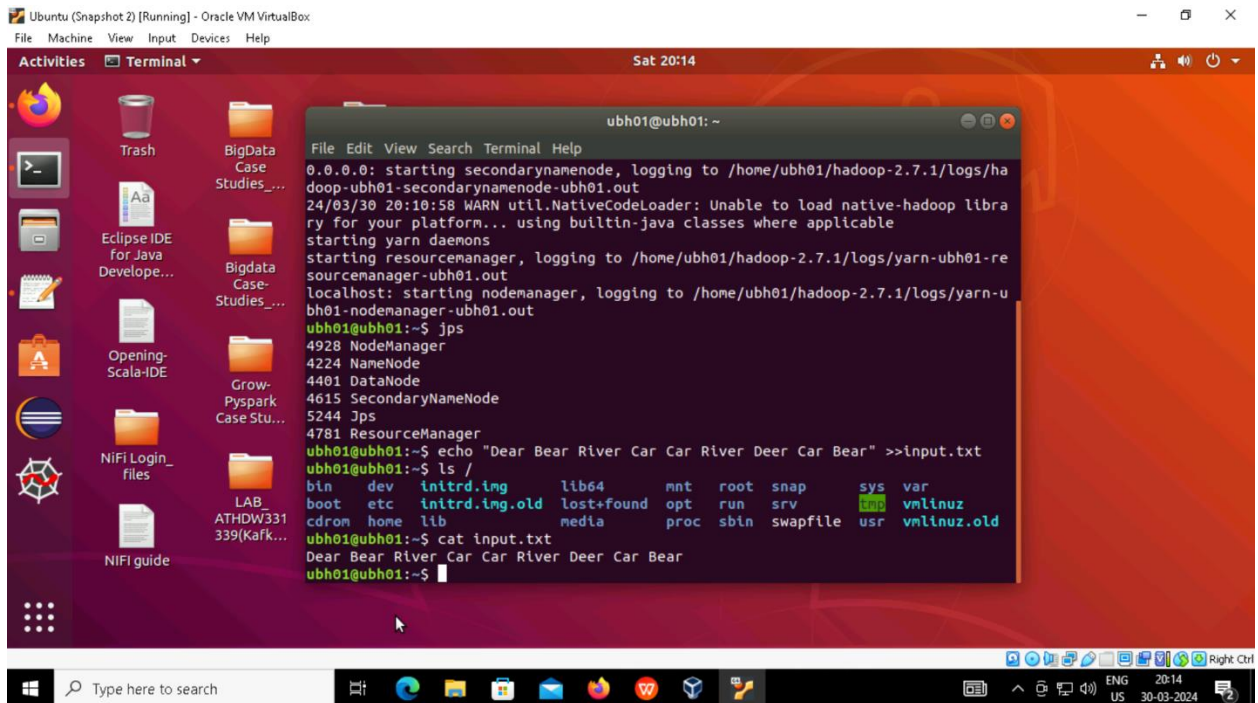# BIG DATA

# ASSIGNMENT 4

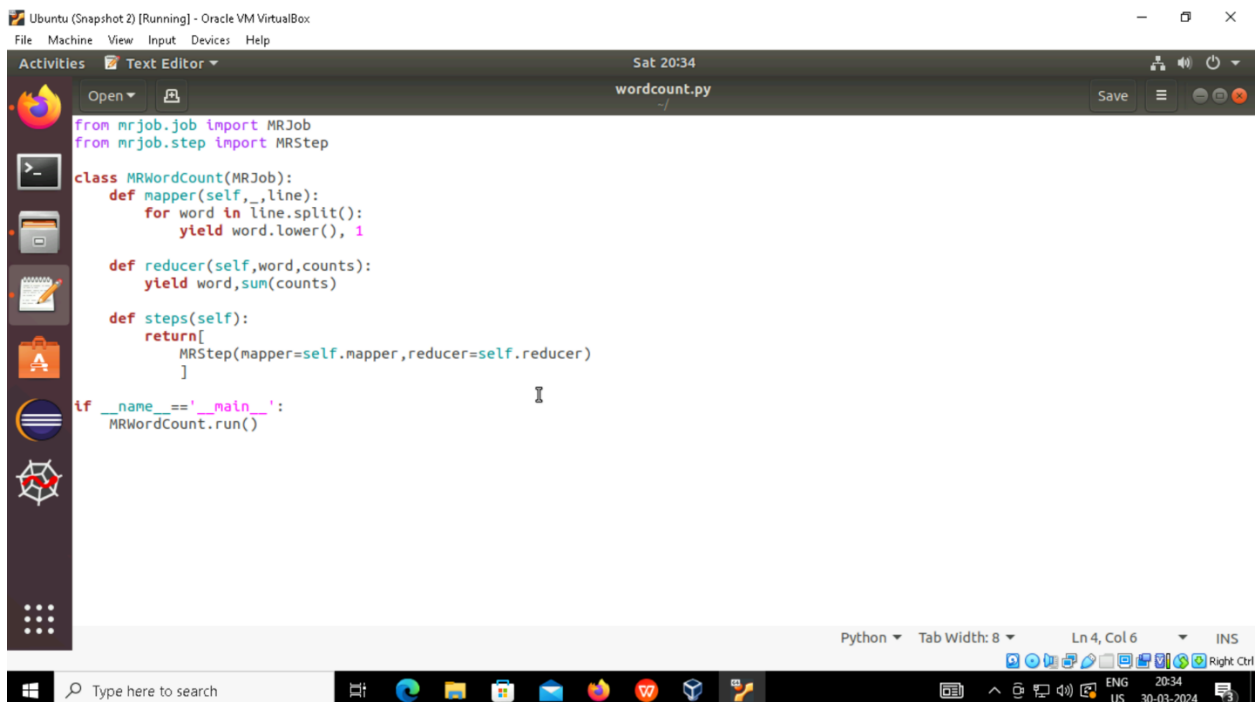# MAP REDUCE

Done by:

SATHIYA PRASAD. L

2319981

# STEP 1:

- Creating input file and wordcount python program file in local.

- Installing mrjob.

- Checking whether the function is working properly in local.

Screenshot 1 — Terminal (Sat 20:22):

```
ubh01@ubh01:~$ pip install mrjob
Collecting mrjob
  Using cached https://files.pythonhosted.org/packages/8e/58/fc28ab743aba16e9073
6ad4e29694bd2adaf7b879376ff149306d50c4e90/mrjob-0.7.4-py2.py3-none-any.whl
Collecting PyYAML>=3.10 (from mrjob)
  Using cached https://files.pythonhosted.org/packages/ba/d4/3cf562876e0cda0405e
65d351b835077ab13990e5b92912ef2bf1a2280e0/PyYAML-5.4.1-cp27-cp27mu-manylinux1_x8
6_64.whl
Installing collected packages: PyYAML, mrjob
Successfully installed PyYAML-5.4.1 mrjob-0.7.4
ubh01@ubh01:~$
```



Screenshot 2 — Terminal (Sat 20:23):

```
Collecting mrjob
  Using cached https://files.pythonhosted.org/packages/8e/58/fc28ab743aba16e9073
6ad4e29694bd2adaf7b879376ff149306d50c4e90/mrjob-0.7.4-py2.py3-none-any.whl
Collecting PyYAML>=3.10 (from mrjob)
  Using cached https://files.pythonhosted.org/packages/ba/d4/3cf562876e0cda0405e
65d351b835077ab13990e5b92912ef2bf1a2280e0/PyYAML-5.4.1-cp27-cp27mu-manylinux1_x8
6_64.whl
Installing collected packages: PyYAML, mrjob
Successfully installed PyYAML-5.4.1 mrjob-0.7.4
ubh01@ubh01:~$ python wordcount.py input.txt
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory /tmp/wordcount.ubh01.20240330.145322.549087
Running step 1 of 1...
job output is in /tmp/wordcount.ubh01.20240330.145322.549087/output
Streaming final output from /tmp/wordcount.ubh01.20240330.145322.549087/output..
.
"car"   3
"river" 2
"dear"  1
"deer"  1
"bear"  2
Removing temp directory /tmp/wordcount.ubh01.20240330.145322.549087...
ubh01@ubh01:~$
```

# STEP 2:

- Moving both the files to Hdfs.
- Executing the command to get the wordcount output.

Ubuntu (Snapshot 2) [Running] - Oracle VM VirtualBox

File   Machine   View   Input   Devices   Help

Activities      Terminal ▾                                    Sat 20:30

ubh01@ubh01: ~

File  Edit  View  Search  Terminal  Help

```
        Total committed heap usage (bytes)=547880960
        Virtual memory (bytes) snapshot=5750063104
    Shuffle Errors
        BAD_ID=0
        CONNECTION=0
        IO_ERROR=0
        WRONG_LENGTH=0
        WRONG_MAP=0
        WRONG_REDUCE=0
job output is in hdfs:///user/ubh01/tmp/mrjob/wordcount.ubh01.20240330.145938.05
5111/output
Streaming final output from hdfs:///user/ubh01/tmp/mrjob/wordcount.ubh01.2024033
0.145938.055111/output...
"bear"  2
"car"   3
"dear"  1
"deer"  1
"river" 2
STDERR: 24/03/30 20:30:15 WARN util.NativeCodeLoader: Unable to load native-hado
op library for your platform... using builtin-java classes where applicable
Removing HDFS temp directory hdfs:///user/ubh01/tmp/mrjob/wordcount.ubh01.202403
30.145938.055111...
Removing temp directory /tmp/wordcount.ubh01.20240330.145938.055111...
ubh01@ubh01:~$
```