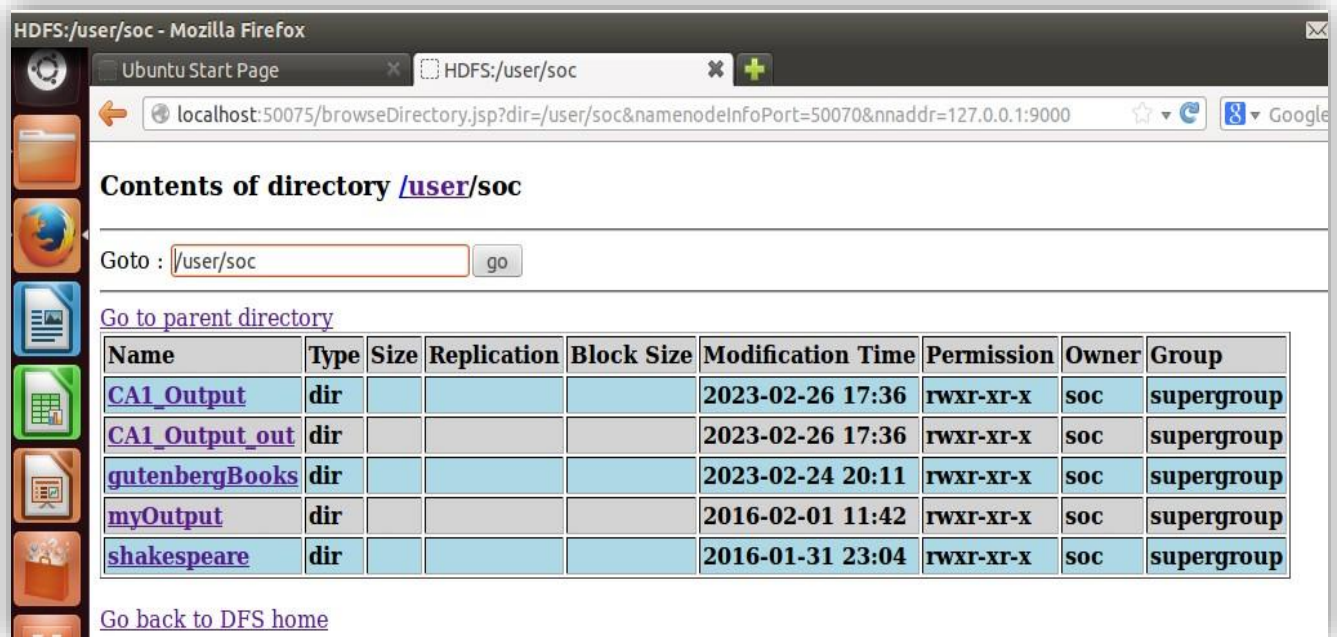
3.2.7 Job 2: Reducer Process

Key steps in the second (Job2) Reducer process are;

- The `reduce` function in the `FreqDistribReducer` class performs a simple summation of the values in the iterable Integer list for each individual key. There would be no real sorting or grouping required from Hadoop.
- This summation includes the value for language specific characters such as *ä*, *ö*, and *ü* (for example), which are grouped under a single Key value named '**Local_Char**'.
- A static Key called '**Total_Chars**' is fed into the Reducer and this summation of the Integer list for this Key will generate the value representing the **total** number of characters in the language book file inputs.
- The `reduce` function does not call the `context.write` routine, as configured by the Driver class. Instead, it incrementally fills a Java HashMap with each Key/Value pair, including the value for the count of all characters.
- The Job 2 Reducer has a `cleanup` routine that executes after main processing is complete. This cycles through the Java HashMap created by the `reduce` function and generates a distribution value for each character, based on a division of character totals by total numbers of characters.
- The Key value for final input is concatenated with the name of the language, which has been parsed from the command line input into a Job 2 configuration custom variable called '*language.text*'.
- The `cleanup` routine output is written to a HDFS directory called **CA1_output_out**.

3.2.8 Job 2: MR Process Outputs

Job 2 output is stored under in the **CA_Output_out** directory, created based on the configuration set up in the Driver class.

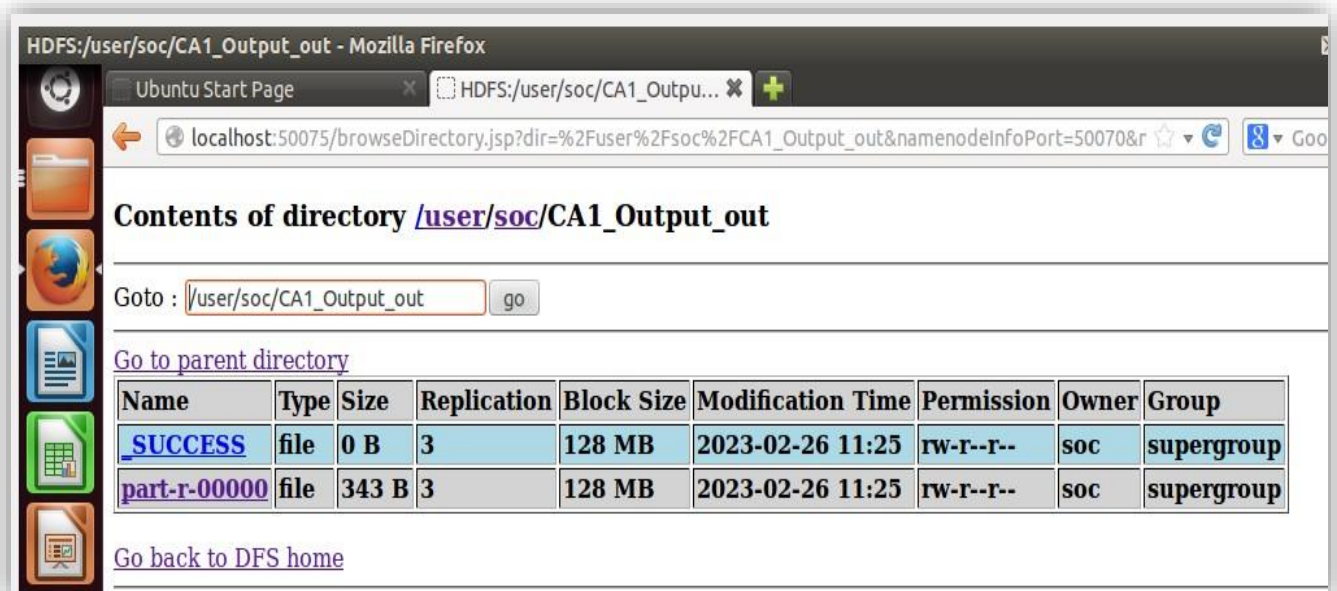


Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
CA1_Output	dir				2023-02-26 17:36	rwXr-Xr-X	soc	supergroup
CA1_Output_out	dir				2023-02-26 17:36	rwXr-Xr-X	soc	supergroup
gutenbergBooks	dir				2023-02-24 20:11	rwXr-Xr-X	soc	supergroup
myOutput	dir				2016-02-01 11:42	rwXr-Xr-X	soc	supergroup
shakespeare	dir				2016-01-31 23:04	rwXr-Xr-X	soc	supergroup

[Go back to DFS home](#)

Figure – CA1_Output_out directory – Job 2 Output

Job 2 output is one file.



Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
_SUCCESS	file	0 B	3	128 MB	2023-02-26 11:25	rw-r--r--	soc	supergroup
part-r-00000	file	343 B	3	128 MB	2023-02-26 11:25	rw-r--r--	soc	supergroup

[Go back to DFS home](#)

Figure – CA1_Output_out directory – Job 2 Output File

The '**part-r-0000**' file contains the final breakdown of average letter frequencies and the description of the language (German example represented below).

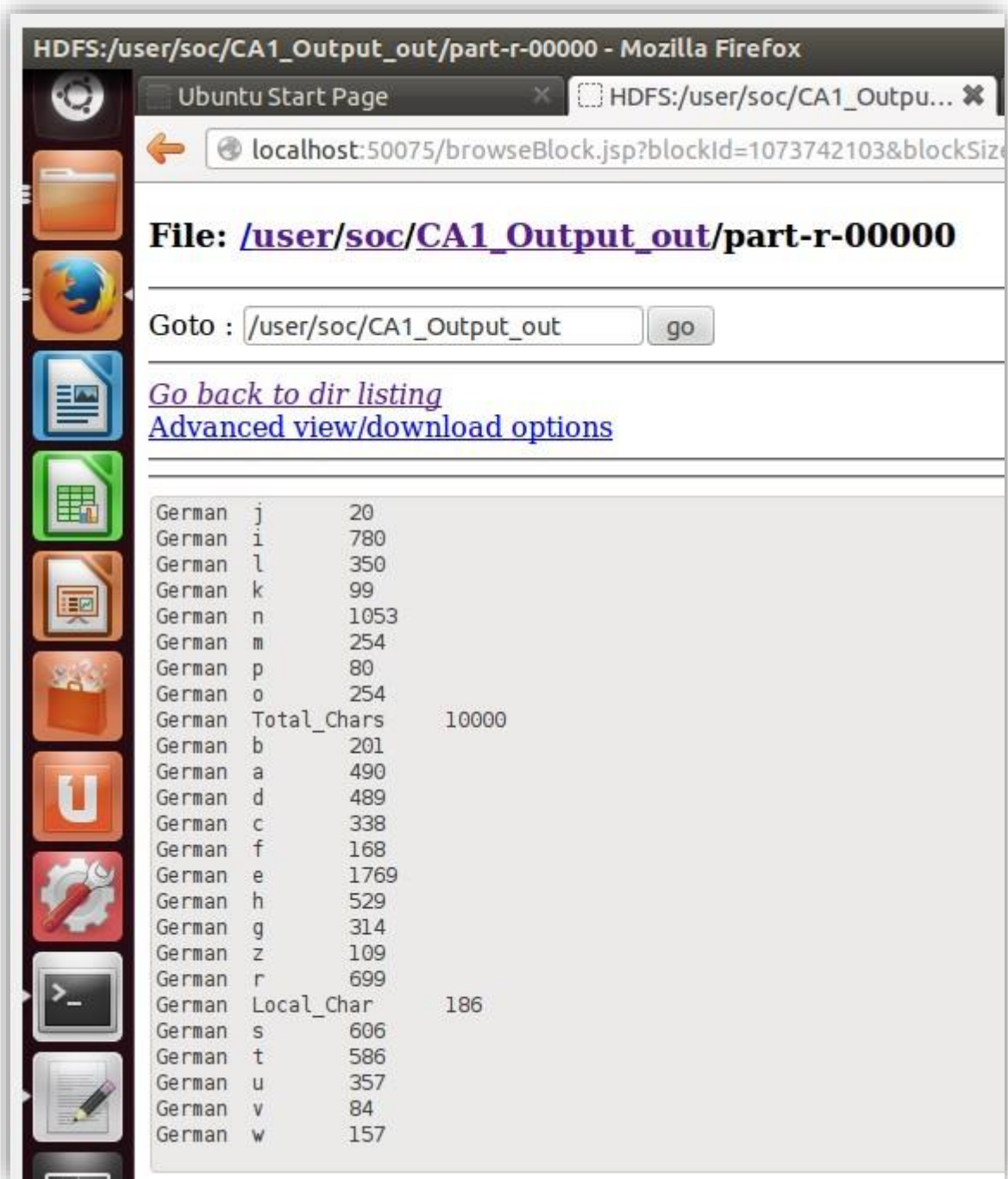


Figure – CA1_Output_out directory – Job 2 Final Output

To maintain the Integer data type for the Job 2 Reducer output, the average letter frequency is multiplied by 10000, The display value is adjusted to a > 1 decimal value during the rendering of the Python bar charts.

3.3 Key Design Decisions and Assumptions

In the design and implementation of PBD CA1 the following were the key design assumptions and decisions;

- Only one jar file is used for the Map-Reduce processing, and this is potentially able to manage multiple files across multiple languages. No hard coding of language names or named file subdirectories in HDFS is used in the Java code. The Java project would not need to be re-compiled to accommodate additional languages (limited to three for this assignment).
- Command line text and file name suffixes are used to alter the selected inputs to the Map-Reduce process and effect the content of the end result formats. All language file names must be in the format:
<bookname_language>, as seen in Section 2.3 of this report. For example, the file *IL_Diavolo_Italian* is one of the input books to the MR process in this assignment.
- Globbs with a wildcard based on the filename language will pick up all language files for the MR process. In the assignment there are two book files per language but there is no limitation to adding more (and no need to amend source code).
- Two chained jobs are required. One to calculate the total occurrences of each character in the input text, and the second to use the output of the first job to calculate the total number of characters and hence the average letter frequencies (and output).
- The Python code is written in a Jupyter Notebook and also applies the same principle of having no hard coded references to languages. No file name conventions are required. The Python code reads any number of EXCEL files in the local directory of the notebook and calls a graph function. For simplicity, only three graphs are specifically rendered for this assignment, but the code could easily accommodate more language inputs.



- The output from Job2 is downloaded from the HDFS data store after each individual language MR process. The text file is transferred from the VM onto the host machine. The inbuilt EXCEL data ingestion wizard takes the text file and is used to create a three-column output of: ***language – character – frequency count.***
- No data manipulation takes place during the generation of the XL file, and no headers are added.
- The Python graph function takes two parameters; the dataframe with the language letter frequency information loaded from XL; and the name of the language, which is read directly from the dataframe.

4 Map-Reduce: Java Code

4.1 Driver

A description of the operations in this Java code class are given in Section 3.2.1.

```
AverageLetterFrequency.java AverageLetterMapper.java AverageLetterPartitioner.java AverageLetterReducer.java FreqDistribMapper.java
1
2 import org.apache.hadoop.conf.Configuration;
13
14
15 // The driver process for the average character frequency MapReduce process
16 // A single MR process is created and input files on HDFS are read based on
17 // the command line (glob) input and the naming convention of the book files
18 // stored on HDFS.
19
20
21 // The book files are all stored with a naming convention of <bookname>_<language>
22 // For example 'TheFox_English'.
23 public class AverageLetterFrequency {
24
25     public static void main(String[] args) throws Exception {
26
27         // check command line is correctly - two parameters
28         if (args.length != 2) {
29             System.err.println("Usage: AverageLetterFrequency <input path> <output path>");
30             System.exit(-1); }
31
32         // Set Up language from input command line glob
33         String cmdInputArg0 = args[0];
34         int iUndrScrIndex = cmdInputArg0.indexOf('_');
35         String strLang = cmdInputArg0.substring(iUndrScrIndex+1).trim();
36         strLang = strLang + '\t';
37
38         // Check Parse of Input globs for language description
39         System.out.println("\n\nInput glob... args[0] - language substring is: " + strLang + "\n\n");
40
41         Configuration conf = new Configuration();
42
43         // Set up Job1
44         Job job1 = Job.getInstance(conf, "AverageLetterFrequency");
45
46         job1.setJarByClass(AverageLetterFrequency.class);
47
48
49         // Set Up Mapper and Combiner
50         job1.setMapperClass(AverageLetterMapper.class);
51         job1.setCombinerClass(AverageLetterReducer.class);
52         job1.setMapOutputKeyClass(Text.class);
53         job1.setMapOutputValueClass(IntWritable.class);
54
```



```
AverageLetterFrequency.java AverageLetterMapper.java AverageLetterPartitioner.java AverageLetterReducer.java FreqDistribMapper.java
55
56 // Set Up Partitioner
57 job1.setPartitionerClass(AverageLetterPartitioner.class);
58
59 // Set Up Reducer
60 job1.setReducerClass(AverageLetterReducer.class);
61
62 // Set number of reducer tasks
63 job1.setNumReduceTasks(2);
64
65
66 // Input and Output format for data
67 job1.setOutputKeyClass(Text.class);
68 job1.setOutputValueClass(IntWritable.class);
69
70 // Check Input globs
71 System.out.println("\n\nInput glob... args[0]: " + args[0] + "\n\n");
72 FileInputFormat.addInputPath(job1, new Path(args[0]));
73 FileOutputFormat.setOutputPath(job1, new Path(args[1]));
74
75
76
77 ControlledJob cJob1 = new ControlledJob(conf);
78 cJob1.setJob(job1);
79
80
81 // Set up Job2
82 Job job2 = Job.getInstance(conf, "AverageLetterFrequency Two");
83
84 job2.setJarByClass(AverageLetterFrequency.class);
85
86 // Set Up Mapper - no need for Combiner; data volumes are much lower
87 job2.setMapperClass(FreqDistribMapper.class);
88 job2.setReducerClass(FreqDistribReducer.class);
89
90 // Input and Output format for data
91 job2.setOutputKeyClass(Text.class);
92 job2.setOutputValueClass(IntWritable.class);
93
94 // Set a custom configuration property that will track the language
95 // being analyzed for character frequency distribution in this
96 // Map Reduce process
97 job2.getConfiguration().set("language.text", strLang);
98
```



```
AverageLetterFrequency.java  AverageLetterMapper.java  AverageLetterPartitioner.java  AverageLetterReducer.java  FreqDistribMapper.java
97     job2.getConfiguration().set("language.text", strLang);
98
99     ChainMapper.addMapper(job2, FreqDistribMapper.class, LongWritable.class, Text.class, Text.class, IntWritable.class);
100     ChainReducer.setReducer(job2, FreqDistribReducer.class, Text.class, IntWritable.class, Text.class, IntWritable.class);
101
102
103     // Job2 reads the files produced on HDFS by Job 1
104     FileInputFormat.addInputPath(job2, new Path(args[1]));
105     // Job2 output is marked with a suffix
106     FileOutputFormat.setOutputPath(job2, new Path(args[1] + "_out"));
107
108     ControlledJob cJob2 = new ControlledJob(conf);
109     cJob2.setJob(job2);
110
111
112     // Set up Job sequence execution
113     // The purpose of the chaining is that Job1 must first generate a count of each character
114     // read from the language book files.
115     //
116     // That data must be in place before Job2 is invoked to generate a count of ALL characters
117     // read from the language book files and then calculate the average character distribution
118     Job[] jobs = new Job[]{job1, job2};
119     for (Job job : jobs) {
120         boolean success = job.waitForCompletion(true);
121
122         // Exit process if one of the jobs fails
123         if (!success) {
124             System.exit(1);
125         }
126     }
127
128     // Add language description to the final output from the
129     // Reducer job after it has completed.
130
131     System.out.println("\nPreparing to end Driver process...\n");
132
133     // No jobs have failed - exit process successfully
134     System.exit(0);
135
136
137 }
138
139 }
```

4.2 Job 1

4.2.1 The Mapper

A description of the operations in this Java code class are given in Section 3.2.2.

```

AverageLetterFrequency.java | AverageLetterMapper.java | AverageLetterPartitioner.java | AverageLetterReducer.java | FreqDistribMapper.java
1 import java.io.IOException;
10
11 public class AverageLetterMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
12     private static final Log LOG = LogFactory.getLog(AverageLetterMapper.class);
13
14     static enum LangCharCount{LANG_CHAR_COUNT};
15
16     public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
17
18         String s = value.toString();
19         String strSplit = "";
20         String sRegex = "^\\p{L}";
21         HashMap<Character, Integer> charCountMap = new HashMap<Character, Integer>();
22
23         System.out.println("In CA1:AverageLetterFrequency Mapper now!");
24
25         // Mapper reads HDFS input line by line
26         for (String sInput : s.split(strSplit)) { // Split input into characters
27             // Filter out grammar symbols but allow local language characters such as - ä, ö, ü (for example)
28             if (sInput.length() > 0 && sInput.matches(sRegex)) {
29
30                 // Convert to lower case - analysis will not differentiate based on case
31                 String strLCaseWord = sInput.toLowerCase();
32
33                 char ch = strLCaseWord.charAt(0);
34
35                 // Check that character is valid alpha character - Log output for Hadoop dashboard
36                 LOG.info("Pre-screening for Character.isLetter() Check.: " + ch);
37
38                 if (Character.isLetter(ch)) { // Exclude non-alphabet characters
39
40                     // If the character is already in the map, increment its count...
41                     if (charCountMap.containsKey(ch)) {
42                         charCountMap.put(ch, charCountMap.get(ch) + 1);
43                     }
44                     // ...otherwise, add the character to the map with a count of 1
45                     else {
46                         charCountMap.put(ch, 1);
47                     }
48                 }
49             }
50         }
51     }
52
53     // Loop through the count of characters in the input read by the Map process
54     for (Character ch : charCountMap.keySet()) {
55
56         // Convert character to string for Mapper Output
57         String sCharInWord = String.valueOf(ch);
58         int ichrCnt = charCountMap.get(ch);
59
60         // Set up data for logging output to first Mapper job
61         LOG.info("Mapper output key: " + sCharInWord);
62         LOG.info("Mapper output values: " + ichrCnt);
63
64         // Output from Mapper is the count of each character in the word
65         context.write(new Text(sCharInWord), new IntWritable(ichrCnt));
66     }
67 }
68
69 }
70
71 }

```


4.2.2 The Partitioner

A description of the operations in this Java code class are given in Section 3.2.3,

```

AverageLetterFrequency.java  AverageLetterMapper.java  AverageLetterPartitioner.java  AverageLetterReducer.java  FreqDistribMapper.java
1
2 import org.apache.hadoop.io.IntWritable;
3
4
5
6
7
8
9 public class AverageLetterPartitioner extends Partitioner<Text, IntWritable>{
10     private static final Log LOG = LogFactory.getLog(AverageLetterPartitioner.class);
11     String partitionKey;
12
13 @Override
14     // Input is partitioned into different Reducers based on a check as to
15     // whether the character is a vowel or not.
16     //
17     // Vowels are more commonly used but are a small subset of a Latin alphabet
18     // The Reducers are split between processing vowels and non-vowels
19     public int getPartition(Text key, IntWritable value, int numPartitions){
20
21         // Set up vowel string
22         String vowelCharsToMatch = "aeiou";
23
24         if(numPartitions == 2){
25             String partitionKey = key.toString();
26
27             // Compare input key character against vowel string
28             if (vowelCharsToMatch.contains(String.valueOf(partitionKey.charAt(0)))){
29                 LOG.info("numPars: 2 - Partition 0"); // Send Output to Log to indicate Partitioner is working
30                 return 0; // Partition 0 will be used for vowels
31             }
32             else {
33                 LOG.info("numPars: 2 - Partition 1"); // Send Output to Log to indicate Partitioner is working
34                 return 1; // Partition 1 will be used for non-vowels
35             }
36         } else if (numPartitions == 1) { // Default
37             return 0;
38         }
39         else {
40             System.err.println("AverageLetterPartitioner can only handle 1 or 2 partitions");
41             return 0;
42         }
43     }
44 }
45
```

4.2.3 The Reducer

A description of the operations in this Java code class are given in Section 3.2.4.

```

AverageLetterFrequency.java  AverageLetterMapper.java  AverageLetterPartitioner.java  AverageLetterReducer.java
1
2+ import java.io.IOException;
11
12
13 public class AverageLetterReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
14     private static final Log LOG = LogFactory.getLog(AverageLetterReducer.class);
15
16     public void reduce(Text key, Iterable<IntWritable> values, Context context)
17         throws IOException, InterruptedException {
18
19         List<IntWritable> valueList = new ArrayList<IntWritable>();
20
21
22         // This code block takes the Iterable list input and reads into
23         // a new list
24         for (IntWritable val : values) {
25             valueList.add(new IntWritable(val.get()));
26         }
27
28         // The purpose of this code is to read the key/value inputs into
29         // new variables and use the LOG function to output to the Hadoop
30         // dashboard to show that the Mapper input to the first Reducer
31         // process is correct.
32         StringBuilder sb = new StringBuilder();
33         for (IntWritable v : valueList) {
34             sb.append(v.get());
35             sb.append(", ");
36         }
37         if (sb.length() > 0) {
38             sb.delete(sb.length() - 2, sb.length());
39         }
40
41         // Output to log file in Hadoop dashboard
42         String valuesAsString = sb.toString();
43         LOG.info("Reducer input key: " + key);
44         LOG.info("Reducer input values: " + valuesAsString);
45
46
47         // Initialize counter variable
48         int iChrCount = 0;
49
50
51         // Count the total number of characters in the book
52         // The use of the Iterable list in the LOG output requires
53         // the use of the 'new' list to provide Reducer function
54         // output.
55         //
56         // (Once the Mapper Iterable list is processed for the LOG
57         // it cannot be reset).
58         for (IntWritable vi : valueList) {
59
60             // Read the Mapper/Combiner output and increment counter
61             iChrCount += vi.get();
62         }
63
64
65         // Write Reducer Output - ignore very infrequent characters
66         // If, after reading entire books in the language, the character
67         // count is less than 5 then the character can be considered spurious
68         if (iChrCount > 5) {
69             context.write(key, new IntWritable(iChrCount));
70         }
71     }
72
73 }

```

4.3 Job 2

4.3.1 The Mapper

A description of the operations in this Java code class are given in Section 3.2.6.

```

AverageLetterFrequency.java  AverageLetterMapper.java  AverageLetterPartitioner.java  AverageLetterReducer.java  FreqDistribMapper.java
1+ import java.io.IOException;
9
10 // This is the second (chained) Mapper process
11 // This reads the output from the first MapReduce Job, which is a file with
12 // a list of each character found in the language book files on HDFS, along
13 // with the count of the occurrence of those characters.
14 public class FreqDistribMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
15     private static final Log LOG = LogFactory.getLog(FreqDistribMapper.class);
16
17     // The purpose of this Map process is to create a Key/Value pair that uses a static key text to
18     // record a value that represents the total count of ALL characters and which will be used for
19     // average frequency distribution calculations in the Job 2 Reducer function.
20     public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
21
22         String s = value.toString();
23         String strSplit = "\\W+";
24         String strTotalLabel = "Total_Chars"; // Static Key Value set to record ALL characters
25
26         // Each line in the Job 1 output is read.
27         // Each line represents a character and character count
28         String[] parts = s.split(strSplit);
29
30         // The first line input represents the individual language character
31         String sCharCapture = parts[0].trim();
32         // The second line represents the count of the character
33         int iCharCntInt = Integer.parseInt(parts[1].trim());
34
35         // 'Local' (non-English) characters will not be read - but their count is
36         // The blank character is interpreted as one of these non-English characters
37         // (such as ä, ö, ü) and they are grouped under a default key value
38         if (sCharCapture.isEmpty()) {
39             sCharCapture = "Local_Char";
40             LOG.info("Checking for empty key-value: " + sCharCapture + "-" + iCharCntInt);
41         }
42
43         // This line rewrites the character and count to preserve this data for the Job 2 Reducer
44         context.write(new Text(sCharCapture), new IntWritable(iCharCntInt));
45         // This line writes out a static key value to ensure the Reducer is fed a count of
46         // total characters read from the language book files.
47         context.write(new Text(strTotalLabel), new IntWritable(iCharCntInt));
48     }
49 }
50 }
```

4.3.2 The Reducer

A description of the operations in this Java code class are given in Section 3.2.7,

```

AverageLetterFrequency.java  AverageLetterMapper.java  AverageLetterPartitioner.java  FreqDistribReducer.java  A
1
2  import java.io.IOException;
11
12
13 // This is the second (chained) Reducer job
14 // The Job 2 Mapper function uses a static key value to allow for a count of all
15 // characters. This allows the Job 2 Reducer function here to calculate the
16 // average frequent for each character and generate the final output
17 public class FreqDistribReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
18     private static final Log LOG = LogFactory.getLog(FreqDistribReducer.class);
19     private HashMap<Text, Integer> totals = new HashMap<Text, Integer>();
20     private String strTotalLabel = "Total_Chars";
21     private Text languageText;
22
23     // The Driver function sets up a language text output for the final HSFS Reducer
24     // output based on the command line (glob) used to run the MR jar file
25     protected void setup(Context context) throws IOException, InterruptedException {
26         // Read the language description used as the parameter to run the jar file
27         // from the command line.
28         languageText = new Text(context.getConfiguration().get("language.text"));
29     }
30
31     public void reduce(Text key, Iterable<IntWritable> values, Context context)
32         throws IOException, InterruptedException {
33
34         // This is similar to the reducer in Job 1 as a total is counted
35         // of all occurrences of the Key values
36
37         // For Job2, one of the inputs reflects the count of ALL characters,
38         // which has been recorded against a static 'Total_Chars' key.
39         int sum = 0;
40
41         for (IntWritable value : values) {
42             sum += value.get();
43         }
44
45         if (totals.containsKey(key)) {
46             sum += totals.get(key);
47         }
48
49         totals.put(new Text(key.toString()), sum);
50     }
51

```



```

AverageLetterFrequency.java  AverageLetterMapper.java  AverageLetterPartitioner.java  FreqDistribReducer.java
53      // This function is run to process the end result Reducer data and perform the
54      // calculations to generate the frequency distribution of characters
55  protected void cleanup(Context context) throws IOException, InterruptedException {
56      float iTotalCharsCnt = 0;
57
58      // Create a HashMap with the Reducer dataset
59      for (Map.Entry<Text, Integer> entry : totals.entrySet()) {
60          String sChkTotal = entry.getKey().toString().trim();
61
62          // Store the value that represents the total count of ALL
63          // characters read by the MapReduce process
64          if (sChkTotal.equals(strTotalLabel)){
65              iTotalCharsCnt = entry.getValue();
66          }
67      }
68
69      // Use the count of ALL characters to loop through the count for individual
70      // characters and produce a frequency distribution value for each one and
71      // output the result as the final Reducer result to HDFS.
72      for (Map.Entry<Text, Integer> entry2 : totals.entrySet()) {
73          float iEnt = (entry2.getValue());
74          float iCalField = ((iEnt) / iTotalCharsCnt);
75          int iDistrib = (int) (iCalField * 10000);
76
77          // Output LOG data for the Hadoop dashboard logs on the second (Job2)
78          // reducer output.
79          LOG.info("Working through second for-loop - Key-Value-Calc-Int: "
80                  + entry2.getKey() + "-"
81                  + entry2.getValue() + "-"
82                  + iCalField + "-" + iDistrib);
83
84          // Datatype manipulation to format output of Reducer job
85          String sTxt1 = languageText.toString();
86          String sTxt2 = entry2.getKey().toString().trim();
87          String sTxt3 = sTxt1+sTxt2;
88          Text keyText = new Text(sTxt3);
89
90          // If character is very infrequent it can be ignored as a spurious
91          // inclusion in the language documents
92          if (iDistrib > 10) { // Appears less than 0.1%
93              //Generate final Job 2 Reducer outputs
94              context.write(keyText, new IntWritable(iDistrib));
95          }
96      }
97  }
98 }
```

5 Python: Graph Analysis of MR Outputs

5.1 Set Up Excel From HSFs Output

The final output from the Map-Reducer process is downloaded on the VM and then copied to the host machine.

The text file is loaded into Excel and the inbuilt wizard creates an *.xlsx file output with the data in three columns (no headers).

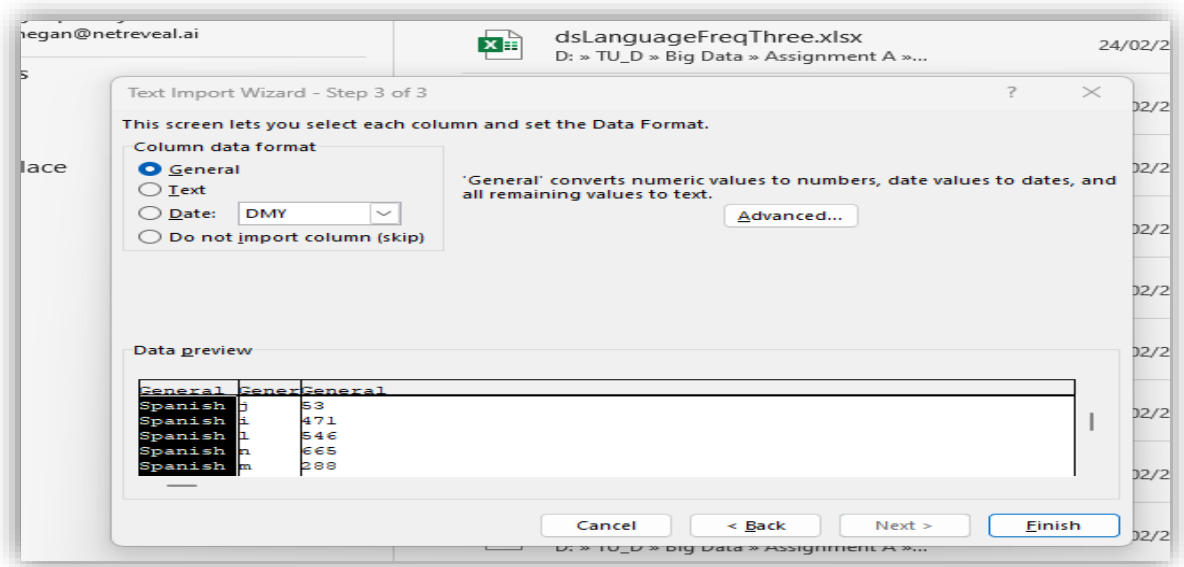


Figure – Loading HDFS Output Into Excel

All *.xlsx files are then copied to the same local directory from which the Python Jupyter Notebook written for the assignment is stored.

5.2 Python Jupyter Notebook

The Python code reads all locally stored *.xlsx files (all of which are assumed to be in the HDFS output format) and loads this language data into separate Panda dataframes.

The notebook generates a character frequency graph for each language dataframe.

Below is the Jupyter Notebook code with generated output included;

Programming for Big Data - Continuous Assessment 1 - Class Group: TU060

MSc in Computer Science - Data Science (Part Time - Yr 2)

Student Name: Ciaran Finnegan

Student No: D21124026

Student EMail: D21124026@mytudublin.ie

February 2023

This Python expects all MapReduce output files to be converted to EXCEL and stored in the local programme directory.

Language names are read from the file content (not hard coded in the file name).

The Python code is written to read any number of input language files, but for simplicity only three language graphs are created.

Table of Contents

- ▼ [1 Import Python Libraries](#)
 - ▼ [1.1 Import Libraries](#)
 - [1.1.1 Panda + Numpy](#)
 - [1.1.2 Plotly Libraries](#)
 - [1.1.3 Library for file loading using Wildcards](#)
- ▼ [2 Language Data Analysis](#)
 - [2.1 Set Up Graph Function](#)
 - ▼ [2.2 Ingest Language Files into Panda Dataframes](#)
 - [2.2.1 Read Language Data Outputs from MapReducer Processes](#)
 - ▼ [2.3 Display Language Dataframes](#)
 - [2.3.1 Display Character Frequency Dataframes](#)
 - ▼ [2.4 Display Language Graphs](#)
 - [2.4.1 Display Bar Chart Graph For Spanish](#)
 - [2.4.2 Display Bar Chart Graph For Italian](#)
 - [2.4.3 Display Bar Chart Graph For German](#)

1 Import Python Libraries

1.1 Import Libraries

Note: It is assumed that the matplotlib and plotly library has been installed on the user's machine before it is available for import.

```
> pip install plotly
```

1.1.1 Panda + Numpy

```
In [1]: # Python Libraries for data manipulation - dataframes
import pandas as pd
```

1.1.2 Plotly Libraries

```
In [2]: # For interactive graphics
import plotly.express as px
```

1.1.3 Library for file loading using Wildcards

```
In [3]: import glob
```

2 Language Data Analysis

What are the letter frequencies across languages?

2.1 Set Up Graph Function

Single code block - function is called multiple times to display a language graph.

```
In [4]: ## This is a function to generate the character distribution graph.

## ALL Language outputs from the MR processes are identical in format.
## Hence a single function can be written to generate all graphs.

def generate_CharFreqGrpah(df_lang, language): ## Language read from XL file contents directly

    # Sort the dataframe so that the x-axis is in alphabetical order
    df_lang = df_lang.sort_values('Character')

    # Bar Chart on Letter frequency
    figCharFreq = px.bar(df_lang,
                        x='Character',
                        y=df_lang['Frequency']/10000, # Adjust format of Y Axis
                        template="simple_white" # Use a clearer template
                        )

    figCharFreq.update_yaxes(# Set up the y-axis format
                            title=dict(
                                font_size=20,
                                text="Letter Frequency",
                            ),
                            tickformat='.2f',
                            showgrid=True
                            )

    figCharFreq.update_xaxes(# Set up the x-axis format
                             title=dict(
                                 font_size=20
                             ),
                             ticks="outside",
                             showgrid=True
                             )

    figCharFreq.update_layout(# Customize font and Legend orientation & position
                              font_family="Rockwell",
                              height=700,
                              title=dict(# Dynamic graph title
                                          text="Average Distribution of Letter Frequency for the " + language + " language",
                                          font_size=20
                              ),
                              margin=dict(
                                  l=10,
                                  r=60,
                                  b=10,
                                  t=150
                              )
                              )

    # Display Bar Chart
    figCharFreq.show()

    return 1
```



2.2 Ingest Language Files into Panda Dataframes

2.2.1 Read Language Data Outputs from MapReducer Processes

```
In [5]: # Set up custom colum names for dataframe
names=['Language','Character', 'Frequency']

In [6]: # Define the wildcard cpattern to search for all Excel files in Local directory
xlLanguage_files = glob.glob('*.*xlsx')

In [7]: # Set up List of Language dataframes
dfs_lanuagelist = []

In [8]: # Loop through the EXCEL files ouput from the MAPReducer processes
# and read each one into a dataframe list
for xl_file in xlLanguage_files:
    df = pd.read_excel(xl_file, header=None, names=names)
    dfs_lanuagelist.append(df)
```

2.3 Display Language Dataframes

2.3.1 Display Character Frequency Dataframes

The output from the dataframe list shows the format of the data after being read from the EXCEL file source.

The dataframe has the following structure:

- > Each line has the name of the 'Language'. The upper most value in the dataframe is read by the graph function for dynamic title display.
- > The 'Total Character' value used in the calaculations second (chained) MapReduce process in Hadoop is removed to avoid distorting the graph.
- > The dataframe is put in an English language centric order (a,b,c, etc.). This is done to make visual comparisons across langauges graphs more straightforward.
- > The 'Frequency' represents an average distribution number of the character. It is calculated by the division of the number of occurences of an individual charcter recorded in the MR process by the total number of characters read in from the language book inputs. The value has been scaled up into a larger Integer value in the MR process (to simplify the last Reduce process output) but is adjusted during graph generation.



In [9]: # Loop through Language dataframe for some basic data clean up and sorting

```
for df_language in dfs_languagelist:
    # Remove the Total Chars row from dataframes - this was used in the MapReduce processing but is not required for graph
    df_language.drop(df_language.index[(df_language['Character'] == "Total_Chars")],axis=0,inplace=True)

    # Sort the dataframe so that the x-axis is in alphabetical order
    df_language = df_language.sort_values('Character')

    # Display dataframe for each Language - break space included to separate Language files
    print(df_language)
    print("\n")
    print("\n")
```

	Language	Character	Frequency
20	Spanish	Local_Char	298
9	Spanish	a	1176
8	Spanish	b	173
11	Spanish	c	372
10	Spanish	d	507
13	Spanish	e	1322
12	Spanish	f	47
15	Spanish	g	102
14	Spanish	h	112
1	Spanish	i	471
0	Spanish	j	53
2	Spanish	l	546
4	Spanish	m	288
3	Spanish	n	665
6	Spanish	o	926
5	Spanish	p	231
18	Spanish	q	163
19	Spanish	r	636
21	Spanish	s	723
22	Spanish	t	412
23	Spanish	u	456
24	Spanish	v	118
16	Spanish	y	143
17	Spanish	z	38

	Language	Character	Frequency
18	Italian	Local_Char	101
8	Italian	a	1063
7	Italian	b	91
10	Italian	c	461
9	Italian	d	391
12	Italian	e	1153
11	Italian	f	124
14	Italian	g	198
13	Italian	h	124
0	Italian	i	1061
1	Italian	l	604
3	Italian	m	272
2	Italian	n	724
5	Italian	o	960
4	Italian	p	265
16	Italian	q	53
17	Italian	r	632
19	Italian	s	527
20	Italian	t	590
21	Italian	u	314
22	Italian	v	201
15	Italian	z	55

	Language	Character	Frequency
19	German	Local_Char	186
10	German	a	490
9	German	b	201
12	German	c	338
11	German	d	489
14	German	e	1769
13	German	f	168
16	German	g	314
15	German	h	529
1	German	i	780
0	German	j	20
3	German	k	99
2	German	l	350
5	German	m	254
4	German	n	1053
7	German	o	254
6	German	p	80
18	German	r	699
20	German	s	606
21	German	t	586
22	German	u	357
23	German	v	84
24	German	w	157
17	German	z	109

2.4 Display Language Graphs

This Python code will read any number of language file inputs.

However for the purpose of this assignment only the three following language files are displayed:

- > Spanish.
- > Italian.
- > German.

The reference to 'Local Character' reflects the decision to group non-English characters such as ä, ö, ü in a single category.

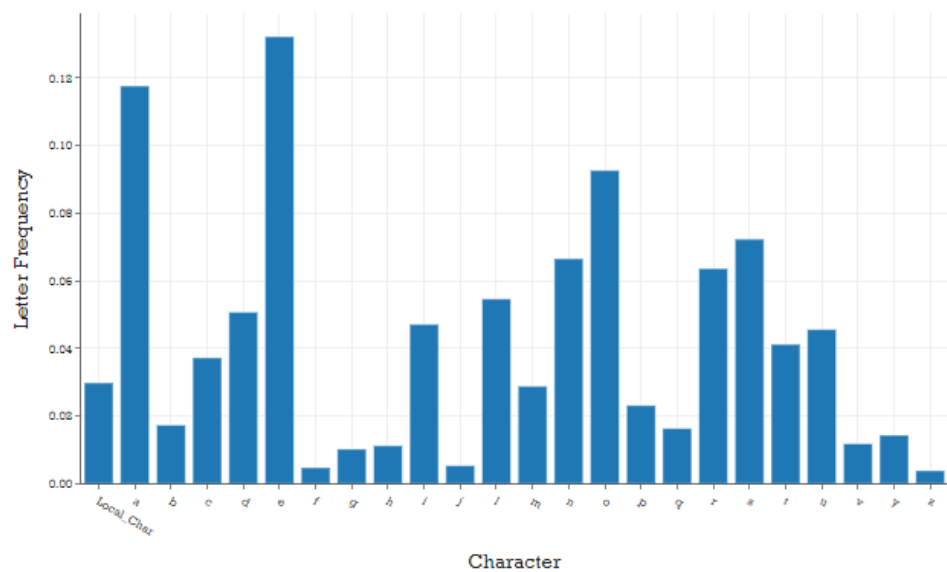
These local language characters are grouped under a single heading in the MR process and the average distribution calculated accordingly.

This decision does remove certain data granularity but makes cross language visual comparisons of the language graphs more straightforward.

2.4.1 Display Bar Chart Graph For Spanish

```
In [10]: # First Language Dataframe - Invoke Graph generation function with Language dataframe
result = generate_CharFreqGraph(dfs_languagelist[0], dfs_languagelist[0].iloc[0,0])
```

Average Distribution of Letter Frequency for the Spanish language

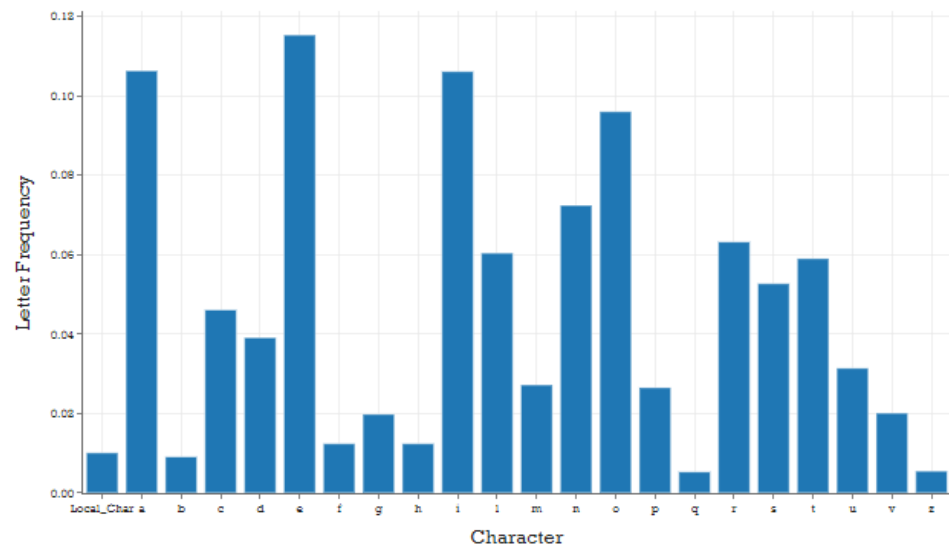


The frequency of the 'a', 'e', and 'o' characters in Spanish is not unlike English. (An English language character distribution graph is included in the brief for this assignment).

2.4.2 Display Bar Chart Graph For Italian

```
In [11]: # Second Language Dataframe - Invoke Graph generation function with language dataframe
result = generate_CharFreqGrpah(dfs_lanuagelist[1], dfs_lanuagelist[1].iloc[0,0])
```

Average Distribution of Letter Frequency for the Italian language

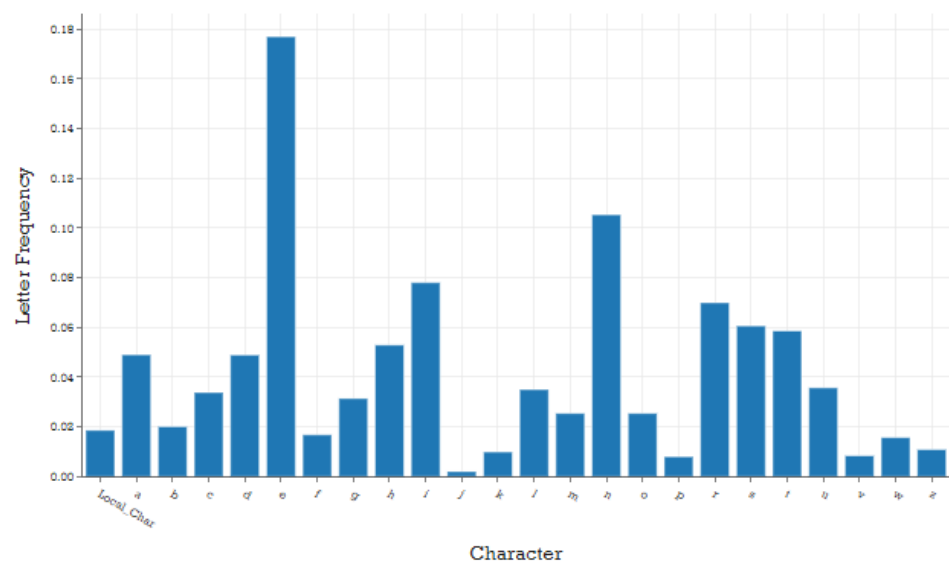


The frequency of 'i', and 'o' characters in Italian is more common than in English.

2.4.3 Display Bar Chart Graph For German

```
In [12]: # Third Language Dataframe - Invoke Graph generation function with language dataframe
result = generate_CharFreqGrpah(dfs_lanuagelist[2], dfs_lanuagelist[2].iloc[0,0])
```

Average Distribution of Letter Frequency for the German language



The frequency of 'e', and 'n' characters in German is significantly more common than in English.

<image>



5.3 Analysis Commentary

Commentary on each language graph is included in the output in Section 2.4 of the Python Jupyter Notebook.

A brief description is provided on how the distribution of characters for each language varies from English. The average letter frequency for English given in the brief for PBD Assignment 1 is used as the basis for this comparison.

Given that the assignment brief shows a letter frequency based on the Oxford English dictionary, it would be necessary to greatly increase the size of input for 'foreign' languages in our VM to give a more accurate comparison. Resource limitation for this CA meant that it was necessary to restrict the size of the Map-Reduce input.