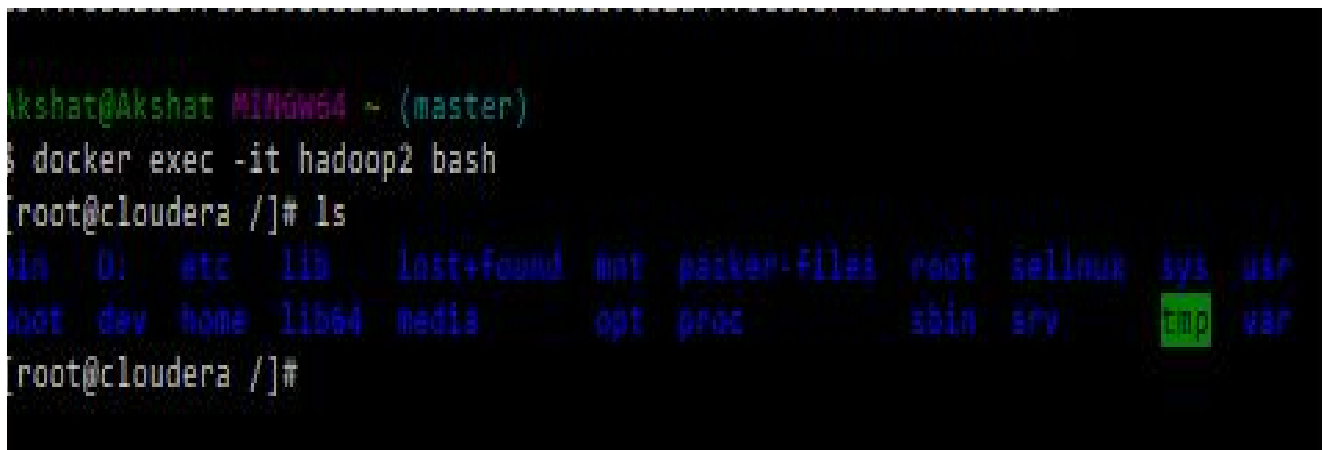# Hadoop HDFS & MapReduce

## Part I

The wordcount program from the class was modified by created a tuple containing word pairs, and printing the pairs as the output of the mapper. The reducer then counts these tuples and aggregates them. The job was performed using the Hadoop streaming service. First, we run docker on Linux machine after setting it up. Then to start Hadoop, we first enter the following command to look for a hadoop image. If the image doesn't exist, then a new image is downloaded.

*docker run --hostname=quickstart.cloudera --name=hadoop2 --privileged=true -t -d -p 80:80 -p 8080:8080 -p 8888:8888 -v /c/Users/:home cloudera/quickstart /usr/bin/docker-quickstart*

We then enter the following command to execute hadoop-

*docker exec -it hadoop2 bash*

This gives the following output -



To run the file in Hadoop, I must iterate through to the directory containing my mapper and reducer, and execute the command for it to run. Then, I will check Hue to locate my output file.

The mapping file, *mapper.py*, is as follows -

```python
#!/usr/bin/env python
"""A more advanced Mapper, using Python iterators and generators."""

import sys
import re

def read_input(file):
    for line in file:
        # split the line into words
        line = re.sub('[^a-zA-Z0-9]+',' ', line)
        yield line.split()

def main(separator='\t'):
    # input comes from STDIN (standard input)
    data = read_input(sys.stdin)
    for words in data:
        # write the results to STDOUT (standard output);
        # what we output here will be the input for the
        # Reduce step, i.e. the input for reducer.py
        #
        # tab-delimited; the trivial word count is 1
        if(len(words) == 1):
            print '%s%s%d' % (word, separator, 1)
        else:
            creator = zip(words,words[1:])
            for word in creator:
                print '%s%s%d' % (word, separator, 1)

if __name__ == "__main__":
    main()
```

The reducer file, *reduce.py*, is as follows -

```python
#!/usr/bin/env python
"""A more advanced Reducer, using Python iterators and generators."""

from itertools import groupby
from operator import itemgetter
import sys

def read_mapper_output(file, separator='\t'):
    for line in file:
        yield line.rstrip().split(separator, 1)

def main(separator='\t'):
    # input comes from STDIN (standard input)
    data = read_mapper_output(sys.stdin, separator=separator)
    # groupby groups multiple word-count pairs by word,
    # and creates an iterator that returns consecutive keys and their group:
    #   current_word - string containing a word (the key)
    #   group - iterator yielding all ["&lt;current_word&gt;", "&lt;count&gt;"] items
    for current_word, group in groupby(data, itemgetter(0)):
        try:
            total_count = sum(int(count) for current_word, count in group)
            print "%s%s%d" % (current_word, separator, total_count)
        except ValueError:
            # count was not a number, so silently discard this item
            pass

if __name__ == "__main__":
    main()
```

I have setup my mapper and reducer in the folder hadoop, located in C:\Users\Akshat\Desktop\hadoop. After iterating through this path in hadoop, I must enter the following command to submit the job to hadoop -

hadoop jar /usr/lib/hadoop-0.20-mapreduce/contrib/streaming/hadoop-streaming-2.6.0-mr1-cdh5.7.0.jar -mapper "python mapper.py" -reducer "python reducer.py" -file mapper.py -file reducer.py -input /user/cloudera/adventures.txt -output /user/cloudera/output8

The text file to use, *adventures.txt*, was placed in user/cloudera through Hue. This would run the job in Hadoop. The jar file would be created and the output would be positioned in the path specified in the command. In case of any errors in the code, the job would not complete and this would be intimated through the console and Hue's job tracker.

The output of the command above is as follows -

```
         Reduce input records=1095695
         Reduce output records=81397
         Spilled Records=2191390
         Shuffled Maps =2
         Failed Shuffles=0
         Merged Map outputs=2
         GC time elapsed (ms)=604
         CPU time spent (ms)=11290
         Physical memory (bytes) snapshot=597385216
         Virtual memory (bytes) snapshot=3919175680
         Total committed heap usage (bytes)=385228800
      Shuffle Errors
         BAD_ID=0
         CONNECTION=0
         IO_ERROR=0
         WRONG_LENGTH=0
         WRONG_MAP=0
         WRONG_REDUCE=0
      File Input Format Counters
            Bytes Read=6621217
      File Output Format Counters
            Bytes Written=946057
18/02/22 07:05:54 INFO streaming.StreamJob: Output directory: /user/cloudera/out
put3
```

Our program compiled and ran on hadoop successfully. On typing http://192.168.99.100:8888 in the web browser, Hue can be accessed. The output in Hue is as follows-



Here, part-0000 is the output file, located in the folder output8. It contains the tuple containing pair of words, separated by their frequencies in the text file.

## Part II

**Running a version of Hadoop (docker or otherwise), create a directory on HDFS with your first name (e.g. mine will be 'juan'). Submit a screen grab of the output of a Hadoop fs listing showing your home directory and your new directory in it.**

The version of Hadoop can be specified with the following command -

*hadoop version*

The output is as follows -



To create a new directory in Hadoop FS, the following command is executed -

*hadoop fs -mkdir akshat_tyagi*

This would create a directory named akshat_tyagi in the present working directory. To look at the files and directories in the current working directory, we may use the following command -

*hadoop fs -ls*

Basically, the command line for hadoop fs follows linux shell line, just that come of the important commands in linux end up as switches for hadoop fs. The outputs and the directories that we created have been shown in the screenshot of the terminal attached below -

**Submit a screen grab of Hue showing the same as (a)**

The screenshots have been attached below -



As shown above, the two directories that we created are present in /user/root. The root directory is located as follows -

| | | | | | hdfs | supergroup | drwxr-xr-x | April 05, 2016 07:26 PM |
|---|---|---|---|---|---|---|---|---|

**Submit a screen grab of your program running or completed in Hadoop. Hue has a jobs status page, use that one. Or use the command: 'hadoop jar' . See the documentation.**

The screengrab has been attached below -



The status for the job for our wordcount program is as follows -