# INDEX

# INTRODUCTION

**Data Generator**

We are importing a Faker to generate the random data. So Faker Can be used for things such as filling databases with random information during development of database related applications.
In our code it generates Name, address, city and Salary of people randomly. Our code generates file of 100MB. So to generate 10 GB file, we execute it 100 times and save it in the database. Our code can also execute for 1GB file and more depending on what size (in bytes) we provide it to the data generator script.

**Map Logic**

Map takes a set of data and converts it into another set of data, where individual elements are broken down into tuples (key/value pairs.)
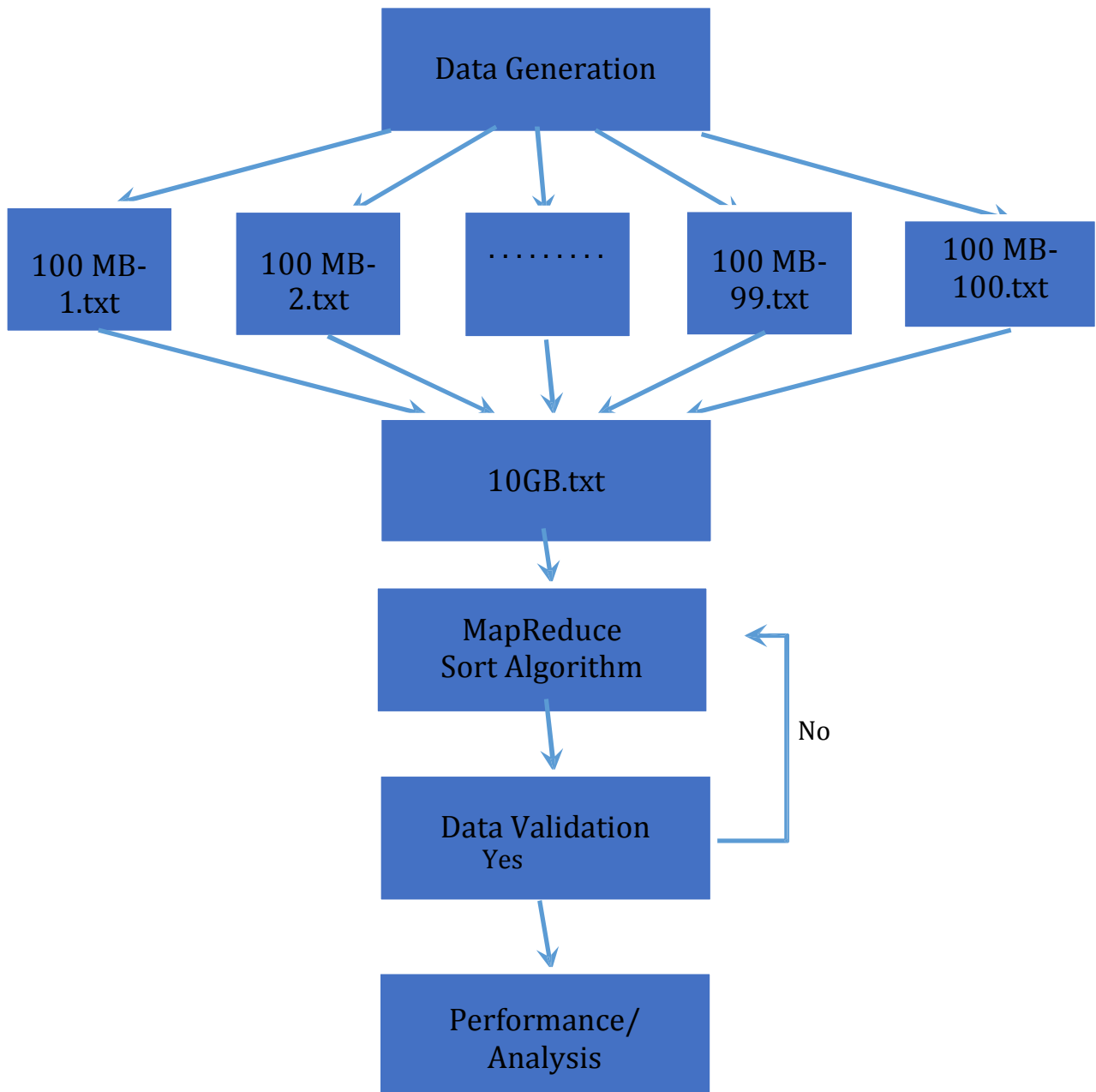So we have considered the Name and City data and we are mapping here.

**Reduce Logic**

The reduce logic we are using here is MapReduce. Reduce task, which takes the output from a map as an input and combines those data tuples into a smaller set of tuples. In the code we are using, Salary and City field is used as Reduce Key. So we are pulling the data as per the salary bracket given in our code.

**HUE**

HUE is a Web interface for analyzing data with apache Hadoop. It's a Free open source. It loads data into Hadoop. Then Views it, Process it and prepare it. Then it analyze, search and visualize it.
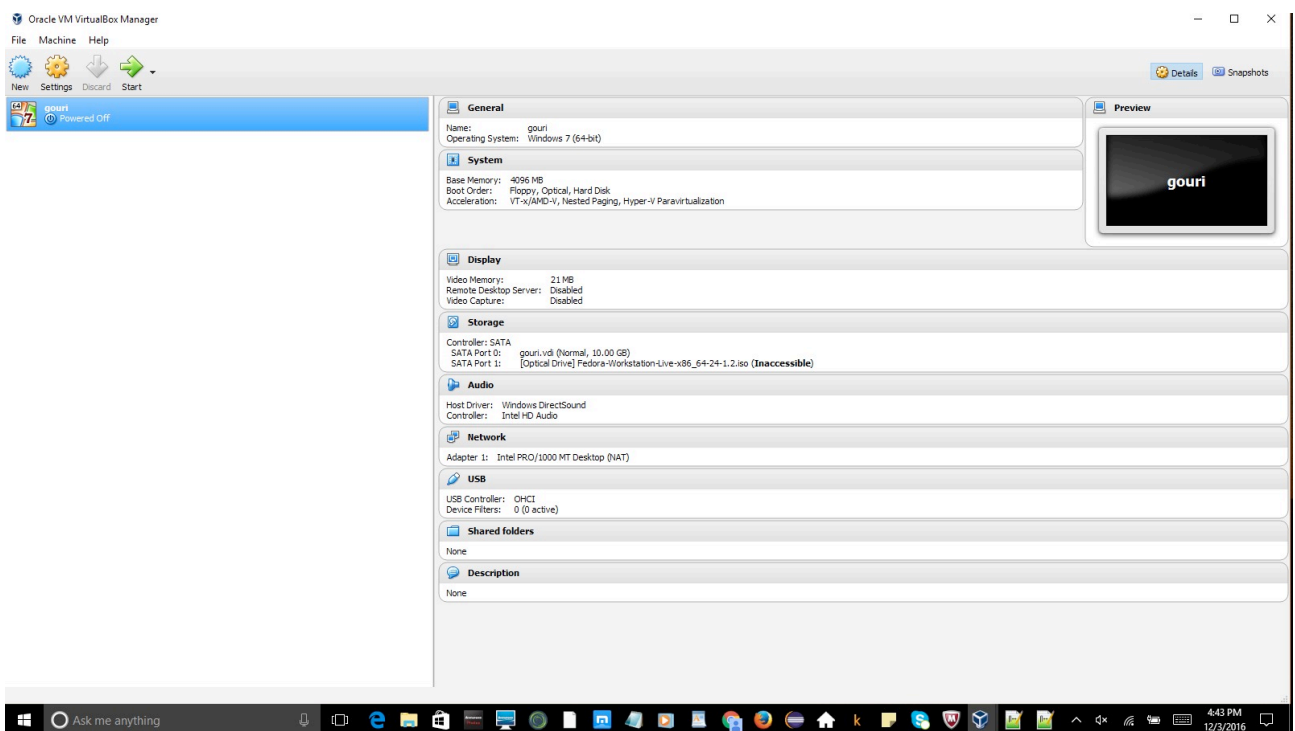
# WORKFLOW

```
┌─────────────────────┐
│   Data Generation   │
└─────────────────────┘

┌──────────┐ ┌──────────┐ ┌──────────┐ ┌──────────┐ ┌──────────┐
│ 100 MB-  │ │ 100 MB-  │ │ ........ │ │ 100 MB-  │ │ 100 MB-  │
│ 1.txt    │ │ 2.txt    │ │          │ │ 99.txt   │ │ 100.txt  │
└──────────┘ └──────────┘ └──────────┘ └──────────┘ └──────────┘

┌─────────────────────┐
│      10GB.txt       │
└─────────────────────┘

┌─────────────────────┐
│     MapReduce       │
│   Sort Algorithm    │ ◄──────┐
└─────────────────────┘        │ No
                               │
┌─────────────────────┐        │
│   Data Validation   │────────┘
│        Yes          │
└─────────────────────┘

┌─────────────────────┐
│    Performance/     │
│     Analysis        │
└─────────────────────┘
```

# SYSTEM CONFIGURATION
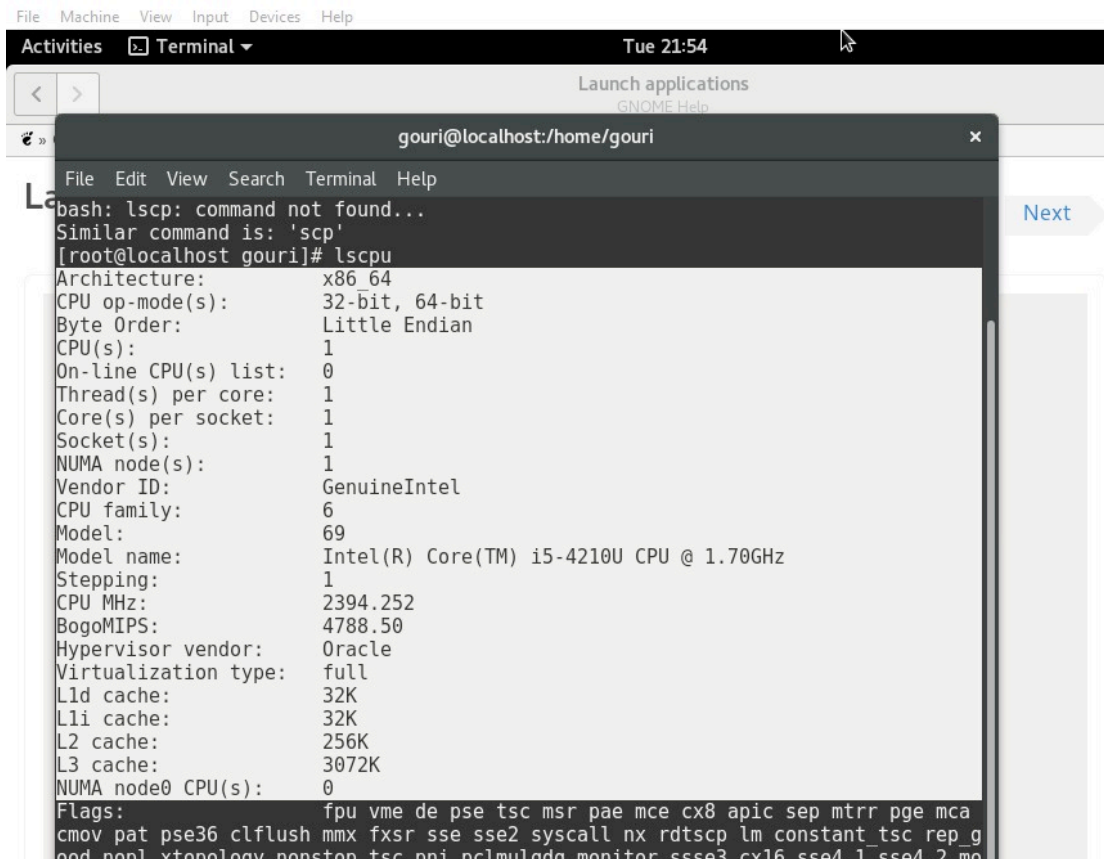
**Approach 1: Windows**

Process to install Hadoop
1. Install virtual box
2. Install fedora
3. Using terminal connect to root
4. Install java
5. Download .rpm file from Cloudera
6. Install Hadoop
7. Format the name node
8. Start the services

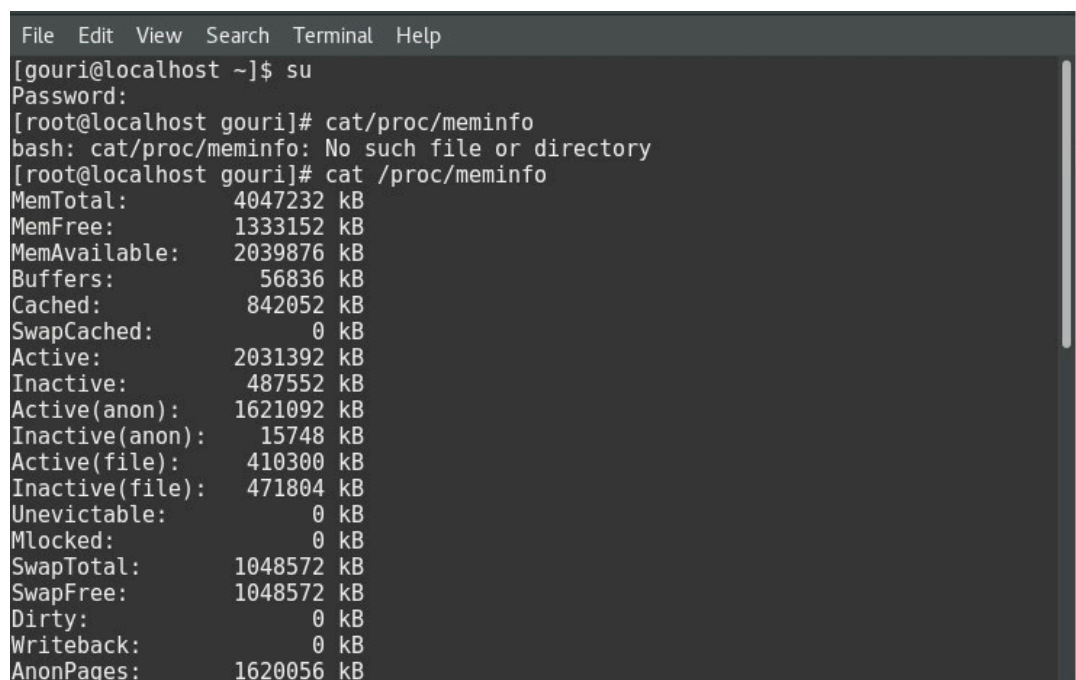After finish install Hadoop, press start to start the services



Using command lscpu to verify system configuration such as # of CPUs, cores, and so on.

Using command cat /pro/meminfo to verify memory size, and etc.



Using command cfdisk /dev/sda to verify HDD size.

Since we had few errors in our configuration files we used second approach i.e Cloudera which has in built hadoop.

**Approach 2: MAC (cloudera-quickstart-vm-5.8.0-0-virtualbox.zip)**

# INSTRUCTION

1.  Install fake-factory-0.7.2 for generating random data and python-2.7.12
    "fake_factory-0.7.2-py2.py3-none-any.whl;
    fake-factory-0.7.2.tar.gz;
    python-2.7.12-macosx10.6.pkg"
    ("https://www.python.org/ftp/python/2.7.12/python-2.7.12-macosx10.6.pkg,
    https://www.python.org/downloads/release/python-2712/.")

2.  Check all the imports in the python scripts are running and getting executed. example on
    python terminal we should be able to execute "import Faker" and so on.

3.  On command prompt run data generator script using "**python DataGenerator.py**"

```
[Taruns-MacBook-Pro:BIGDATA_Project rashmi$ python DataGenerator.py
 starting thread number: 0
 end thread number: 0
 Done
```

Refer 1: Data generation(File name "DataGenerator.py")

```python
import sys
import os
from faker import Faker
import csv
import random
from threading import Thread, Lock

def generateData(thread_num,lock):

    print "starting thread number: " + str(thread_num)
    filename = "data_"+str(thread_num)+".txt"
    fake = Faker()
    f = open( filename, 'w+b' )
    writer = csv.DictWriter(f, fieldnames=["name", "address", "city", "salary"])
    #f.close()
    #size = 1073741824 # 1gb in bytes
    size = 100000000 #100 MB file
    while os.path.getsize(filename) < size:
        writer.writerow(  dict( name=fake.name().replace(","," ").replace("\n","
        address=fake.address().replace(",","").replace("\n"," "),
        city=fake.city().replace(",","").replace("\n"," "),
        salary = random.randrange(40000,150000)))
    f.close()
    print "end thread number: " + str(thread_num)

if __name__ == '__main__':
    lock = Lock()
    threads = [Thread(target = generateData, args = (k,lock,)) for k in range(1)

    for x in threads:
        x.start()

    for x in threads:
        x.join()

    print "Done"
```

4.     This script will generate (data_0.txt)data file with provided size. We can adjust the variable "**for k in range(1)**" for generating either one data_0.txt or more than one data text file (example if we say range(2) then it will generate files data_0.txt and data_1.txt files with random data)

| | | | |
|---|---|---|---|
| data_0.txt | Nov 17, 2016, 6:05 PM | 100 MB | Plain Text |
| data_1.txt | Nov 17, 2016, 5:11 AM | 102 MB | Plain Text |

5.     Create a folder "local_data" on the virtual machine.

6.     Place the data generated files i.e. "data_0.txt" along with the "MapLogic.py" and "ReduceLogic.py" in that folder.

7.     We will be executing the python files through Hadoop **streaming.jar**

8.     On command prompt/ terminal go to the specified folder which has all the data file and the python scripts i.e. in our case "cd Desktop/local_data" (we can double check using "ls" command to see if all the files are available in local_data folder)

9.     Next step is to move all the files to Hadoop space, for that we need to go to Firefox and click on "HUE" bookmark(since we have downloaded cdh . Login to HUE with username/password as "**cloudera/cloudera**".

10.     Go to File Browser tab- create new directory—**input**

11.     Next go back to terminal—enter below command —"**hadoop fs  -copyFromLocal data_0.txt  /user/cloudera/input**"

12.     Once the command is executed successfully you should be able to see the data-0.txt file in the "/user/cloudera/input" location in HUE File browser section.

13.     The **INPUT** file is in format of "username, address, city name and salary"

14. **MAP LOGIC** :

   - Splitting each line by comma; size restriction
   - Getting the salary and name field
   - Between the map and reducer logic the data is automatically sorted with the key
   - output is comma separated.(i.e salary, name)

**MapReduce Sort Algorithm**

Map and Reduce used, why? Data will be stored in HDFS and the data will be processes with MapReduce. Processing that serially from the top to the bottom could take a long time. Instead, MapReduce is designed to process data in parallel. As a result, the file is broken in the chunks, and then processed in parallel.

MapReduce job will sort the salaries. Here is code for MapReduce.

File name MapLogic.py

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

import os
import sys


def mapper(stdin):
    for line in stdin:
        data = line.split(',')

        # perform some error checking
        if len(data) < 4:
            continue

        # unpack data
        salary = int(data[3])
        name = data[0]


        yield "%06d,%s" % (salary, name)

if __name__ == '__main__':
    for output in mapper(sys.stdin):
        print output
```

## 15.    REDUCER LOGIC

- Output from Map is taken as input by Reducer logic.
- Data is first split by comma, assigning values to variable salary and name.
- Later we are trying to match the respective salary with the associated user.

File name ReduceLogic.py

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

import sys

current_salary = None
current_name = ''

def reducer(stdin):
    current_salary = None
    current_name = ''
    for line in stdin:
        data = line.split(',')
        salary = data[0]
        name = data[1].strip()

        if current_salary == salary:
            current_name = current_name + ":" + name
        else:
            if current_salary:
                yield '%s\t%s' % (current_salary, current_name)
            current_salary = salary
            current_name = name

    if current_salary == salary:
        yield '%s\t%s' % (current_salary, current_name)

if __name__ == '__main__':
    for output in reducer(sys.stdin):
        print output
```

16.    Next step is to execute "**hadoop jar /usr/lib/hadoop-0.20-mapreduce/contrib/streaming/hadoop-streaming-2.6.0-mr1-cdh5.8.0.jar -input input/\* -**

**output output_streaming -mapper MapLogic.py -reducer ReduceLogic.py -jobconf mapred.reduce.tasks=1 -file MapLogic.py -file ReduceLogic.py**" command.

17.　　Below is the explanation for the command used:

- Input taken from **-input input/*** path
- Output is placed in **-output output_streaming** path.
- Mapper logic is read from desktop from **-mapper MapLogic.py** file.
- Reducer logic is picked from **-reducer ReduceLogic.py** file.
- **-jobconf mapred.reduce.tasks=1** part of the command says that just run one reducer so that the output will be written in one single file
- **-file MapLogic.py -file ReduceLogic.py** part of command tells us that we are sending the file to every node in the machine.
-

Execute output_console:  Output:

```
INFO mapreduce.Job: Counters: 49
        File System Counters
                FILE: Number of bytes read=29950367
                FILE: Number of bytes written=60260499
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=100008174
                HDFS: Number of bytes written=18382544
                HDFS: Number of read operations=9
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
        Job Counters
                Launched map tasks=2
                Launched reduce tasks=1
                Data-local map tasks=2
                Total time spent by all maps in occupied slots (ms)=121787
                Total time spent by all reduces in occupied slots (ms)=82150
                Total time spent by all map tasks (ms)=121787
                Total time spent by all reduce tasks (ms)=82150
                Total vcore-seconds taken by all map tasks=121787
                Total vcore-seconds taken by all reduce tasks=82150
                Total megabyte-seconds taken by all map tasks=124709888
                Total megabyte-seconds taken by all reduce tasks=84121600
        Map-Reduce Framework
                Map input records=1233781
                Map output records=1233781
                Map output bytes=27482799
                Map output materialized bytes=29950373
                Input split bytes=228
                Combine input records=0
                Combine output records=0
                Reduce input groups=1233683
                Reduce shuffle bytes=29950373
                Reduce input records=1233781
                Reduce output records=109999
                Spilled Records=2467562
                Shuffled Maps =2
                Failed Shuffles=0
                Merged Map outputs=2
                GC time elapsed (ms)=1213
                CPU time spent (ms)=27680
                Physical memory (bytes) snapshot=670920704
                Virtual memory (bytes) snapshot=4512272384
                Total committed heap usage (bytes)=583999488
```

```
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=100007946
File Output Format Counters
        Bytes Written=18382544
```

18.    In case of any errors in executing the above command double check on the location of hadoop jar on the Virtual Machine.
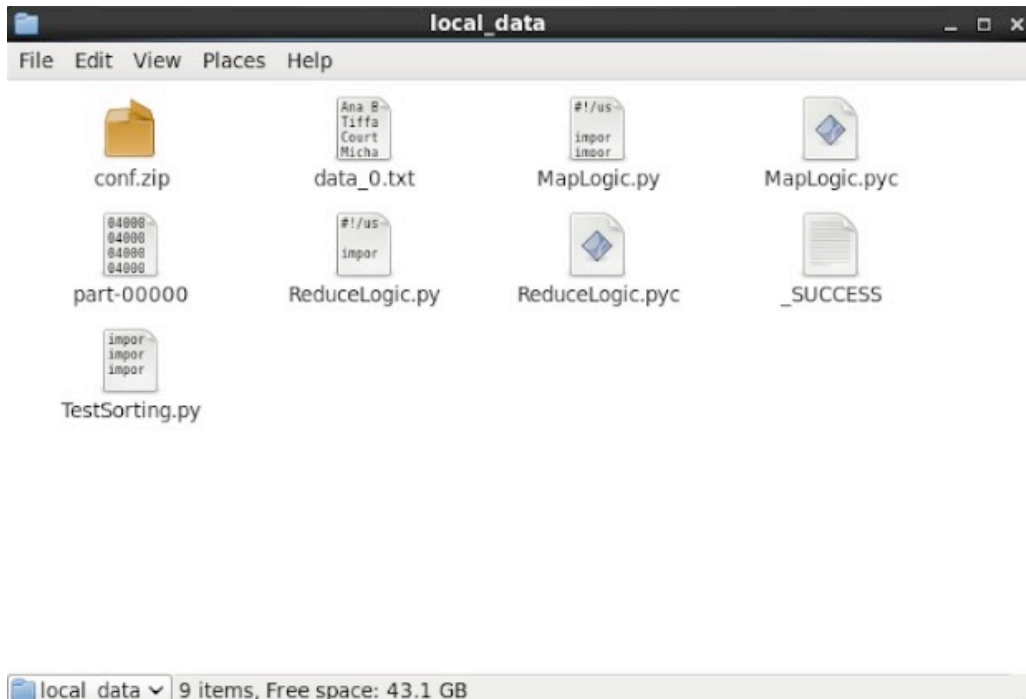
-   PATH:  "**hadoop jar /usr/lib/hadoop-0.20-mapreduce/contrib/streaming/hadoop-streaming-2.6.0-mr1-cdh5.5.0.jar**"
-   Once the command is executed successfully go back to HUE, file browser and check if the output folder is created or not. And we can find the output file which gets created.
-   To get the data from hadoop file system to local file system we can run the below command: "**hadoop fs -copyToLocal  /user/cloudera/output_streaming/\***"
-   If you enter "ls" command we should be able to see part-00000 as code and then we can check the first five or any respective entries in the output file. e.g.: "head -10 part-00000", we should be able to see part of output file where it starts from 040000 salary as the first entry, followed by 040001, 040002..... and so on.
-   To check the records from bottom of the output file we can use command: "tail -10 part-00000" where we can see the last record will have a salary of 149999(As size < 150000)



19.    The output file will be in output_streaming(**/user/cloudera/output_streaming**) folder in HUE.

In case we want to insert more input files we can place it in INPUT folder and execute the command mentioned in step 12. Also we should double check that the output_streaming folder is deleted so that it creates it newly and generates the respective output file.

Below is the local_data folder with moved output file (part-00000) from the Hadoop space and other scripts which help us execute the code using cloudera.



20.     Data Validation:

In this step, we validate the input data with the expected output. If the expected data matches the input Test data then it passes the validation test, if not the failure description is displayed on the console. We are validating both input from mapper and reducer in our validation i.e TestSorting.py code. We place this file in same local_data folder and execute it once our output is generated using the Map and reducer logic. Command to execute it "python TestSorting.py" on terminal/command prompt.

file Name: TestSorting.py

```
import unittest
import MapLogic
import ReduceLogic

class WCMapperTest(unittest.TestCase):
    def test_mapper(self):
        input = ['Ana Brennan,7728 Hall,Williamsberg,45095','Tiffany Wilkins,675
        expected = ['045095,Ana Brennan', '066323,Tiffany Wilkins', '066323,Joe
        result = []
        for output in MapLogic.mapper(input):
            result.append(output)
        self.assertEqual(expected, result)

    def test_reducer(self):
        input = ['045095,Ana Brennan', '066323,Tiffany Wilkins','066323,Joe Turn
        expected = ['045095\tAna Brennan', '066323\tTiffany Wilkins:Joe Turner']
        result = []
        for output in ReduceLogic.reducer(input):
            result.append(output)
        self.assertEqual(expected, result)

if __name__ == '__main__':
    unittest.main()
```

**Mapper**

input = ['Ana Brennan,7728 Hall,Williamsberg,45095','Tiffany Wilkins,6755 Patricia 57191,Ruthchester,66323','Joe Turner,6755 Patricia 57191,Ruthchester,66323']

expected = ['045095,Ana Brennan', '066323,Tiffany Wilkins', '066323,Joe Turner']

**Reducer**

input = ['045095,Ana Brennan', '066323,Tiffany Wilkins','066323,Joe Turner']

expected = ['045095\tAna Brennan', '066323\tTiffany Wilkins:Joe Turner']

**<u>Result</u>**

When input is valid

```
TestSorting.py---VALIDATION:
[cloudera@quickstart local_data]$ python TestSorting.py
..
----------------------------------------------------------------------
Ran 2 tests in 0.000s

OK
```

**Mapper**

input = ['Ana Brennan,7728 Hall,Williamsberg,45095','Tiffany Wilkins,6755 Patricia 57191,Ruthchester,66323','Joe Turner,6755 Patricia 57191,Ruthchester,66323']

expected = ['045095,Ana Brennan', '066323,Tiffany Wilkins', '066323,Joe Turner']

**Reducer**

input = ['045095,Ana Brennan', **'000000066323**,Tiffany Wilkins','066323,Joe Turner']

expected = ['045095\tAna Brennan', **'066323**\tTiffany Wilkins:Joe Turner']

**Result:** When input is invalid

```
-----------------------------TEST 02------------------------
[cloudera@quickstart local_data]$ python TestSorting.py
.F
==================================================================
FAIL: test_reducer (__main__.WCMapperTest)
------------------------------------------------------------------
Traceback (most recent call last):
  File "TestSorting.py", line 20, in test_reducer
    self.assertEqual(expected, result)
AssertionError: ['045095\tAna Brennan', '000000066323\tTiffany Wilkins:Joe Turner'] !=
['045095\tAna Brennan', '066323\tTiffany Wilkins:Joe Turner']

------------------------------------------------------------------

Ran 2 tests in 0.000s

FAILED (failures=1)
```
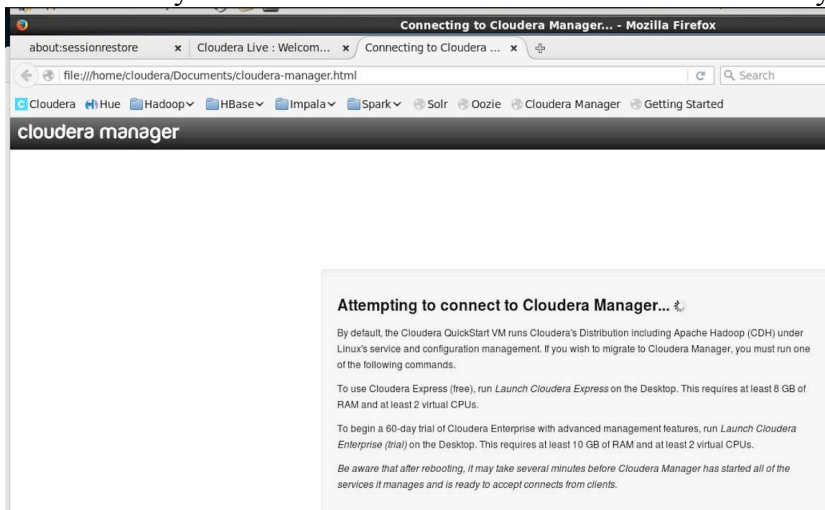
## Performance/Analysis

Performance: latency = start time – end time

As we used cloud era, we had to look at the Cloud Manager for Performance latency, but it did not load for the system. We tried to load it since two continuous days.



## Conclusion/Summary

We were successfully able to install Hadoop on system. We performed the data generation, data sorting and data validation on cloudera which was successful. We followed two approaches to perform the above operations on data, using Hadoop and cloudera. So we gained knowledge about both the platforms.

## References

www.**cloudera**.com/products/apache-hadoop/**hue**.html
https://www.tutorialspoint.com/hadoop/hadoop_mapreduce.htm
www.**cloudera**.com/
**www.hadoop.apache.org**

# APPENDIX

**DataGenerator.py**

```python
import sys
import os
from faker import Faker
import csv
import random
from threading import Thread, Lock

def generateData(thread_num,lock):

    print "starting thread number: " + str(thread_num)
    filename = "data_"+str(thread_num)+".txt"
    fake = Faker()
    f = open( filename, 'w+b' )
    writer = csv.DictWriter(f, fieldnames=["name", "address", "city", "salary"])
    #f.close()
    size = 1073741824 # 1gb in bytes
    #size = 100000000 #100 MB file
    while os.path.getsize(filename) < size:
        writer.writerow(  dict( name=fake.name().replace(","," ").replace("\n"," "),
            address=fake.address().replace(",","").replace("\n"," "),
            city=fake.city().replace(",","").replace("\n"," "),
            salary = random.randrange(40000,150000)))
    f.close()
    print "end thread number: " + str(thread_num)

if __name__ == '__main__':
    lock = Lock()
    threads = [Thread(target = generateData, args = (k,lock,)) for k in range(1)]

    for x in threads:
        x.start()

    for x in threads:
        x.join()

    print "Done"
```

**ReduceLogic.py**

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

import sys

current_salary = None
current_name = ''

def reducer(stdin):
    current_salary = None
    current_name = ''
    for line in stdin:
        data = line.split(',')
```

```python
        salary = data[0]
        name = data[1].strip()

        if current_salary == salary:
            current_name = current_name + ":" + name
        else:
            if current_salary:
                yield '%s\t%s' % (current_salary, current_name)
            current_salary = salary
            current_name = name

    if current_salary == salary:
        yield '%s\t%s' % (current_salary, current_name)

if __name__ == '__main__':
    for output in reducer(sys.stdin):
        print output
```

---

**MapLogic.py**

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

import os
import sys


def mapper(stdin):
    for line in stdin:
        data = line.split(',')

        # perform some error checking
        if len(data) < 4:
            continue

        # unpack data
        salary = int(data[3])
        name = data[0]


        yield "%06d,%s" % (salary, name)

if __name__ == '__main__':
    for output in mapper(sys.stdin):
        print output
```

---

**TestSorting.py**

```python
import unittest
import MapLogic
import ReduceLogic

class WCMapperTest(unittest.TestCase):
    def test_mapper(self):
        input = ['Ana Brennan,7728 Hall,Williamsberg,45095','Tiffany Wilkins,6755 Patricia
57191,Ruthchester,66323','Joe Turner,6755 Patricia 57191,Ruthchester,66323']
        expected = ['045095,Ana Brennan', '066323,Tiffany Wilkins', '066323,Joe Turner']
```

```python
        result = []
        for output in MapLogic.mapper(input):
            result.append(output)
        self.assertEqual(expected, result)

    def test_reducer(self):
        input = ['045095,Ana Brennan', '066323,Tiffany Wilkins','066323,Joe Turner']
        expected = ['045095\tAna Brennan', '066323\tTiffany Wilkins:Joe Turner']
        result = []
        for output in ReduceLogic.reducer(input):
            result.append(output)
        self.assertEqual(expected, result)

if __name__ == '__main__':
    unittest.main()
```