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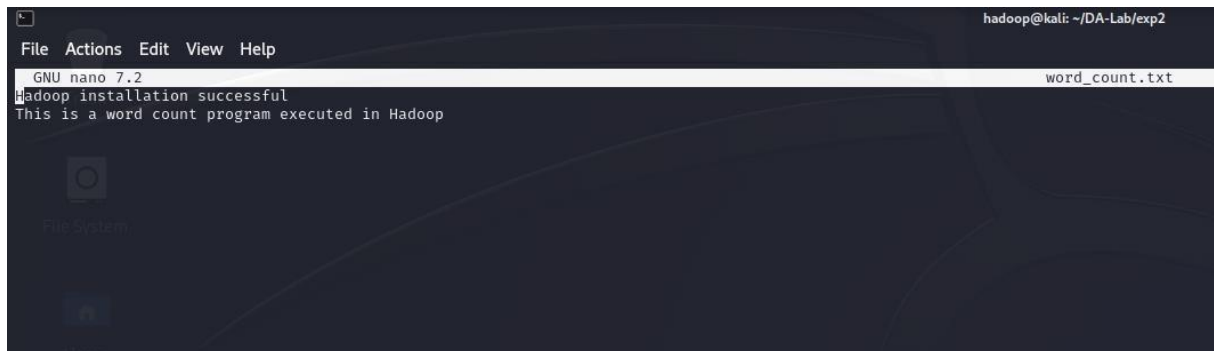
**$mkdir DA-Lab**
**$cd DA-Lab**
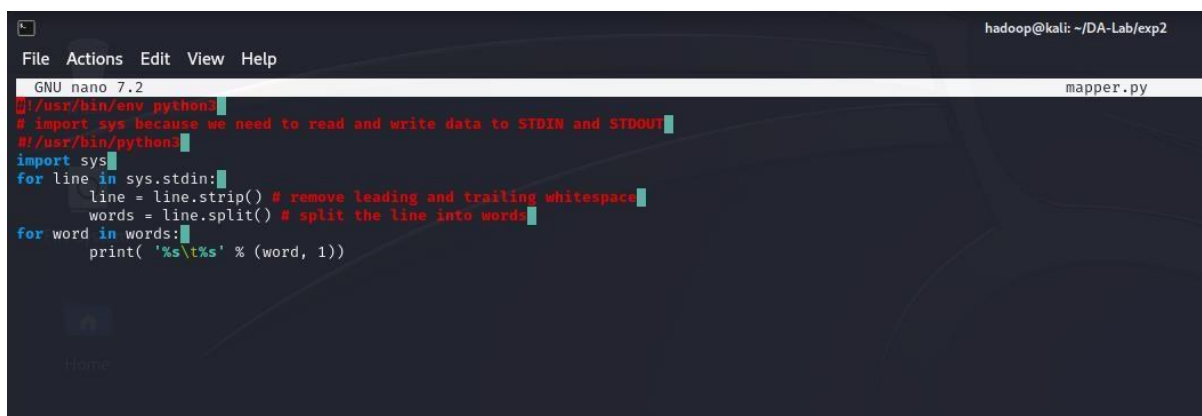**$mkdir exp2**
**$cd exp2**

**$nano word_count.txt**

```
GNU nano 7.2                                                    word_count.txt
Hadoop installation successful
This is a word count program executed in Hadoop
```

**$nano mapper.py**

```
GNU nano 7.2                                                        mapper.py
#!/usr/bin/env python3
# import sys because we need to read and write data to STDIN and STDOUT
#!/usr/bin/python3
import sys
for line in sys.stdin:
        line = line.strip() # remove leading and trailing whitespace
        words = line.split() # split the line into words
for word in words:
        print( '%s\t%s' % (word, 1))
```

**$nano reducer.py**

```
GNU nano 7.2                                                        reducer.py
#!/usr/bin/python3
from operator import itemgetter
import sys
current_word = None
current_count = 0
word = None
for line in sys.stdin:
        line = line.strip()
        word, count = line.split('\t', 1)
        try:
                count = int(count)
        except ValueError:
                continue
        if current_word == word:
                current_count += count
        else:
                if current_word:
                        print( '%s\t%s' % (current_word, current_count))
                current_count = count
                current_word = word
if current_word == word:
        print( '%s\t%s' % (current_word, current_count))
```

**$start-all.sh**



**$ jps**



**$hdfs dfs -mkdir /exp2**

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



**$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...
```