EX 2          IMPLEMENT WORD COUNT PROGRAMS USING MAPREDUCE
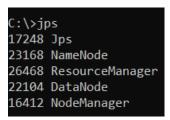
## Aim:

To implement word count/frequency program using mapReduce in Hadoop.

## Procedure:

**1.** Start the Hadoop namenode and datanode using the command
   **start-dfs.cmd**
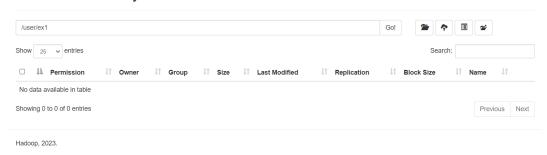   **start-yarn.cmd**

   Check if namenode and datanode are running using the comman
   **jps**

```
C:\>jps
17248 Jps
23168 NameNode
26468 ResourceManager
22104 DataNode
16412 NodeManager
```

2.  Create a directory in the Hadoop filesystem using the command

    **hadoop fs -mkdir /user/ex1**

### Browse Directory

| | /user/ex1 | | | | Go! | | | | |

Show 25 entries                                           Search:

| | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|---|
| No data available in table |

Showing 0 to 0 of 0 entries                         Previous   Next

Hadoop, 2023.

   Empty directory is created.

3.  Insert the input file into the directory using the command

    **hadoop fs -put C:\Users\jawah\OneDrive\Desktop\LathikaDA\input.txt /user/ex1**

    *//input.txt*

```
java
hello
hi
welcome
java
hello
run
execute
run
run
```

4. The MapReduce Program is written to count the frequency of word in the input file.

*//mapper.py*

```python
#!/usr/bin/env python

import sys

# Input comes from STDIN (standard input)

for line in sys.stdin:

    # Remove leading and trailing whitespace

    line = line.strip()

    # Split the line into words

    words = line.split()

    # Output each word with a count of 1

    for word in words:

        print(f'{word}\t1')
```

*//reducer.py*

```python
#!/usr/bin/env python

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(f'{current_word}\t{current_count}')
```

        current_count = count

        current_word = word

    if current_word == word:

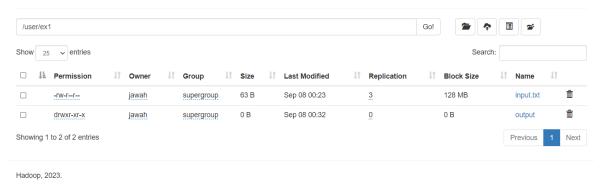        print(f'{current_word}\t{current_count}')

5. The mapper reducer program is executed by the following command

**hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.3.6.jar -input /user/ex1/input.txt -output /user/ex1/output -mapper "python C:\Users\jawah\OneDrive\Desktop\LathikaDA\mapper.py" -reducer "python C:\Users\jawah\OneDrive\Desktop\LathikaDA\reducer.py"**



Thus the output directory is created.

## Browse Directory

| | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|---|
| ☐ | -rw-r--r-- | jawah | supergroup | 63 B | Sep 08 00:23 | 3 | 128 MB | input.txt 🗑 |
| ☐ | drwxr-xr-x | jawah | supergroup | 0 B | Sep 08 00:32 | 0 | 0 B | output 🗑 |

Showing 1 to 2 of 2 entries

Hadoop, 2023.

6. To view the output files

```
C:\>hadoop fs -ls /user/ex1/output
Found 2 items
-rw-r--r--   3 jawah supergroup          0 2024-09-08 00:32 /user/ex1/output/_SUCCESS
-rw-r--r--   3 jawah supergroup         46 2024-09-08 00:32 /user/ex1/output/part-00000
```

**hadoop fs -cat /user/ex1/output/part-00000**

```
C:\>hadoop fs -cat /user/ex1/output/part-00000
execute 1
hello   2
hi      1
java    2
run     3
welcome 1
```

7. Stop the Hadoop namenode and datanode

   **stop-all.cmd**

## Result:

Thus the MapReduce Program to find the word count of a input file is completed successfully