# WORDCOUNT PROBLEM ON HADOOP REPORT

## 1. Set up:



- **Step 1:** Create a file named "WordCount.java"

**Note: the nvim command will both create a file and open a code editor neovim**

- **Step 2:** Copy code from the tutorial into the file
- **Step 3:** Compile the java file into jar file (wc.jar)



- **Step 4:** Create a input folder for our input file
- **Step 5:** Create a input file (input3.txt)

**The content of the input file (input3.txt)**
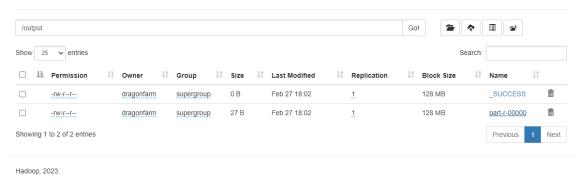
## 2. Run the application:





- **Step 6**: Put the input3.txt file into the input folder
- **Step 7:** Go into the hadoop folder and run the application

## 3. Checking the result:



**Website should have both input folder and output folder**

## Browse Directory

| | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name | |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | -rw-r--r-- | dragonfarm | supergroup | 0 B | Feb 27 18:02 | 1 | 128 MB | _SUCCESS | 🗑 |
| ☐ | -rw-r--r-- | dragonfarm | supergroup | 27 B | Feb 27 18:02 | 1 | 128 MB | part-r-00000 | 🗑 |

/output — Go!

Show 25 entries — Search:

Showing 1 to 2 of 2 entries — Previous 1 Next

Hadoop, 2023.

**Output folder content**

- **Step 8:** Go into your terminal.
- **Step 9:** Enter command **"hdfs dfs -cat /output/part-r-00000".**
- **Step 10:** Now it should show the result like the image below.

```
                    Bytes Written=27
tml_21127642@  > (dragonfarm ~/hadoop/hadoop-3.3.6) hdfs dfs -cat /output/part-r-00000
Goodbye  1
Hadoop   2
Hello    1
tml_21127642@  > (dragonfarm ~/hadoop/hadoop-3.3.6) hdfs dfs -cat /input/input3.txt
Hello Hadoop Goodbye Hadoop
tml_21127642@  > (dragonfarm ~/hadoop/hadoop-3.3.6)
```

**The result**

And there you have it, you've run the Example WordCount problem on hadoop

## 4. Source:

- https://hadoop.apache.org/docs/current/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html#Example%3A_WordCount_v1.0 – ("Apache Hadoop 3.3.6 – MapReduce Tutorial")