# Course: Large Scale Data Management

## Professor: Panagiotis Liakos

## Student name: Vasileios Ilias Drouzas

## Student ID: f3352301

## 1st Assignment

In this project, we will be using the Hadoop Map-Reduce framework. In Part I, we will build a map-reduce application while in Part II we will develop our own map-reduce application.

PART I:

For this part, we will use a book from Agatha Christie, "The mystery of the Blue train". Here, we will be using the Vagrantfile that is provided, but we will do two minor changes, in order to run our own text file:

A) we will alter in the vagrant file the link to our text file and the name of the file to *train.txt*.

B) we will change the name in the java executable source file **Driver.java** to *train.txt* to make sure we will produce the correct .jar file.

The next steps are the following:

(Note: Here, not all the output is presented (only the most important aspects). The whole output is on 'output.txt').

1) Running `vagrant up` we get a successful build at 02:04 min.



2) We will execute `vagrant ssh` to get to the environment of the vagrant and then with `docker ps` we can check if everything went as normal:



We get all the expected containers and they are all set up and running, so we are ready to go!

3) Next, we will change directory and install maven:

4) Now it's time to copy our application inside a docker container and execute it!





Even though we got a (warning) exception in the DataStreamer, our execution is not interrupted and the application is finely produced. Now we can check whether our application is running in http://localhost:8088 :



Our application is up and running! We can also check the first 100 lines of the results:

```
              Bytes Written=124265
vagrant@vagrant:/vagrant/hadoop-mapreduce-examples$ docker exec namenode hdfs dfs -text /user/hdfs/output/part-r-00000 | head -100
2024-01-30 15:30:35,323 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2024-01-30 15:30:35,451 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
        3592
"'Allo--'allo--yes,     1
"'Journeys      1
"'One   1
"'Viscountess   1
"'_Soupir       1
"A      45
"A-ha!  1
"About  6
"Absurd,        1
"According      1
"Ada    1
"After, 1
"Afterwards,    1
"Against        1
"Ah     2
"Ah!    7
"Ah!"   21
"Ah!--before    1
"Ah,    16
"Ah,"   11
"Ah-ha,"        1
"Ah?"   3
"All    15
"Allow  1
"Always 1
"Always."       1
"Am     2
"Ambrose        1
"Amongst        1
"An     7
"And    93
"And,   1
"Any    3
"Anything       2
"Anyway,"       1
"Apparently--for        1
"Apple  1
"Are    10
"As     14
"Ask    2
```

```
"Breathe          1
"Bring   1
"Business         1
"Business?"       1
"But     58
"But,    6
"But,"   1
"But--but         1
"By      3
"Can     4
"Can't   1
"Care    1
"Cast    1
"Certainly        6
"Certainly,       3
"Certainly,"      1
"Certainly--but 1
"Certainly.       1
"Chubby 1
"Clever 1
"Clothes          1
"Clothes?         1
"Clothilde,       1
"Come    3
"Come,   2
"Comte   1
"Console          1
"Courage,         2
"Cousin,          1
"Dad!"   1
"Dad,    1
"Damn    1
"Dancing          1
"Darling!"        1
"Darling,"        2
"Darned 1
"Day-dreaming,    1
"Dead!" 1
"Death, 1
"Dereek!"         1
"Dereek--I        1
"Dereek--you      1
"Derek   1
"Did     6
"Divorce!"        1
"Divorce."        1
"Do      18
"Does    3
"Dollars?"        2
```

Finally, we can print the *train.txt* file with the following command:

```
docker exec namenode hdfs dfs -text /user/hdfs/input/train.txt
```

A random peak at the data:

kind to attract the men. And, besides, you're getting on. How old are
you now?"

"Thirty-three," Katherine told her.

"Well," remarked Miss Viner doubtfully, "that's not so very bad. You've
lost your first freshness, of course."

"I'm afraid so," said Katherine, much entertained.

"But you're a very nice girl," said Miss Viner kindly. "And I'm sure
there's many a man might do worse than take you for a wife instead of
one of these flibbertigibbets running about nowadays showing more of
their legs than the Creator ever intended them to. Good-bye, my dear,
and I hope you'll enjoy yourself, but things are seldom what they seem
in this life."

Heartened by these prophecies, Katherine took her departure. Half the
village came to see her off at the station, including the little maid
of all work, Alice, who brought a stiff wired nosegay and cried openly.

"There ain't a many like her," sobbed Alice when the train had finally
departed. "I'm sure when Charlie went back on me with that girl from
the Dairy, nobody could have been kinder than Miss Grey was, and though
particular about the brasses and the dust, she was always one to notice
when you'd give a thing an extra rub. Cut myself in little pieces for
her, I would, any day. A real lady, that's what I call her."

Such was Katherine's departure from St. Mary Mead.


                    8. Lady Tamplin Writes a Letter


"Well," said Lady Tamplin, "well."

She laid down the continental _Daily Mail_ and stared out across the
blue waters of the Mediterranean. A branch of golden mimosa, hanging
just above her head, made an effective frame for a very charming
picture. A golden-haired, blue-eyed lady in a very becoming negligee.
That the golden hair owed something to art, as did the pink-and-white
complexion, was undeniable, but the blue of the eyes was Nature's gift,
and at forty-four Lady Tamplin could still rank as a beauty.

Charming as she looked, Lady Tamplin was, for once, not thinking of
herself. That is to say, she was not thinking of her appearance. She
was intent on graver matters.

… and in the end of the file:

```
($1 to $5,000) are particularly important to maintaining tax exempt
status with the IRS.

The Foundation is committed to complying with the laws regulating
charities and charitable donations in all 50 states of the United
States. Compliance requirements are not uniform and it takes a
considerable effort, much paperwork and many fees to meet and keep up
with these requirements. We do not solicit donations in locations
where we have not received written confirmation of compliance. To SEND
DONATIONS or determine the status of compliance for any particular state
visit www.gutenberg.org/donate.

While we cannot and do not solicit contributions from states where we
have not met the solicitation requirements, we know of no prohibition
against accepting unsolicited donations from donors in such states who
approach us with offers to donate.

International donations are gratefully accepted, but we cannot make
any statements concerning tax treatment of donations received from
outside the United States. U.S. laws alone swamp our small staff.

Please check the Project Gutenberg web pages for current donation
methods and addresses. Donations are accepted in a number of other
ways including checks, online payments and credit card donations. To
donate, please visit: www.gutenberg.org/donate.

Section 5. General Information About Project Gutenberg™ electronic works

Professor Michael S. Hart was the originator of the Project
Gutenberg™ concept of a library of electronic works that could be
freely shared with anyone. For forty years, he produced and
distributed Project Gutenberg™ eBooks with only a loose network of
volunteer support.

Project Gutenberg™ eBooks are often created from several printed
editions, all of which are confirmed as not protected by copyright in
the U.S. unless a copyright notice is included. Thus, we do not
necessarily keep eBooks in compliance with any particular paper
edition.

Most people start at our website which has the main PG search
facility: www.gutenberg.org.

This website includes information about Project Gutenberg™,
including how to make donations to the Project Gutenberg Literary
Archive Foundation, how to help produce our new eBooks, and how to
subscribe to our email newsletter to hear about new eBooks.
```

PART II

In this part, we implemented two java classes to implement the application. You can see them in the corresponding .java files, which include the code and some explanatory comments. The results are in the 'output_partb.txt'.

Note that I have changed the name of the 'wordcount' to 'danceability', and in order to work I changed the name of the class, the package, the .java file and the name of the class in the pom.xml file. So for your convenience, I will be attaching .xml file too.



Figure 1: Types of dancing[1]

---