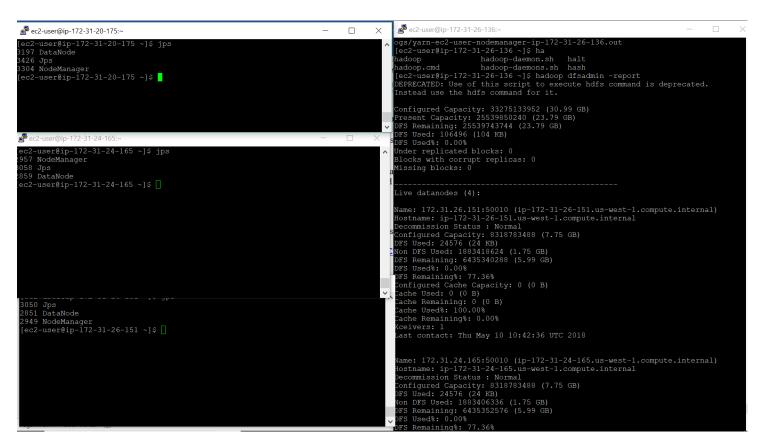
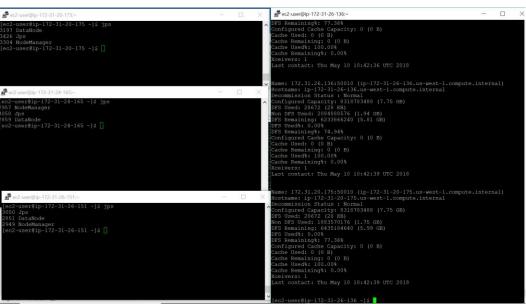
Final Project Phase-1 Nisarg Patel

1) Creating 4 Node cluster





Word Count time taken for 4 node cluster compare to 1 node cluster

```
real 0m36.320s
user 0m3.824s
sys 0m0.195s
[ec2-user@ip-172-31-26-136 ~]$

Bytes Written=20056175

real 1m20.660s
user 0m4.240s
sys 0m0.164s
[ec2-user@ip-172-31-26-246 ~]$
```

2) Time Taken to run query 1.2, 1.3, 2.1

```
2018-05-12 07:30:04,807 Stage-2 map = 0%,
2018-05-12 07:30:13,250 Stage-2 map = 25%,
                                                                                                                                                                                                  2018-05-12 07:39:12,041 Stage-2 map = 0%, reduce = 0%
2018-05-12 07:39:18,298 Stage-2 map = 25%, reduce = 0%, C
mulative CPU 2.83 sec
2018-05-12 07:39:19,332 Stage-2 map = 50%, reduce = 0%, C
                                                                                                                                             reduce = 0%. Cu
                                                                                                                                                                                                mulative CPU 6.43 sec
2018-05-12 07:39:22,451 Stage-2 map = 75%, reduce = 0%, Cu
mulative CPU 9.99 sec
2018-05-12 07:39:23,517 Stage-2 map = 100%, reduce = 0%, Cu
mulative CPU 14.86 sec
2018-05-12 07:39:25,594 Stage-2 map = 100%, reduce = 100%,
Cumulative CPU 16.25 sec
MapReduce Total cumulative CPU time: 16 seconds 250 msec
Ended Job = job_1526109783349_0002
MapReduce Jobs Launched:
Stage-Stage-2: Map: 4 Reduce: 1 Cumulative CPU: 16.25 se
C HDFS Read: 594378990 HDFS Write: 11 SUCCESS
Total MapReduce CPU Time Spent: 16 seconds 250 msec
2018-05-12 07:30:17,383 Stage-2 map = 75%, reduce = 0%, Cu
  2018-05-12 07:30:19,455 Stage-2 map = 100%,
                                                                                                                                                reduce = 100%,
 MapReduce Total cumulative CPU time: 16 seconds 360 msec
Ended Job = job_1526109783349_0001
MapReduce Jobs Launched:
 Stage-Stage-2: Map: 4 Reduce: 1 Cumulative CPU: 16.36 se
9110 HDFS Write: 12 SUCCESS
Total MapReduce CPU Time Spent: 16 seconds 360 msec
                                                                                                                                                                                                 OK
4435791464
 Time taken: 34.389 seconds, Fetched: 1 row(s)
--Q1.2 Simplified to remove expression in sum
select sum(lo_extendedprice) as revenue
from lineorder, dudate
where lo_orderdate = d_datekey
and d_yearmonth = 'Jan1993'
and lo_discount between 5 and 6
and lo_quantity between 25 and 35;
                                                                                                                                                                                                 --Q1.3 Simplified to remove expression in sum
select sum(lo_extendedprice) as revenue
from lineorder, dudate
where lo_orderdate = d_datekey
and d_weeknuminyear = 6 and d_year = 1994
and lo_discount between 5 and 8
and lo_quantity between 36 and 41;
                                                                                                        361416497
                                                                                                                                                                                                 MFGR#128
                                                                                                                                                                                                MFGR#129
                                                                                                          ime taken: 108.893 seconds, Fetched: 280 row(s)
                                                                                                       hive>
                                                                                                       --Q2.1 No simpifications
select sum(lo_revenue), d_year, p_brand1
from lineorder, dwdate, part, supplier
where lo_orderdate = d_datekey
                                                                                                       where lo_orderoate = d_dateki
and lo_partkey = p_partkey
and lo_suppkey = s_suppkey
and p_category = 'MFGR#12'
and s_region = 'AMERICA'
group by d_year, p_brand1
order by d_year, p_brand1;
```

Adding the file :-ADD FILE /home/ec2-user/10char.py

Now we will do select transform :-

INSERT OVERWRITE TABLE customer3 SELECT

TRANSFORM(c_custkey,c_name,c_address,c_city,c_nation,c_region,c_phone,c_mktsegment)
USING 'python 10char.py' AS

(c_custkey,c_name,c_address,c_city,c_nation,c_region,c_phone,c_mktsegment) FROM customer;

10Char.py python file :-

Describe Table:-

```
ive> describe customer;
                         varchar(25)
varchar(10)
__city
_nation
_phone
                          varchar(15)
                          varchar(10)
 mktseament
ime taken: 0.041 seconds, Fetched: 8 row(s)
hive> describe customer3;
 name
                         varchar (25)
varchar (12)
 address
 _city
_nation
                          varchar (12)
 phone
                          varchar (15)
 __mktsegment
 ime taken: 0.038 seconds, Fetched: 8 row(s)
```

Output:

```
ec2-user@ip-172-31-26-136:~/apache-hive-2.0.1-bin
     select c address, c city from customer3 where c custkey
5JsirBM9P
                 MOROCCO
                 JORDAN
487LW1dovn
fkRGN8n ARGENTINA 7
4u58h f EGYPT
nwBtxkoBF
gls,pzDen
OkMVLQ1dK
                 SAUDI ARA 2
,pZ,Qp,qt
gIql8H6zo
                 INDIA
  mQ6Ug9U
                 ETHIOPIA 9
                 UNITED KI 3
                 JORDAN
b4qxKs7
z3ax0D5Hn
                 CANADA
                 ARGENTINA 0
y4KK4CcfN
                 UNITED KI 0
 .
2IQMff18e
 yukcsqIxl
                 FRANCE
 OOXPkiuSW
                              Fetched: 19 row(s
```

3)

lorder = LOAD '/data/lineorder.tbl' USING PigStorage('|') AS
(lo_orderkey:int,lo_linenumber:int,lo_custkey:int,lo_partkey:int,lo_suppkey:int,lo_orderdate:int,lo_
orderpriority:chararray,lo_shippriority:chararray,lo_quantity:int,lo_extendedprice:int,lo_ordertotalp
rice:int,lo_discount:int,lo_revenue:int,lo_supplycost:int,lo_tax:int,lo_commitdate:int,lo_shipdate:ch
ararray);

QUERY 0.1

grouplorder = group lorder all; avgrevenue = FOREACH grouplorder GENERATE AVG(lorder.lo_revenue);

```
2018-05-13 00:10:53,206 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success! 2018-05-13 00:10:53,207 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not ge 2018-05-13 00:10:53,215 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process: 1 2018-05-13 00:10:53,215 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to proc (3634300.709514323) grunt>
```

Average lo_revenue = 3634300.7095

Time to run query =

```
ecutionEngine - Connecting to hadoop file system at: hdfs://1/2.31.26.136/
2018-05-14 03:48:31,798 [main] INFO org.apache.hadoop.conf.Configuration.deprecation
- fs.default.name is deprecated. Instead, use fs.defaultFS
2018-05-14 03:48:31,915 [main] INFO org.apache.pig.Main - Pig script completed in 1 s
econd and 998 milliseconds (1998 ms)
[ec2-user@ip-172-31-26-136 pig-0.15.0]$
```

QUERY 0.2

groupdiscount = group lorder BY lo_discount;
outputdiscount = FOREACH groupdiscount GENERATE lorder.lo_discount, COUNT(
lorder.lo_extendedprice);

```
10), (10), (10), (10), (10), (10), (10), (10), (10), (
, (10), (10), (10), (10), (10), (10), (10), (10), (10)
(10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10), (10)
```

Time taken to run query =

```
ecutionEngine - Connecting to hadoop file system at: hdfs://172.31.26.136/
2018-05-14 03:56:23,282 [main] INFO org.apache.hadoop.conf.Configuration.deprecation
- fs.default.name is deprecated. Instead, use fs.defaultFS
2018-05-14 03:56:23,397 [main] INFO org.apache.pig.Main - Pig script completed in 1 s
econd and 889 milliseconds (1889 ms)
```

OUERY 0.3

```
grunt> filterdiscount = FILTER lorder BY lo_discount < 3;
grunt> groupquantity = group filterdiscount BY lo_quantity;
grunt> outputquantity = FOREACH groupquantity GENERATE(filterdiscount.lo_quantity), SUM(filterdiscount.lo_revenue)
grunt>
```

```
(49), (49), (49), (49), (49), (49), (49), (49), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50), (50),
```

Time to run query =

```
- fs.default.name is deprecated. Instead, use fs.defaultFS
2018-05-14 03:53:43,674 [main] INFO org.apache.pig.backend.hadoop.executionengine.HEx
ecutionEngine - Connecting to hadoop file system at: hdfs://172.31.26.136/
2018-05-14 03:53:44,583 [main] INFO org.apache.hadoop.conf.Configuration.deprecation
- fs.default.name is deprecated. Instead, use fs.defaultFS
2018-05-14 03:53:44,708 [main] INFO org.apache.pig.Main - Pig script completed in 1 s
econd and 887 milliseconds (1887 ms)
```

Query 0.2 Mapper

```
import sys
import fileinput

for line in sys.stdin:
    data = line.strip().split('|')
    lo_discount = data[11]
    lo_extendedprice = data[9]
    lo_discount = int(lo_discount)
    print "{0}\t{1}".format(lo_discount,1)
```

Reducer

```
GNU nano 2.5.3 File: countreducer.py

import sys
import fileinput

curr_count = 0
currentkey = None

for line in sys.stdin:
    line = line.strip().split('\t')
    key, count = line
    count=int(count)

    if currentkey == line:
        curr_count += count

    else:
        if currentkey:
            print "{0}\t{1}".format(currentkey, curr_count)

    if currentkey == key:
        print "{0}\t{1}" (currentkey, curr_count))
```

Query 0.3

Mapper

Reducer