**EXP NO: 2     RUN A BASIC WORD COUNT MAP REDUCE PROGRAM TO**
**UNDERSTAND MAP REDUCE PARADIGM**
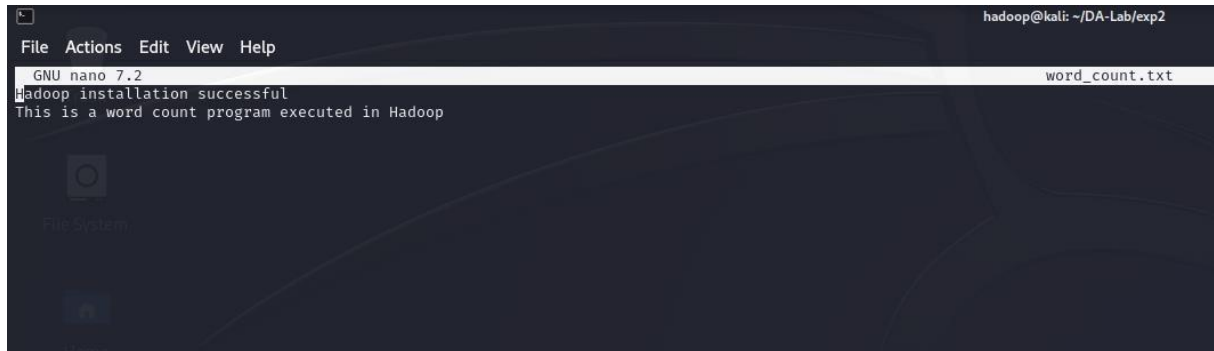
**$mkdir DA-Lab**
**$cd DA-Lab**
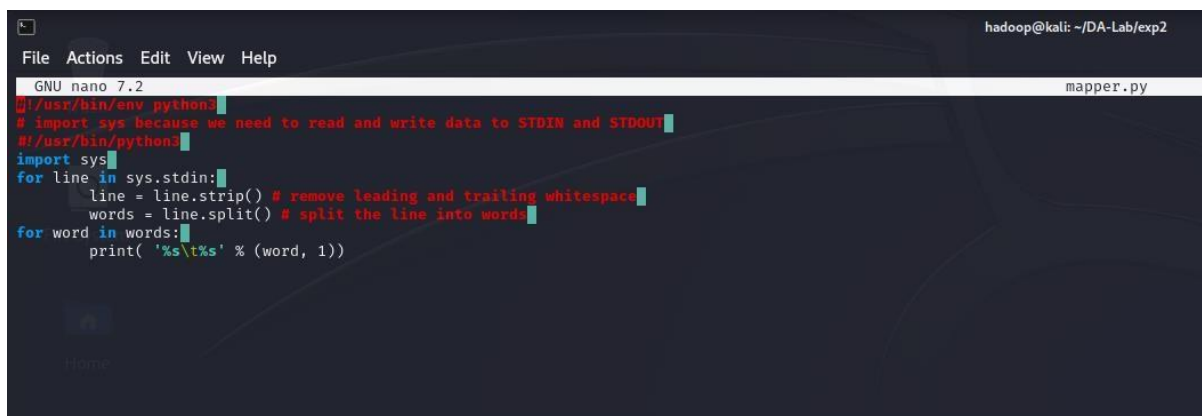**$mkdir exp2**
**$cd exp2**

**$nano word_count.txt**



**$nano mapper.py**



**$nano reducer.py**

**$start-all.sh**

```
                                                                              hadoop@kali: ~
File Actions Edit View Help
  ┌──(hadoop㉿kali)-[~]
  └─$ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as hadoop in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [kali]
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
2024-09-11 04:59:16,429 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Starting resourcemanager
Starting nodemanagers
```

## $ jps

```
  ┌──(hadoop㉿kali)-[~]
  └─$ jps
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
14436 NodeManager
16772 Jps
13830 SecondaryNameNode
14311 ResourceManager
13597 DataNode
13471 NameNode

  ┌──(hadoop㉿kali)-[~]
  └─$ 
```

**$hdfs dfs -mkdir /exp2**

**$hdfs dfs -copyFromLocal ~/DA-Lab/exp2/word_count.txt /exp2**

```
  ┌──(hadoop㉿kali)-[~/hadoop/bin]
  └─$ ./hdfs dfs -ls /exp2
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
2024-09-21 00:05:07,404 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Found 2 items
drwxr-xr-x   - hadoop supergroup          0 2024-09-13 01:05 /exp2/output
-rw-r--r--   1 hadoop supergroup         80 2024-09-13 01:02 /exp2/word_count.txt
```

**$chmod 777 mapper.py reducer.py**

**$hadoop jar $HADOOP_STREAMING -input /exp2/word_count.txt -output /exp2/output -mapper ~/DA-Lab/exp2/mapper.py -reducer ~/DA-Lab/exp2/reducer.py**

**$hdfs dfs -cat /exp2/output/\***

```
┌─(hadoop◉kali)-[~/hadoop/bin]
└─$ ./hdfs dfs -cat /exp2/output/*
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
2024-09-21 00:07:24,178 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform...
Hadoop  1
This    1
a       1
count   1
executed        1
in      1
is      1
program 1
word    1
```

```
┌─(hadoop◉kali)-[~/hadoop/bin]
└─$ ./hdfs dfs -cat /exp2/output/*
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
2024-09-21 00:07:24,178 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform...
```