

IST769 Homework Submission

Name: Sathish Kumar Rajendiran

SUID: 666555028

Email: srajendi@syr.edu

Due Date: 08/17/2021

Task: Hive and Pig

Homework #:6

Exercise(s):

1. Load the comma-delimited HDFS dataset at [clickstream/iplookup](#) into a relation with an explicit schema. Use filter logic to remove the first row (which contains a header) then sort the output by IP and dump a comma-delimited data set to [clickstream/iplookup_noheader](#). Record all your Pig commands required to complete your transformation.

Solution:

--list available folders

fs -ls clickstream/iplookup

--view the contents of the file

```
fs -cat clickstream/iplookup/ip_lookup.csv
```

--Load the file using PigStorage with explicit schema

```
iplookup_File = LOAD 'clickstream/iplookup/ip_lookup.csv' USING  
PigStorage(',') AS  
(ip:chararray,country:chararray,state:chararray,city:chararray,lat:double,lan:  
double);
```

--Remove the header row from the file

iplookup File Noheader = FILTER iplookup File by ip! ='IP':

--Sort the file

iplookup File Noheader sorted = ORDER iplookup File Noheader by ip ASC:

--Store the file into HDFS

STORE iplookup File Noheader sorted into

'/user/cloudera/clickstream/iplookup_noheader' USING PigStorage(',');

--view the contents of the file

```
fs -cat clickstream/iplookup noheader/part-r-00000
```

Evidence:

769-Win10Docker-srajendi

```
cloudera@quickstart:~$ ls clickstream/iplookup_noheader
Found 2 items
-rw-r--r--  1 cloudera cloudera          0 2021-08-13 21:19 clickstream/iplookup_noheader/_SUCCESS
-rw-r--r--  1 cloudera cloudera  1185 2021-08-13 21:19 clickstream/iplookup_noheader/part-r-00000
cloudera@quickstart:~$ cat clickstream/iplookup_noheader/part-r-00000
128.122.140.238,USA,NY,New York,40.712784,-74.005941
128.230.122.180,USA,NY,Syracuse,43.048122,-76.147424
155.100.169.152,USA,UT,Salt Lake City,40.760779,-111.891047
172.189.252.8,USA,VA,Dulles,38.955855,-77.447819
215.82.23.2,USA,OH,Columbus,39.961176,-82.998794
38.68.15.223,USA,TX,Dallas,32.776664,-96.796988
54.114.107.209,USA,NJ,Jersey City,40.728157,-74.077642
56.216.127.219,USA,NC,Raleigh,35.77959,-78.638179
68.199.40.156,USA,NY,Freeport,40.657602,-73.583184
70.209.14.54,USA,FL,Tampa,27.950575,-82.457178
74.111.18.59,USA,NY,Syracuse,43.048122,-76.147424
74.111.6.173,USA,VA,Arlington,38.87997,-77.10677
8.37.70.112,USA,CA,Los Angeles,34.052234,-118.243685
8.37.70.170,USA,CA,Los Angeles,34.052234,-118.243685
8.37.70.226,USA,CA,Los Angeles,34.052234,-118.243685
8.37.70.77,USA,CA,Los Angeles,34.052234,-118.243685
8.37.70.99,USA,CA,Los Angeles,34.052234,-118.243685
8.37.71.25,USA,CA,Los Angeles,34.052234,-118.243685
8.37.71.43,USA,CA,Los Angeles,34.052234,-118.243685
8.37.71.57,USA,CA,Los Angeles,34.052234,-118.243685
8.37.71.69,USA,CA,Los Angeles,34.052234,-118.243685
8.37.71.9,USA,CA,Los Angeles,34.052234,-118.243685
98.29.25.44,USA,OH,Cleveland,41.49932,-81.694361
grun
```

2. Write Pig commands to produce a count of IP Addresses by state codes, sorted by the count with highest values first, like this:
(CA, 10)
(NY, 4)
(VA, 2)
Etc....

Record all your Pig commands required to complete your transformation.

Solution:

```
--Group by State
iplookup_by_state = GROUP iplookup_File_Noheader_sorted by
state;
--view the schema
describe iplookup_by_state;
--count number of ip's by state
iplookup_by_state_n_counts = FOREACH iplookup_by_state
GENERATE group as state,
COUNT(iplookup_File_Noheader_sorted.ip) as counts;
--sort number of ip's by state
```

iplookup_by_state_n_counts_sorted = ORDER

iplookup_by_state_n_counts DESC;

--view the output

DUMP iplookup_by_state_n_counts_sorted;

Evidence:

```

206701 [jobcontrol] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
2867691 [jobcontrol] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths (combined) to process : 1
2868899 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - JobId: job_1628829248144_0007
2868899 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases iplookup_by_state
2868899 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - detailed information at: M: iplookup_by_state[11,20] C: R:
2868899 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - detailed information at: http://localhost:50636/jobdetails.jsp?jobid=job_1628829248144_0007
2879524 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 87% complete
2881893 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
288221 [main] INFO org.apache.pig.tools.pigstats.SimplePigStats - Script Statistics:
HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.6.0-cdh5.7.0 0.12.0-cdhs7.0 cloudera 2021-08-13 21:44:07 2021-08-13 21:45:44 GROUP_BY,BY_ORDER_BY,FILTER

Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime
job_1628829248144_0004 1 0 3 3 3 3 n/a n/a n/a iplookup_File,iplookup_File_Noheader
job_1628829248144_0005 1 1 3 3 3 3 3 3 3 iplookup_File_Noheader_sorted
AMPLER
job_1628829248144_0006 1 1 3 3 3 3 3 3 3 iplookup_File_Noheader_sorted
ROVER
job_1628829248144_0007 1 1 3 3 3 3 3 3 3 iplookup_by_state
GROUP_BY hdfs://quickstart.cloudera:8020/tmp/temp-196607337/tmp-1545452863,
Input(s):
Successfully read 24 records (1658 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/clickstream/iplookup/ip.lookup.csv"

Output(s):
Successfully stored 9 records (1415 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-196607337/tmp-1545452863"

Counters:
Total records written : 9
Total records written : 1415
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1628829248144_0004 -> job_1628829248144_0005,
job_1628829248144_0005 -> job_1628829248144_0006,
job_1628829248144_0006 -> job_1628829248144_0007,
```

```

769-Win10Docker-srajendi
[Enforce]

cloudera@quickstart:~$ grunts describe iplookup_by_state;
iplookup_by_state: {group: chararray,iplookup_File_Noheader_sorted: {(ip: chararray,country: chararray,state: chararray,city: chararray,lat: double,lon: double)}}
grunts iplookup_by_state_n_counts = FOREACH iplookup_by_state GENERATE group as state,COUNT(iplookup_File_Noheader_sorted.ip) as counts;
grunts DUMP iplookup_by_state_n_counts;
3899315 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - Pig features used in the script: GROUP_BY_ORDER_BY_FILTER
3899315 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, DuplicateForEachColumnRewrite, GroupByConstP
zen, LoadTypeCastInserter, MergeFilter, MergeForEach, NewPartitionFilterOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter], RULES_DISABLED=[Fi
lmerizer]}
3899344 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.CombinerOptimizer - Choosing to move algebraic branch to combiner
3899350 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 6
3899351 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 6
3899353 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
3899354 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
3899374 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
3899446 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - jar file Job174910676869897926.jar created
3899448 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
38994481 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] will not generate code
3899450 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Starting process to move generated code to distributed cache
3899451 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple] with classes to deserialize []
38994504 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Pig script settings are added to the job
38994548 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
38994553 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - MapReduceLauncher - HadoopJobId: job_1628829248144_0003
3899500 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - More information at: http://localhost:50030/jobdetails.jsp?jobid=job_16288
38995004 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995005 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995006 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995007 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995008 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995009 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995010 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995011 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995012 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995013 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995014 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995015 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995016 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995017 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995018 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995019 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995020 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995021 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995022 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995023 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995024 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995025 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995026 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995027 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995028 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995029 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995030 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995031 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995032 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995033 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Job completed
38995034 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
38995327 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Setting key [pig.schematuple.classes] with classes to deserialize []
38995472 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
38995473 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - HadoopJobId: job_1628829248144_0009
38995478 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLauncher - Processing aliases iplookup_File_Noheader_sorted

```

```

769-Win10Docker-srajendi
[Enforce]

cloudera@quickstart:~$ Job DAG:
job_1628829248144_0008 --> job_1628829248144_0009,
job_1628829248144_0009 --> job_1628829248144_0010,
job_1628829248144_0010 --> job_1628829248144_0011,
job_1628829248144_0011 --> job_1628829248144_0012,
job_1628829248144_0012 --> job_1628829248144_0013,
job_1628829248144_0013

4068842 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAILED 2 time(s)
4068842 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
4068848 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
4068868 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
(CA,18)
(NY,4)
(WA,2)
(OH,2)
(UT,1)
(TX,1)
(NJ,1)
(NC,1)
(FL,1)
grunts>

```

3. Use pig to load the web log files from **clickstream/logs** using the following schema:
`reqdate:chararray, reqtime:chararray, x1:int, method:chararray, uri:chararray, x2:int ,x3:int, x4:int ,ipaddress:chararray, useragent:chararray, filter any rows which begin with a "#" (these are header rows and should be removed, then writes out the reqdate, reqtime, method, uri, ipaddress and useragent columns to a tab-delimited data set in HDFS clickstream/logs_noheader. HINT: The data is space delimited.`

Solution:

```

--view the log files
fs -ls clickstream/logs
--view the content of log file
fs -cat clickstream/logs/u_ex160211.log
--LOAD the content of log files
clickstream_logs = LOAD 'clickstream/logs/*' USING PigStorage(' ') AS
(reqdate:chararray,reqtime:chararray,X1:int,method:chararray,uri:chararray,X
2:int,X3:int,X4:int,ipaddress:chararray,useragent:chararray);
--remove the header rows with "#"

```

```

clickstream_logs_clean = FILTER clickstream_logs BY reqdate! ='#';
--prepare the final file
clickstream_logs_final = FOREACH clickstream_logs_clean GENERATE reqdate,reqtime,method,uri,ipaddress,useragent;
--store the file into HDFS
STORE clickstream_logs_final INTO
'/user/cloudera/clickstream/logs_noheader' USING PigStorage('\t');
--view the file created
fs -cat clickstream/logs_noheader/part-m-00000

```

Evidence:

```

cloudera@quickstart:~$ grunt fs -s clickstream
grunt> clickstream logs_clean = FILTER clickstream.logs_BY reqdate! ='#';
grunt> describe clickstream.logs_final = FOREACH clickstream.logs_clean GENERATE reqdate,reqtime,method,uri,ipaddress,useragent;
clickstream.logs_clean: chararray,reqtime:chararray,method:chararray,uri:chararray,ipaddress:chararray,useragent:chararray
grunt> STORE clickstream.logs_final INTO '/user/cloudera/clickstream/logs_noheader' USING PigStorage('\t');
1533701 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig Features used in the Script: FILTER
1533701 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - MR plan size before optimizations: 1
1533705 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - MR plan size after optimizations: 1
1533737 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
1533741 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - map Reduce tasks: 1
1534105 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mreded.job.reduces.marker.percent is not set, set to default 0.3
1534105 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - creating jar file Job3081162826026095241905.jar
1538797 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - jar file Job3081162826026095241905.jar created
1538801 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - pigStorage store job3081162826026095241905
1538809 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig:schematuple] is false, will not generate code
1538809 [main] INFO org.apache.pig.data.SchemaTupleFrontend - starting process to move generated code to distributed cache
1538813 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - MapReduce jobs initialized []
1538829 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
1538838 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Total input paths to process: 1
1539031 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Job3081162826026095241905 assigned to local host
1539031 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Local host IP address: 127.0.0.1
1539031 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Local host port: 5080
1539031 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Local host path: /user/cloudera/clickstream_logs_clean.clickstream_logs_final[7,25].clickstream_logs_final[0]
1539031 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - More information at: http://localhost:5080/jobdetails.jsp?jobid=job_1628829248144_0015
1539468 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - %0 completed
1539468 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - %100 complete
1539468 [main] INFO org.apache.pig.tools.pigstats.SimplePigStat - Script Statistics:
1554465 [main] INFO org.apache.pig.tools.pigstats.SimplePigStat - 
HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.6.0-cm05.7.0 0.12.0-cdh5.7.0 cloudera 2021-08-13 23:01:05 2021-08-13 23:01:25 FILTER
Success!
Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime
job_1628829248144_0015 3 3 3 3 n/a n/a n/a clickstream_logs_clean.clickstream_logs_final
40_ONLY /user/cloudera/clickstream/logs_noheader,
Inputs():
Successfully read 1167 records (321710 bytes) from: 'hdfs://quickstart.cloudera:8020/user/cloudera/clickstream/logs/*'
```

```

cloudera@quickstart:~$ grunt fs -s clickstream
Found 4 items
drwxr-xr-x . - cloudera cloudera 0 2021-08-08 17:44 clickstream/plookup
drwxr-xr-x . - cloudera cloudera 0 2021-08-13 21:19 clickstream/plookup_noheader
drwxr-xr-x . - cloudera cloudera 0 2021-08-08 17:45 clickstream/logs
drwxr-xr-x . - cloudera cloudera 0 2021-08-13 23:01 clickstream/logs_noheader
grouped by clickstream/logs_noheader
round 2 items
-rw-r--r-- 1 cloudera cloudera 0 2021-08-13 23:01 clickstream/logs_noheader..SUCCESS
-rw-r--r-- 1 cloudera cloudera 227447 2021-08-13 23:01 clickstream/logs_noheader/part-m-00000
#cat clickstream/logs_noheader/part-m-00000
#Software: Microsoft Internet Information Services
#Version: 1.0
#Date: 2010-02-11
#HTTP: s-ip CS-method CS-username C-ip
#HTTP: s-username C-ip
2010-02-11 17:16:13 GET / 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Content/jquery-ui-themes/smoothness/jquery-ui-1.10.3.custom.min.css 215.82.23.2
2010-02-11 17:16:13 GET /Themes/DefaultClean/Content/css/styleless.css 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Scripts/jquery-migrate-1.2.1.min.js 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Scripts/public-common.js 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Plugins/widgets/nivoSlider/nivoSlider.js 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Scripts/jquery-ui-1.10.3.custom.min.js 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Content/images/logo.png 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Themes/DefaultClean/Content/images/icon75px.png 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Content/images/thumbs/00000009_apparel_450.jpg 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Content/images/thumbs/00000072_25-virtual-gift-card_415.jpg 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Content/images/thumbs/00000020_build-your-own-computer_415.jpg 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Content/images/thumbs/0000024_apple-macbook-pro-13-inch_415.jpg 215.82.23.2
2011-02-09 17:16:13 GET /Themes/DefaultClean/Content/images/rating1.png 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2011-02-11 17:16:13 GET /Themes/DefaultClean/Content/images/rating1.png 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Themes/DefaultClean/Content/images/rating2.png 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Themes/DefaultClean/Content/images/shopping-bag.png 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Themes/DefaultClean/Content/images/social-spiral.png 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Themes/DefaultClean/Content/images/zooming.png 215.82.23.2 Mozilla/5.0+(Windows+NT+10.0;+WOW64;+rv:43.0)+Gecko/20100101+Firefox/43.0
2010-02-11 17:16:13 GET /Content/jquery-ui-themes/smoothness/images/ui-bg_flat_75_ffffff_40x100.png 215.82.23.2

```

4. Use hive to create two external tables for the **clickstream/logs_noheader** and **clickstream/iplookup_noheader** files you created in the previous steps. These tables should be named **weblogs** and **iplookup** respectively and should be placed in the **clickstream** database. Be sure to record all HQL steps to complete the operations.

Solution:

```
--view the log files
hdfs dfs -ls clickstream/*header

--connect to Hive
beeline -u jdbc:hive2://localhost:10000/default -u cloudera -p cloudera

--create 'clickstream' database
CREATE database clickstream;

--view the log files on HDFS from Hive
!sh hdfs dfs -ls clickstream/*header

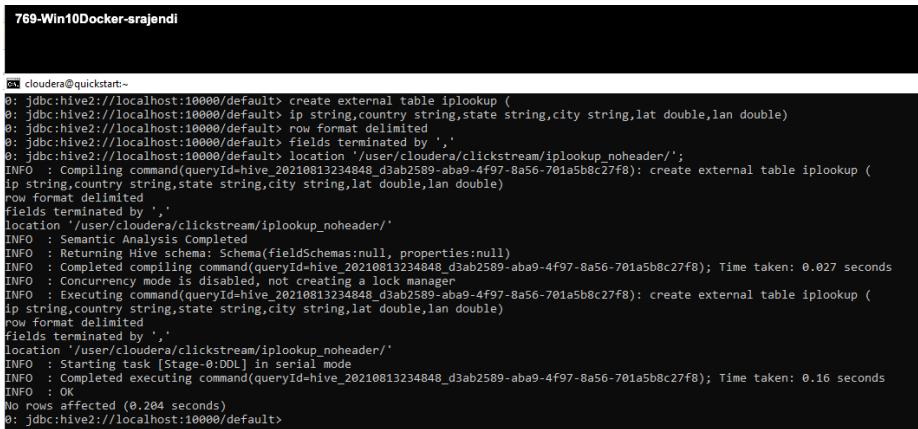
--create external table for 'iplookup'
create external table iplookup (ip string, country string, state string, city string,
lat double, lon double) row format delimited fields terminated by ',' location
'/user/cloudera/clickstream/iplookup_noheader/';

--create external table for 'weblogs'
create external table weblogs (reqdate date, reqtime timestamp, method
string, uri string, ipaddress string, useragent string) row format delimited fields
terminated by '\t' location '/user/cloudera/clickstream/logs_noheader/';

--view the hive tables schema
describe iplookup;
describe weblogs;

--view the data hive tables
select * from iplookup limit 5;
select * from weblogs limit 5;
```

Evidence:



```
769-Win10Docker-srajendi

cloudera@quickstart:~
```

```
0: jdbc:hive2://localhost:10000/default> create external table iplookup (
0: jdbc:hive2://localhost:10000/default> ip string,country string,state string,city string,lat double,lon double)
0: jdbc:hive2://localhost:10000/default> row format delimited
0: jdbc:hive2://localhost:10000/default> fields terminated by ','
0: jdbc:hive2://localhost:10000/default> location '/user/cloudera/clickstream/iplookup_noheader/';
INFO : Compiling command(queryId=hive_20210813234848_d3ab2589-ab09-4f97-8a56-701a5b8c27f8): create external table iplookup (
ip string,country string,state string,city string,lat double,lon double)
row format delimited
Fields terminated by ','
location '/user/cloudera/clickstream/iplookup_noheader/'
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fieldsSchemas:null, properties:null)
INFO : Completed compiling command(queryId=hive_20210813234848_d3ab2589-ab09-4f97-8a56-701a5b8c27f8); Time taken: 0.027 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20210813234848_d3ab2589-ab09-4f97-8a56-701a5b8c27f8): create external table iplookup (
ip string,country string,state string,city string,lat double,lon double)
row format delimited
Fields terminated by ','
location '/user/cloudera/clickstream/iplookup_noheader/'
INFO : Starting task [Stage-0:DDL] in serial mode
INFO : Completed executing command(queryId=hive_20210813234848_d3ab2589-ab09-4f97-8a56-701a5b8c27f8); Time taken: 0.16 seconds
INFO : OK
No rows affected (0.204 seconds)
0: jdbc:hive2://localhost:10000/default>
```

```

0: jdbc:hive2://localhost:10000/default> create external table weblogs (
    regdate date,reqtime timestamp,method string,url string,ipaddress string,useragent string)
    row format delimited
    fields terminated by '\t'
    location '/user/cloudera/clickstream/logs_noheader/'
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fields=null, properties=null)
INFO : Compiling command(queryId=hive_20210813235353_a49da797-ee6b-40d8-96eb-41a9ed56b1a7); Time taken: 0.027 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20210813235353_a49da797-ee6b-40d8-96eb-41a9ed56b1a7); create external table weblogs (
    regdate date,reqtime timestamp,method string,url string,ipaddress string,useragent string)
    row format delimited
    fields terminated by '\t'
    location '/user/cloudera/clickstream/logs_noheader/'
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fields=null, properties=null)
INFO : Compiling command(queryId=hive_20210813235353_a49da797-ee6b-40d8-96eb-41a9ed56b1a7); Time taken: 0.027 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20210813235353_a49da797-ee6b-40d8-96eb-41a9ed56b1a7); create external table weblogs (
    regdate date,reqtime timestamp,method string,url string,ipaddress string,useragent string)
    row format delimited
    fields terminated by '\t'
    location '/user/cloudera/clickstream/logs_noheader/'
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fields=null, properties=null)
INFO : Compiling command(queryId=hive_20210813235353_a49da797-ee6b-40d8-96eb-41a9ed56b1a7); Time taken: 0.027 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20210813235353_a49da797-ee6b-40d8-96eb-41a9ed56b1a7); Time taken: 0.027 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Starting task [Stage-0] in serial mode
INFO : Completed executing command(queryId=hive_20210813235353_a49da797-ee6b-40d8-96eb-41a9ed56b1a7); Time taken: 0.061 seconds
INFO : OK
No rows affected (0.187 seconds)
0: jdbc:hive2://localhost:10000/default>

```

```

cloudera@quickstart:-
0: jdbc:hive2://localhost:10000/default> describe iplookup;
+-----+-----+-----+
| col_name | data_type | comment |
+-----+-----+-----+
| ip       | string   |          |
| country  | string   |          |
| city     | string   |          |
| city     | string   |          |
| lat      | double   |          |
+-----+-----+-----+
0 rows selected (0.144 seconds)

0: jdbc:hive2://localhost:10000/default> describe weblogs;
+-----+-----+-----+
| col_name | data_type | comment |
+-----+-----+-----+
| regdate | date    |          |
| reqtime | timestamp|          |
| method  | string   |          |
| method  | string   |          |
| ipaddress| string   |          |
| useragent| string   |          |
+-----+-----+-----+
0 rows selected (0.148 seconds)

0: jdbc:hive2://localhost:10000/default>

```

```

769-Win10Docker-srajendi
Enforce US Keyboard Layout View Fullscreen Send Ctrl+Alt+Delete
cloudera@quickstart:-
0: jdbc:hive2://localhost:10000/default> select * from iplookup limit 5;
INFO : Compiling Command(queryId=hive_20210813235454_8193abaa_fcc4-4a93-b5e4-b246f470e948); select * from iplookup limit 5
INFO : Semantic Analysis Completed
INFO : Executing command(queryId=hive_20210813235454_8193abaa_fcc4-4a93-b5e4-b246f470e948); select * from iplookup limit 5
INFO : FieldsSchema:[FieldSchema(name:iplookup.ip, type:string, comment:null), FieldSchema(name:iplookup.state, type:string, comment:null), FieldSchema(name:iplookup.city, type:string, comment:null), FieldSchema(name:iplookup.lat, type:double, comment:null), FieldSchema(name:iplookup.lng, type:double, comment:null)], properties=null
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20210813235454_8193abaa_fcc4-4a93-b5e4-b246f470e948); select * from iplookup limit 5
INFO : Completed executing command(queryId=hive_20210813235454_8193abaa_fcc4-4a93-b5e4-b246f470e948); Time taken: 0.082 seconds
INFO : OK
+-----+-----+-----+-----+-----+
| iplookup.ip | iplookup.state | iplookup.city | iplookup.lat | iplookup.lng |
+-----+-----+-----+-----+-----+
| 128.122.160.238 | USA | NY | New York | 40.72784 | -74.000041 |
| 128.238.232.188 | USA | NY | Syracuse | 43.04812 | -76.147424 |
| 128.238.232.189 | USA | NY | Albany | 43.06807 | -74.88007 |
| 172.189.252.8 | USA | VA | Dulles | 38.95585 | -77.447819 |
| 255.82.23.3 | USA | OH | Columbus | 39.961176 | -82.998794 |
+-----+-----+-----+-----+-----+
5 rows selected (0.482 seconds)

0: jdbc:hive2://localhost:10000/default> select * from weblogs limit 5;
INFO : Compiling Command(queryId=hive_20210813235454_8193abaa_fcc4-4a93-b5e4-b246f470e948); select * from weblogs limit 5
INFO : Semantic Analysis Completed
INFO : Executing command(queryId=hive_20210813235454_8193abaa_fcc4-4a93-b5e4-b246f470e948); select * from weblogs limit 5
INFO : FieldsSchema:[FieldSchema(name:weblogs.regdate, type:date, comment:null), FieldSchema(name:weblogs.reqtime, type:timestamp, comment:null), FieldSchema(name:weblogs.method, type:string, comment:null), FieldSchema(name:weblogs.url, type:string, comment:null), FieldSchema(name:weblogs.ipaddress, type:string, comment:null), FieldSchema(name:weblogs.useragent, type:string, comment:null)], properties=null
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId=hive_20210813235454_8193abaa_fcc4-4a93-b5e4-b246f470e948); select * from weblogs limit 5
INFO : Completed executing command(queryId=hive_20210813235454_8193abaa_fcc4-4a93-b5e4-b246f470e948); Time taken: 0.081 seconds
INFO : OK
+-----+-----+-----+-----+-----+-----+
| weblogs.regdate | weblogs.reqtime | weblogs.method | weblogs.url | weblogs.ipaddress | weblogs.useragent |
+-----+-----+-----+-----+-----+-----+
| NULL | NULL | Information | Services | / | Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/81.0.4044.138 Safari/537.36 |
| NULL | NULL | NULL | NULL | / | Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/81.0.4044.138 Safari/537.36 |
| NULL | NULL | s-IP | cs-method | cs-username | c-ip |
| 2005-09-11 | NULL | GET | / | 255.82.23.3 | Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/81.0.4044.138 Safari/537.36 |
+-----+-----+-----+-----+-----+-----+
5 rows selected (0.14 seconds)

0: jdbc:hive2://localhost:10000/default>

```

5. Write an HQL query to display the name of the city and the number of HTTP requests from that city (NOTE: each row in the web logs is an HTTP request). Order the output so cities with the most requests are at the top. If you complete the query correctly, you should see Syracuse has 272-page requests and Los Angeles has 24.

Solution:

```

--HQL to select cities with most ip requests
select i.city, count(w.method) as http_requests
from iplookup as i join weblogs as w
on i.ip = w.ipaddress
group by i.city
order by http_requests desc;

```

Evidence:

```
769-Win10Docker-srajendi

cloudera@quickstart:~
```

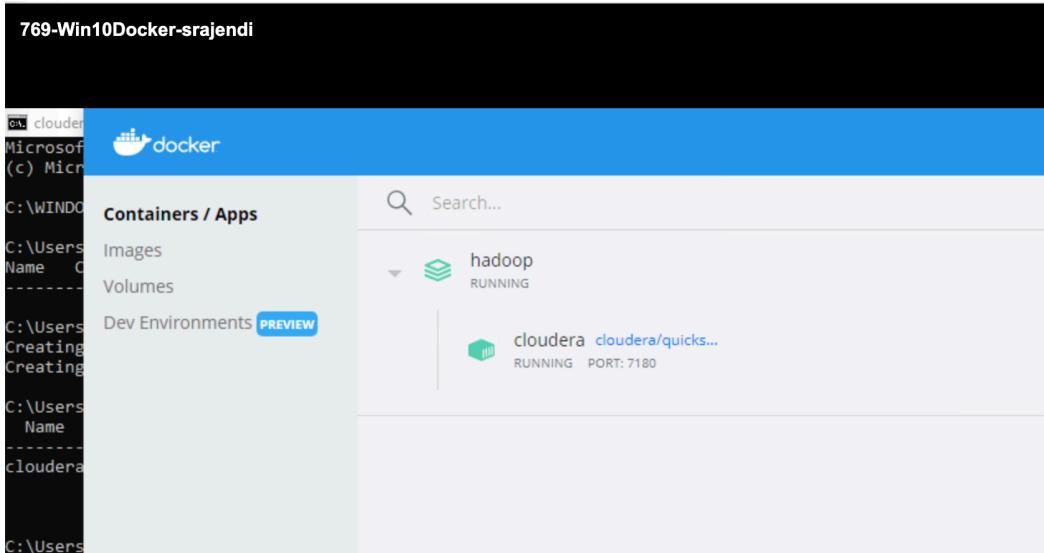
0: jdbc:hive2://localhost:10000/default> select
0: jdbc:hive2://localhost:10000/default> i.city,count(w.method) as http_requests
0: jdbc:hive2://localhost:10000/default> from iplookup as i join weblogs as w
0: jdbc:hive2://localhost:10000/default> on i.ip = w.ipaddress
0: jdbc:hive2://localhost:10000/default> group by i.city
0: jdbc:hive2://localhost:10000/default> order by http_requests desc;

```
INFO : MapReduce Jobs Launched:  
INFO : Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 3.62 sec HDFS Read: 237563 HDFS Write: 456 SUCCESS  
INFO : Stage-Stage-3: Map: 1 Reduce: 1 Cumulative CPU: 2.45 sec HDFS Read: 5093 HDFS Write: 164 SUCCESS  
INFO : Total MapReduce CPU Time Spent: 6 seconds 70 msec  
INFO : Completed executing command(queryId=hive_20210814000404_239f1460-9b8f-4cbc-a59e-337aa1ebd637); Time taken: 47.451 seconds  
INFO : OK  
+-----+-----+  
| i.city | http_requests |  
+-----+-----+  
| Syracuse | 272 |  
| Columbus | 152 |  
| Dulles | 91 |  
| Jersey City | 86 |  
| Dallas | 85 |  
| Freeport | 78 |  
| Arlington | 72 |  
| New York | 68 |  
| Raleigh | 54 |  
| Tampa | 54 |  
| Salt Lake City | 51 |  
| Cleveland | 35 |  
| Los Angeles | 24 |  
+-----+-----+  
13 rows selected (47.911 seconds)  
0: jdbc:hive2://localhost:10000/default>
```

Appendix

769-Win10Docker-srajendi

```
cloudera@quickstart:~  
Microsoft Windows [Version 10.0.19041.1110]  
(c) Microsoft Corporation. All rights reserved.  
C:\WINDOWS\system32>cd C:\Users\LocalAdmin\srajendi\adv-db-labs\hadoop  
C:\Users\LocalAdmin\srajendi\adv-db-labs\hadoop>docker-compose ps  
Name      Command     State      Ports  
-----  
cloudera /usr/bin/docker-quickstart   Up      0.0.0.0:7180->7180/tcp,:::7180/tcp,  
          0.0.0.0:8080->80/tcp,:::8080->80/tcp,  
          0.0.0.0:8888->8888/tcp,:::8888/tcp  
C:\Users\LocalAdmin\srajendi\adv-db-labs\hadoop>docker-compose up -d  
Creating network "hadoop_default" with the default driver  
Creating cloudera ... done  
C:\Users\LocalAdmin\srajendi\adv-db-labs\hadoop>docker-compose ps  
Name      Command     State      Ports  
-----  
cloudera /usr/bin/docker-quickstart   Up      0.0.0.0:7180->7180/tcp,:::7180/tcp,  
          0.0.0.0:8080->80/tcp,:::8080->80/tcp,  
          0.0.0.0:8888->8888/tcp,:::8888/tcp  
C:\Users\LocalAdmin\srajendi\adv-db-labs\hadoop>docker-compose exec cloudera bash -c "su -l cloudera"  
[cloudera@quickstart ~]$ hdfs dfs -ls  
Found 12 items  
drwxr-xr-x - cloudera cloudera 0 2021-08-08 17:42 clickstream  
drwxr-xr-x - cloudera cloudera 0 2021-08-10 00:32 fudgemart-clothing  
drwxr-xr-x - cloudera cloudera 0 2021-08-10 18:57 fudgemart-clothing_1  
drwxr-xr-x - cloudera cloudera 0 2021-08-10 18:58 fudgemart-clothing_2  
drwxr-xr-x - cloudera cloudera 0 2021-08-10 19:03 fudgemart-clothing_4  
drwxr-xr-x - cloudera cloudera 0 2021-08-10 19:09 fudgemart-clothing_5  
drwxr-xr-x - cloudera cloudera 0 2021-08-10 19:14 fudgemart-clothing_6  
drwxr-xr-x - cloudera cloudera 0 2021-08-10 18:48 fudgemart-clothing_new  
drwxr-xr-x - cloudera cloudera 0 2021-08-10 00:23 fudgemart-products-by-clothing  
drwxr-xr-x - cloudera cloudera 0 2021-08-09 23:40 sotu2016  
drwxr-xr-x - cloudera cloudera 0 2021-08-06 06:34 text  
drwxr-xr-x - cloudera cloudera 0 2021-08-10 00:47 tweets  
[cloudera@quickstart ~]$
```



769-Win10Docker-srajendi

```
[cloudera@quickstart:~]$ pig
log4j:WARN No appenders could be found for logger (org.apache.hadoop.util.Shell).
log4j:WARN Please initialize the log4j system properly.
log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for more info.
2021-08-13 00:21:12,196 [main] INFO org.apache.pig.Main - Apache Pig version 0.12.0-cdh5.7.0 (reported) compiled Mar 23 2016, 11:34:31
2021-08-13 00:21:12,196 [main] INFO org.apache.pig.Main - Logging error messages to: /home/cloudera/pig_1628814071168.log
2021-08-13 00:21:12,228 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file /home/cloudera/.pigbootup not found
2021-08-13 00:21:12,834 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapped.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-13 00:21:12,835 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-13 00:21:12,836 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://quickstart.cloudera:8020
2021-08-13 00:21:14,014 [main] INFO org.apache.pig.Main - Configuration deprecation - mapreduce.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-13 00:21:14,009 [main] INFO org.apache.pig.Main - Configuration deprecation - mapreduce.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-13 00:21:14,017 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-13 00:21:14,073 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-13 00:21:14,078 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapped.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-13 00:21:14,081 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-13 00:21:14,137 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapped.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-13 00:21:14,192 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-13 00:21:14,194 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapped.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-13 00:21:14,260 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-13 00:21:14,262 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapreduce.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-13 00:21:14,301 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-13 00:21:14,308 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapped.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-13 00:21:14,371 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-13 00:21:14,380 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapped.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-08-13 00:21:14,461 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-13 00:21:14,565 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-08-13 00:21:14,571 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapped.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
grunt>
```

769-Win10Docker-srajendi

```
grunt> quit
[cloudera@quickstart:~]$ echo "log4j.rootLogger=fatal" > noglog.conf
[cloudera@quickstart:~]$ pig -4 noglog.conf
log4j:WARN No appenders could be found for logger (org.apache.hadoop.util.Shell).
log4j:WARN Please initialize the log4j system properly.
log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for more info.
64 [main] INFO org.apache.pig.Main - Apache Pig version 0.12.0-cdh5.7.0 (reported) compiled Mar 23 2016, 11:34:31
66 [main] INFO org.apache.pig.Main - Logging error messages to: /home/cloudera/pig_1628814117998.log
69 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file /home/cloudera/.pigbootup not found
866 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://quickstart.cloudera:8020
1772 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to map-reduce job tracker at: localhost:8021
grunt>
```

769-Win10Docker-srajendi

```
cloudera@quickstart:~
```

```
grunt> fs -ls
Found 12 items
drwxr-xr-x  - cloudera cloudera          0 2021-08-08 17:42 clickstream
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:32 fudgemart-clothing
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 18:57 fudgemart-clothing_1
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 18:58 fudgemart-clothing_2
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 19:03 fudgemart-clothing_4
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 19:09 fudgemart-clothing_5
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 19:14 fudgemart-clothing_6
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 18:48 fudgemart-clothing_new
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:23 fudgemart-products-by-clothing
drwxr-xr-x  - cloudera cloudera          0 2021-08-09 23:40 sotu2016
drwxr-xr-x  - cloudera cloudera          0 2021-08-06 06:34 text
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:47 tweets
grunt> fs -mkdir testing
grunt> fs -ls testing
grunt> -
```

769-Win10Docker-srajendi

```
Enforce US Keyboard Layout | View Fullscreen
```

```
cloudera@quickstart:~$ echo "log4j.rootLogger=fatal" > noglog.conf
[cloudera@quickstart:~]$ pig -4 noglog.conf
log4j:WARN No appenders could be found for logger (org.apache.hadoop.util.Shell).
log4j:WARN Please initialize the log4j system properly.
log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for more info.
52 [main] INFO org.apache.pig.Main - Apache Pig version 0.12.0-cdh5.7.0 (reported) compiled Mar 23 2016, 11:34:31
54 [main] INFO org.apache.pig.Main - Logging error messages to: /home/cloudera/pig_1628813987388.log
89 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file /home/cloudera/.pigbootup not found
591 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://quickstart.cloudera:8020
1888 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to map-reduce job tracker at: localhost:8021
grunts> fs -ls
Found 12 items
drwxr-xr-x  - cloudera cloudera          0 2021-08-08 17:42 clickstream
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:32 fudgemart-clothing
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 18:57 fudgemart_clothing_1
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 18:58 fudgemart_clothing_2
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 19:03 fudgemart_clothing_4
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 19:09 fudgemart_clothing_5
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 19:14 fudgemart_clothing_6
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 18:48 fudgemart_clothing_new
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:23 fudgemart_products-by-clothing
drwxr-xr-x  - cloudera cloudera          0 2021-08-09 23:40 sotu2016
drwxr-xr-x  - cloudera cloudera          0 2021-08-06 06:34 text
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:47 tweets
grunt> fs -mkdir testing
grunt> fs -ls testing
ls: Unknown command
Did you mean -ls? This command begins with a dash.
grunt> fs -ls testing
grunt> -
```

```

cloudera@quickstart:~$ clear
[cloudera@quickstart ~]$ hdfs dfs -ls
Found 7 items
drwxr-xr-x  - cloudera cloudera          0 2021-08-08 17:42 clickstream
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:32 fudgemart-clothing
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:23 fudgemart-products-by-clothing
drwxr-xr-x  - cloudera cloudera          0 2021-08-09 23:40 sotu2016
drwxr-xr-x  - cloudera cloudera          0 2021-08-13 00:23 testing
drwxr-xr-x  - cloudera cloudera          0 2021-08-06 06:34 text
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:47 tweets
[cloudera@quickstart ~]$ hdfs dfs -ls clickstream/iplookup
Found 1 items
-rw-r--r--  1 cloudera cloudera     1251 2021-08-08 17:44 clickstream/iplookup/ip_lookup.csv
[cloudera@quickstart ~]$ hdfs dfs -cat clickstream/iplookup/ip_lookup.csv
IP,Country,State,City,ApproxLat,ApproxLng
172.189.252.8,USA,VA,Dulles,38.955855,-77.447819
215.82.23.2,USA,OH,Columbus,39.961176,-82.998794
98.29.25.44,USA,OH,Cleveland,41.49932,-81.694361
68.199.40.156,USA,NY,Freeport,40.657602,-73.583184
155.100.169.152,USA,UT,Salt Lake City,40.760779,-111.891047
38.68.15.223,USA,TX,Dallas,32.776664,-96.796988
70.209.14.54,USA,FL,Tampa,27.950575,-82.457178
74.111.6.173,USA,VA,Arlington,38.87997,-77.10677
128.230.122.180,USA,NY,Syracuse,43.048122,-76.147424

```

```

769-Win10Docker-srajend
[cloudera@quickstart:~$ echo "log4j.rootLogger=fatal" > nolog.conf
[cloudera@quickstart ~]$ pig -A nolog.conf
log4j:WARN No appenders could be found for logger (org.apache.hadoop.util.Shell).
log4j:WARN Please initialize the log4j system properly.
log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for more info.
85 [main] INFO org.apache.pig.Main - Apache Pig version 0.12.0-cdh5.7.0 (reexported) compiled Mar 23 2016, 11:34:31
87 [main] INFO org.apache.pig.Main - Logging error messages to: /home/cloudera/pig_1628888256766.log
112 [main] INFO org.apache.pig.impl.util.Utils - Default bootup file /home/cloudera/.pigbootup not found
672 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://quickstart.cloudera:8020
1746 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to map-reduce job tracker at: localhost:8021
grunt> fs -ls
Found 7 items
drwxr-xr-x  - cloudera cloudera          0 2021-08-08 17:42 clickstream
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:32 fudgemart-clothing
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:23 fudgemart-products-by-clothing
drwxr-xr-x  - cloudera cloudera          0 2021-08-09 23:40 sotu2016
drwxr-xr-x  - cloudera cloudera          0 2021-08-13 00:23 testing
drwxr-xr-x  - cloudera cloudera          0 2021-08-06 06:34 text
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:47 tweets
grunt> fs -ls clickstream/iplookup
Found 2 items
drwxr-xr-x  - cloudera cloudera          0 2021-08-08 17:44 clickstream/iplookup
drwxr-xr-x  - cloudera cloudera          0 2021-08-08 17:45 clickstream/logs
grunt> fs -ls clickstream/iplookup
Found 1 items
-rw-r--r--  1 cloudera cloudera     1251 2021-08-08 17:44 clickstream/iplookup/ip_lookup.csv
grunt> fs -cat clickstream/iplookup/ip_lookup.csv
IP,Country,State,City,ApproxLat,ApproxLng
172.189.252.8,USA,VA,Dulles,38.955855,-77.447819
215.82.23.2,USA,OH,Columbus,39.961176,-82.998794
98.29.25.44,USA,OH,Cleveland,41.49932,-81.694361
68.199.40.156,USA,NY,Freeport,40.657602,-73.583184
155.100.169.152,USA,UT,Salt Lake City,40.760779,-111.891047
38.68.15.223,USA,TX,Dallas,32.776664,-96.796988
70.209.14.54,USA,FL,Tampa,27.950575,-82.457178
74.111.6.173,USA,VA,Arlington,38.87997,-77.10677
128.230.122.180,USA,NY,Syracuse,43.048122,-76.147424
128.222.140.238,USA,NY,New York,40.712784,-74.005941
56.216.127.219,USA,NC,Raleigh,35.77959,-78.638197

```

```

cloudera@quickstart:~$ clear
grunt> iplookup.File = LOAD 'clickstream/iplookup/ip_lookup.csv' USING PigStorage(',');
>> iplookup.File: org:pchararray,country:chararray,state:chararray,city:chararray,lat:double,lan:double;
java.lang.Error: ERROR org.apache.pig.tools.grunt.Grunt: ERROR 1200: <line 6, column 6> : mismatched input 'iplookup.File' expecting SEMI_COLON
details at: logFile: /tmp/cloudera/pig_1628888259659.log
grunt> clear

grunt> iplookup.File = LOAD 'clickstream/iplookup/ip_lookup.csv' USING PigStorage(',');
>> AS (pchararray,country,chararray,state:chararray,city:chararray,lat:double,lan:double);
java.lang.Error: ERROR org.apache.pig.tools.grunt.Grunt: ERROR 1200: <line 6, column 6> : mismatched input 'AS' expecting SEMI_COLON
details at: logFile: /tmp/cloudera/pig_1628888259659.log
grunt> iplookup.File.NoHeader_sorted = ORDER iplookup.File_by_ip ASC;
grunt> STORE iplookup.File.NoHeader_sorted INTO '/user/cloudera/clickstream/iplookup_noheader' USING PigStorage('');
1228807 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer  - [RULES_ENABLED=[AddOrForEach, ColumnMapKeyPrune, DuplicateForEachColumnRewrite, GroupByConstParallelSetter, ImplicitSplitInserter, LimitOptimizer, LoadySpecInserter, MergeFilter, MergeGroupBy, Partitioner, PushdownFilter, PushdownGroupBy, PushdownTypeFilter], RULES_DISABLED=[FilterLogicExpressionSimplifier, PartitionFilterOptimizer]}
1228837 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler  - File contains too many partitions: threshold=100 (limits=1) False
1228837 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler  - MR plan size before optimization: 3
1228848 [main] INFO org.apache.pig.tools.pigscript.parser.PigScriptParser  - Pig script settings are added to the job
1228897 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRControlCompiler  - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
1234552 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRControlCompiler  - jar file:job7232152999980217632.jar created
1234573 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRControlCompiler  - Setting up single store job
1234583 [main] INFO org.apache.pig.data.SchemaTupleFrontend  - Starting process to move generated code to distributed cache
1234583 [main] INFO org.apache.pig.data.SchemaTupleFrontend  - Setting key [pig.schematuple.classes] with classes to deserialize []

```

```

cloudera@quickstart:~$ cloudera@quickstart:~$ 
1288280 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
1288442 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
1288442 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total output paths (calculated) to process : 1
1288787 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Hadoop job: job_1628829248144_0003
1288787 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases iplookup_File_Noheader_sorted
1288787 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - detailed locations: M: iplookup_File_Noheader_sorted[8,32] C: R:
1288510 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 83% complete
1288510 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
1288444 [main] INFO org.apache.pig.tools.pigstats.ScriptStatistics - Script Statistics:
HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.6.0-cdh5.7.0 0.12.0-cdh5.7.0 cloudera 2021-08-13 21:18:05 2021-08-13 21:19:25 ORDER_BY,FILTER
Success!
Job Stats (time in seconds):
JobId Map Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime
job_1628829248144_0001 1 0 3 3 3 n/a n/a n/a iplookup_File_iplookup_File_Noheader
AP_ONLY
job_1628829248144_0002 1 1 3 3 3 3 3 3 iplookup_File_Noheader_sorted
AMPLER
job_1628829248144_0003 1 1 4 4 4 4 3 3 3 iplookup_File_Noheader_sorted
ORDER_BY /user/cloudera/clickstream/iplookup_noheader,
Input(s):
Successfully read 24 records (1650 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/clickstream/iplookup/ip_lookup.csv"
Output(s):
Successfully stored 23 records (1185 bytes) in: "/user/cloudera/clickstream/iplookup_noheader"

```

769-Win10Docker-srajendi

```

cloudera@quickstart:~$ 
Successfully read 24 records (1650 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/clickstream/iplookup/ip_lookup.csv"
Output(s):
Successfully stored 23 records (1185 bytes) in: "/user/cloudera/clickstream/iplookup_noheader"

Counters:
Total records written : 23
Total bytes written : 1185
Spillable Memory Manager spill count : 0
total bags proactively spilled: 0
total records proactively spilled: 0

Job DAG:
job_1628829248144_0001 -> job_1628829248144_0002,
job_1628829248144_0002 -> job_1628829248144_0003,
job_1628829248144_0003

1317161 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAILED 2 time(s).
1317161 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
grunts> fs -ls clickstream
grunts> 3 items
drwxr-xr-x - cloudera cloudera 0 2021-08-08 17:44 clickstream/iplookup
drwxr-xr-x - cloudera cloudera 0 2021-08-13 21:19 clickstream/iplookup_noheader
drwxr-xr-x - cloudera cloudera 0 2021-08-08 17:45 clickstream/logs

```

769-Win10Docker-srajendi

```

cloudera@quickstart:~$ grunt> fs -ls clickstream/iplookup_noheader
Found 2 items
-rw-r--r-- 1 cloudera cloudera 0 2021-08-13 21:19 clickstream/iplookup_noheader/_SUCCESS
-rw-r--r-- 1 cloudera cloudera 1185 2021-08-13 21:19 clickstream/iplookup_noheader/part-r-00000
grunts> fs -cat clickstream/iplookup_noheader/part-r-00000
128.122.140.238,USA,NY,New York,40.712784,-74.005941
128.230.122.180,USA,NY,Syracuse,43.048122,-76.147424
155.100.169.152,USA,UT,Salt Lake City,40.768779,-111.891047
172.189.252.8,USA,VA,Dulles,38.955855,-77.447819
215.82.23.2,USA,OH,Columbus,39.961176,-82.998794
38.68.15.223,USA,TX,Dallas,32.776664,-96.796988
54.114.107.209,USA,NJ,Jersey City,40.728157,-74.077642
56.216.127.219,USA,NC,Raleigh,35.77959,-78.638179
68.199.40.156,USA,NY,Freeport,40.657602,-73.583184
70.209.14.54,USA,FL,Tampa,27.950575,-82.457178
74.111.18.59,USA,NY,Syracuse,43.048122,-76.147424
74.111.6.173,USA,VA,Arlington,38.87997,-77.10677
8.37.70.112,USA,CA,Los Angeles,34.052234,-118.243685
8.37.70.170,USA,CA,Los Angeles,34.052234,-118.243685
8.37.70.226,USA,CA,Los Angeles,34.052234,-118.243685
8.37.70.77,USA,CA,Los Angeles,34.052234,-118.243685
8.37.70.99,USA,CA,Los Angeles,34.052234,-118.243685
8.37.71.25,USA,CA,Los Angeles,34.052234,-118.243685
8.37.71.43,USA,CA,Los Angeles,34.052234,-118.243685
8.37.71.57,USA,CA,Los Angeles,34.052234,-118.243685
8.37.71.69,USA,CA,Los Angeles,34.052234,-118.243685
8.37.71.9,USA,CA,Los Angeles,34.052234,-118.243685
98.29.25.44,USA,OH,Cleveland,41.49932,-81.694361
grunts> 

```

```

$ cloudera@quickstart:
$ ./bin/hadoop jar $HADOOP_HOME/share/hadoop/tools/lib/hadoop-mapreduce-client-tools.jar org.apache.hadoop.mapred.lib.MrJobConfigTool -D mapreduce.job.reduces=1 /tmp/quickstart/cloudera/clickstream/iplookup/ip_lookup.csv /tmp/quickstart/cloudera/clickstream/iplookup/ip_lookup_by_state
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Job ID: job_1628829248144_0007
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Processing aliases iplookup_by_state
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Configuration options: M: iplookup@state[12,20] C: R:
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: More information at: http://localhost:50030/jobdetails.jsp?jobid=job_1628829248144_0007
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: 100% complete
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Success!
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Job Stats (time in seconds):
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: medianReductionTime Alias Feature Outputs
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: job_1628829248144_0004 1 0 3 3 3 3 n/a n/a n/a n/a iplookup_File,iplookup_File_Noheader
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: AP_ONLY
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: job_1628829248144_0005 1 1 3 3 3 3 3 3 3 iplookup_File_Noheader_sorted
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: _REDUCER
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: job_1628829248144_0006 1 1 3 3 3 3 3 3 3 iplookup_File_Noheader_sorted
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: REDER_BY
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: job_1628829248144_0007 1 1 3 3 3 3 3 3 3 iplookup_by_state
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: GROUP_BY hdfs://quickstart.cloudera:8020/tmp/temp-196607337/tmp-1545452863,
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Input(s):
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Successfully read 24 records (1650 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/clickstream/iplookup/ip_lookup.csv"
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Output(s):
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Successfully stored 9 records (1415 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-196607337/tmp-1545452863"
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Counters:
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Total records written : 9
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Total bytes written : 1415
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Spillable Memory Manager spill count : 0
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Total bags proactively spilled: 0
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: total records proactively spilled: 0
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: Job DAG:
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: job_1628829248144_0004 -> job_1628829248144_0005,
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: job_1628829248144_0005 -> job_1628829248144_0006,
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: job_1628829248144_0006 -> job_1628829248144_0007,
2021-08-13 21:45:44,462 INFO org.apache.hadoop.mapred.lib.MrJobConfigTool: job_1628829248144_0007

```

```
789-Win10Docker-arajend| Enforce US Keyboard Layout | View Fullscreen | Send Ctrl+Alt+Delete

cloudera@quickstart:~
```

```
cloudera@quickstart:~  
  
grunt describe iplookup_by_state;  
iplookup_by_state: {group: chararray,iplookup_File_Noheader_sorted: {{ip: chararray,country: chararray,state: chararray,city: chararray,lat: double,lon: double}}}  
grunt iplookup_by_state_n_counts = FOREACH iplookup_by_state GENERATE group as state,COUNT(iplookup_File_Noheader_sorted.ip) as counts;  
grunt iplookup_by_state_n_counts_sorted = ORDER iplookup_by_state_n_counts by counts DESC;  
grunt>
```

```

769-Win10Docker-srajendi

cloudera@quickstart:~$ grunt> describe iplookup_by_state;
iplookup_by_state: {group: chararray,iplookup_File_Noheader_sorted: {((ip: chararray,country: chararray,state: chararray,city: chararray,lat: double,lon: double))}}
group> iplookup_by_state_n_counts = FOREACH iplookup_by_state GENERATE group as state,COUNT(iplookup_File_Noheader_sorted.ip) as counts;
group> iplookup_by_state_n_counts_sorted = ORDER iplookup_by_state_n_counts by counts DESC;
group> DUMP iplookup_by_state_n_counts_sorted;
3899313 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP_BY_ORDER_BY,FILTER
3899315 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - (RULES_ENABLED-[AddForEach, ColumnMapKeyPrune, DuplicateForEachColumnRewrite, GroupByConstProp, LoadTypeCastInserter, MergeForEach, NewPartitionOptimizer, PushDownForEachFlatIter, PushFilter, SplitFilter, StreamTypeCastInserter], RULES_DISABLED-[])
3899326 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCCompiler - File concatenation threshold: 100 optimistic? false
3899344 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.CombinerOptimizer - Choosing to move algebraic foreach to combiner
3899350 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 6
3899351 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 6
3899383 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
3899384 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Job has been submitted. Job reduce.mapReducers.buffForPercent is not set, set to default 0.3
3899674 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCCompiler - creating jar file Job714910676869897826.jar created
38994466 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - jar file Job714910676869897826.jar created
38994481 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
38994481 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
38994481 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
38994481 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - MR plan size after optimization: 6
38994504 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
38994548 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
38994553 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths (combined) to process : 1
38995009 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - HadoopJobId: job_1628829248144_0008
38995009 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases iplookup_File.iplookup_File_Noheader
38995009 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - detailed locations: M: iplookup_File[4,16],iplookup_File[-1,-1]
38995009 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - More information at: http://localhost:50030/jobdetails.jsp?jobid=job_1628829248144_0008
38995058 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - % complete
3916537 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 8% complete
3920213 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 16% complete
3929305 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
3929376 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mpared: job.reduce.mapReducers.buffForPercent is not set, set to default 0.3
3929376 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reducers phase detected, estimating # of required reducers
3929378 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Using reducer estimator: org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher
3929378 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - BytesPerReducer=1000000000 maxReducers=999 totalInputFileSize=1385
3929378 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting Parallelism to 1
3929378 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - creating jar file Job70204232360865629210.jar created
3929378 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - jar file Job70204232360865629210.jar created
3929378 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Single file store job
3929378 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
3929378 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
3929378 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserialize []
3929378 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
3929472 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
3929472 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths (combined) to process : 1
3929530 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - HadoopJobId: job_1628829248144_0009
3929530 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases iplookup_File_Noheader_sorted

```

```

cloudera@quickstart:~$ 
4028983 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - HadoopJobId: job_1628829248144_0013
4028983 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases iplookup_by_state_n_counts_sorted
4028983 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - detailed locations: M: iplookup_by_state_n_counts_sorted[18,30] C: R:
4028983 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - More information at: http://localhost:50030/jobdetails.jsp?jobid=job_1628829248144_0013
4049130 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 91% complete
4049125 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
4049161 [main] INFO org.apache.pig.tools.pigstats.SimplePigStat - Script Statistics:

Hadoop Version PigVersion UserId StartedAt FinishedAt Features
2.6.0-cdh5.7.0 0.12.0-cdhs5.7.0 cloudera 2021-08-13 22:02:36 2021-08-13 22:05:05 GROUP_BY,ORDER_BY,FILTER

Success!
Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime
job_1628829248144_0008 1 0 3 3 3 n/a n/a n/a iplookup_File.iplookup_File_Noheader
job_1628829248144_0013 SAMPLER
job_1628829248144_0010 1 1 3 3 3 3 3 3 3 iplookup_File_Noheader_sorted
job_1628829248144_0011 1 1 3 3 3 3 3 3 3 iplookup_by_state.iplookup_by_state_n_counts
job_1628829248144_0012 1 1 3 3 3 3 3 3 3 iplookup_by_state_n_counts_sorted
job_1628829248144_0013 1 1 3 3 3 3 3 3 3 iplookup_by_state_n_counts_sorted
ORDER_BY hdfs://quickstart.cloudera:8020/tmp/temp-196607337/tmp1642888536,
Input(s):
Successfully read 24 records (1656 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/clickstream/iplookup/ip_lookup.csv"
Output(s):
Successfully stored 9 records (94 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-196607337/tmp1642888536"

Counters:
Total Records written : 9
Total bytes written : 94
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

```

```

769-Win10Docker-srajendi

cloudera@quickstart:~$ 
Job DAG:
job_1628829248144_0008 --> job_1628829248144_0009
job_1628829248144_0009 --> job_1628829248144_0010
job_1628829248144_0010 --> job_1628829248144_0011
job_1628829248144_0011 --> job_1628829248144_0012
job_1628829248144_0012 --> job_1628829248144_0013
job_1628829248144_0013

4068842 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAILED 2 time(s)
4068842 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
4068848 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Key [pig.schematuple] was not set... will not generate code.
4068860 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
(CA,18)
(NY,4)
(VA,2)
(OH,2)
(UT,1)
(TX,1)
(WI,1)
(NG,1)
(FL,1)
grunt>
```

```

769-Win10Docker-srajendi

cloudera@quickstart:~$ ls -ls
total 2 items
drwxr-xr-x  - cloudera cloudera          0 2021-08-13 21:19 clickstream
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:32 fudgemart-clothing
cloudera@quickstart:~$ pig -4 nolog.conf
log4j:WARN Please initialize the log4j system properly.
log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for more info.
97 [main] INFO org.apache.pig.Main - Apache Pig version 0.12.0-cdh5.7.0 (reported) compiled Mar 23 2016, 11:34:31
98 [main] INFO org.apache.pig.Main - Logging error messages to: /home/cloudera/pig_1628894131487.log
133 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file /home/cloudera/.pigbootstrap not found
760 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://quickstart.cloudera:8020
1859 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to map-reduce job tracker at: localhost:8021
grunt> fs -ls
Found 2 items
drwxr-xr-x  - cloudera cloudera          0 2021-08-13 21:19 clickstream
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:32 fudgemart-clothing
cloudera@quickstart:~$ grunt> fs -ls clickstream
Found 3 items
drwxr-xr-x  - cloudera cloudera          0 2021-08-08 17:44 clickstream/plookup
drwxr-xr-x  - cloudera cloudera          0 2021-08-13 21:19 clickstream/plookup_noheader
drwxr-xr-x  - cloudera cloudera          0 2021-08-08 17:45 clickstream/logs
grunt> fs -ls clickstream/logs
Found 3 items
-rw-r--r--  1 cloudera cloudera 137233 2021-08-08 17:45 clickstream/logs/u_ex160211.log
-rw-r--r--  1 cloudera cloudera 78658 2021-08-08 17:45 clickstream/logs/u_ex160212.log
-rw-r--r--  1 cloudera cloudera 105235 2021-08-08 17:45 clickstream/logs/u_ex160213.log
grunt>
```

```

cloudera@quickstart:
grunt> fs -cat clickstream/logs/u_ex160211.log
#Software: Microsoft Internet Information Services 8.5
#Version: 1.0.0.0
#Date: 2016-08-11 00:01:31
#Fields: date time s-ip cs-method cs-uri-stem cs-uri-query s-port cs-username cc-port cs-useragent cs(Referrer) sc-status sc-substatus sc-idleTime status time-taken
08/16/2016 17:16:13 128.238.247.37 GET / - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; x64; rv:43.0) Gecko/20100101 Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 383
08/16/2016 17:16:13 128.238.247.37 GET /Content/jQuery-ui-themes/smoothness/jquery-ui-1.10.3.custom.min.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0) Gecko/20100101 Firefox/43.0+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 15
08/16/2016 17:16:13 128.238.247.37 GET /Plugins/Widgets.NivoSlider/Content/nivoslider/themes/custom.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 18
08/16/2016 17:16:13 128.238.247.37 GET /Plugins/Widgets.NivoSlider/Content/nivoslider/nivo_slider.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 18
08/16/2016 17:16:13 128.238.247.37 GET /Scripts/jquery-unobtrusive.min.js - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 36
08/16/2016 17:16:13 128.238.247.37 GET /Content/jQueryUI/jQueryUI.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 57
08/16/2016 17:16:13 128.238.247.37 GET /Content/jQueryUI/themes/smoothness/ui.all.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 62
08/16/2016 17:16:13 128.238.247.37 GET /Content/jQueryUI/themes/smoothness/ui.core.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 22
08/16/2016 17:16:13 128.238.247.37 GET /Content/jQueryUI/themes/smoothness/ui.date.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 9
08/16/2016 17:16:13 128.238.247.37 GET /Content/jQueryUI/themes/smoothness/ui.dialog.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 9
08/16/2016 17:16:13 128.238.247.37 GET /Content/jQueryUI/themes/smoothness/ui.slider.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 9
08/16/2016 17:16:13 128.238.247.37 GET /Content/jQueryUI/themes/smoothness/ui.tabs.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 107
08/16/2016 17:16:13 128.238.247.37 GET /Content/jQueryUI/themes/smoothness/ui.widget.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 125
08/16/2016 17:16:13 128.238.247.37 GET /Content/jQueryUI/themes/smoothness/ui.accordion.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 37
08/16/2016 17:16:13 128.238.247.37 GET /Content/jQueryUI/themes/smoothness/ui.button.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 37
08/16/2016 17:16:13 128.238.247.37 GET /Content/jQueryUI/themes/smoothness/ui.select.css - 80 - 215.82.23.2 Mozilla/5.0 (Windows NT 10.0; Win64; rv:43.0)+Gecko/20100101+Firefox/43.0 http://group0.ist722.lschool.syr.edu/ 200 0 0 200
```

```

769-Win10Docker-srajendi

cloudera@quickstart:-
grunt> clickstream_logs = LOAD 'clickstream/logs/*' USING PigStorage(' ')
>> AS (readate:chararray,reqtime:chararray,X1:int,method:chararray,uri:chararray,X2:int,X3:int,X4:int,ipaddress:chararray,useragent:chararray);
grunt> describe clickstream_logs;
clickstream_logs: {readate: chararray,reqtime: chararray,X1: int,method: chararray,uri: chararray,X2: int,X3: int,X4: int,ipaddress: chararray,useragent: chararray}
grunt> DUMP clickstream_logs;
```

```

cloudera@quickstart:
grunt> clickstream_logs = LOAD 'clickstream/logs/*' USING PigStorage(' ')
>> AS (readate:chararray,reqtime:chararray,X1:int,method:chararray,uri:chararray,X2:int,X3:int,X4:int,ipaddress:chararray,useragent:chararray);
clickstream_logs: {readate: chararray,reqtime: chararray,X1: int,method: chararray,uri: chararray,X2: int,X3: int,X4: int,ipaddress: chararray,useragent: chararray}
grunt> DUMP clickstream_logs;
048244 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig Features used in the script: UNKNOWN
048245 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED=[AddOrReplaceColumnMapKeyPrune, DuplicateForEachColumnRewrite, GroupByConstParallelSetter, PigTypeCastInserter, MergeGelfitter, NewPartitionFilterOptimizer, PushDownForEachListner, PushUpFilter, SplitFilter, StreamTypeCastInserter], RULES_DISABLED=[FilterLogicExprOptimizer]}
048455 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic? False
048488 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - MR plan size before optimization: 1
048489 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - MR plan size after optimization: 1
048504 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
048908 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreduces.buffer.percent is not set, set to default 0.3
049495 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - creating jar file Job05340110411575104.jar
049496 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - jar file Job05340110411575104.jar created
054293 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - setting single store job
054303 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
054303 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
054303 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserializable []
054303 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.aliases] with aliases to deserializable []
054788 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process: 1
054819 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths (combined) to process: 1
055395 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases clickstream_logs
055395 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - More information at: http://localhost:50800/jobdetails.jsp?jobid=job_1628829248144_0014
055449 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 0% complete
066617 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 50% complete
071999 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
071181 [main] INFO org.apache.pig.tools.pigstats.SimplePigStats - Script Statistics:
HadoopVersion PigVersion Userid StartedAt FinishedAt Features
2.6.0-cdh5.7.0 0.12.0-cdh5.7.0 cloudera 2021-08-13 22:46:20 2021-08-13 22:46:42 UNKNOWN
Success!
Job stats (time in seconds):
JobID      Map Reduces MaxMapTime   MinMapTime    AvgMapTime   MedianMapTime  MaxReduceTime  MinReduceTime AvgReduceTime
job_1628829248144_0014      0       3        3        3        3        n/a        n/a        n/a        clickstream_logs
AP_ONLY hdfs://quickstart.cloudera:8020/tmp/temp1444246572/tmp677452236,
```

Input(s):
Successfully read 1147 records (32170 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/clickstream/logs/*"

Output(s):
successfully stored 1147 records (252612 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp1444246572/tmp677452236"

769-Win10Docker-srajendi

```
bx cloudera@quickstart:~  
  
grunt> clickstream_logs_clean = FILTER clickstream_logs BY reqdate!='#';  
grunt> clickstream_logs_final = FOREACH clickstream_logs_clean GENERATE reqdate,reqtime,method,uri,ipaddress,useragent;  
grunt> describe clickstream_logs_final;  
clickstream_logs_final: {reqdate: chararray,reqtime: chararray,method: chararray,uri: chararray,ipaddress: chararray,useragent: chararray}  
grunt>
```

```

$ ./clickstream_optimize.sh
grunt clickstream_logs_clean - FILTER Clickstream_logs_BY_Regdate="";
grunt clickstream_logs_final - FOREACH clickstream_logs_clean REGRATE regdate,reftime,method,uri,ipaddress,useragent;
grunt clickstream_logs_final - (regrate,chararray,reftime:chararray,method:chararray,uri:chararray,ipaddress:chararray,useragent:chararray)
153367 [main] INFO org.apache.pig.tools.pigglets.SchemaPigScript - Pig features used in this script: FILTER
153367 [main] INFO org.apache.pig.scripting.PigScriptEngine - (ALL,MAP,REDUCE,ForEach,ColumnMapKeyPrune, DuplicateForEachColumnRewrite, GroupByConstParallelSetter, ImplicitSplitInserter, LimitOptimizer, LoadTypeCastInserter, MergeForEach, NewPartitionOptimizer, PushOnForEachLimiter, PushOnFilter, StreamtypeCastInserter, RULES_DisableDedup[filterlogicExpressionSimplifier, PartitionFilterOptimizer])
153370 [main] INFO org.apache.pig.mapreduce.logical.RuleManager - ColumnPumpForWriten: Columns pruned for clickstream_logs: $3, $5, $6, $7
153370 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.MRCompiler - File concatenation threshold: 100 optimistic? false
153370 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.MultiQueryOptimizer - MR plan size before optimization: 1
153370 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.MultiQueryOptimizer - MR plan size after optimization: 1
153372 [main] INFO org.apache.pig.tools.pigglets.ScriptState - Pig script settings are added to the job
153372 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.JobControlCompiler - mapped.job.reduce.markreset.buffer.percent is not set, set to default 0.3
153372 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.JobControlCompiler - creating jar file /tmp/job_162829248144_0015/jobcontrol.jar
153372 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.JobControlCompiler - JobControlJarPath: /tmp/job_162829248144_0015/jobcontrol.jar created
153372 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.JobControlCompiler - Setting up single store job
1533887 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.JobControlCompiler - Generating code for mapred execution engine
1533889 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is False, will not generate code
1533889 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Generating code for mapred execution engine
1533889 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple_classes] with classes to deserialze []
1533889 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.MapredLauncher - i-map reduce job(s) waiting for submission.
1533889 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.MapredLauncher - i-map reduce job(s) waiting for submission.
1533929 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths (combined) to process : 1
1533931 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.MapredLauncher - HadoopJobId: job_162829248144_0015
1533931 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.MapredLauncher - Processing aliases clickstream_logs.clickstream_logs_clean.clickstream_logs_final
1533931 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.MapredLauncher - Job details: N-clickstream_logs[clickstream_logs[1..19],clickstream_logs[-1..1],clickstream_logs_clean[7..25],clickstream_logs_final[8..25]; C: R;
1533931 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.MapredLauncher - More information at: http://localhost:50030/jobdetails.jsp?jobid=job_162829248144_0015
1534068 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.MapredLauncher - OK complete
1535076 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.MapredLauncher - %OK complete
1535446 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapredExecutioner.MapredLauncher - 100% complete
1534465 [main] INFO org.apache.pig.tools.pigglets.Simplifier - Script statistics.

HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.6.0-cdh7.0 0.12.0-cdh7.0 cloudera 2021-08-13 23:01:05 2021-08-13 23:01:25 FILTER
Success!

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime
clickstream_logs.clickstream_logs_clean.clickstream_logs_final 162829248144_0015 3 3 3 n/a n/a n/a clickstream_logs.clickstream_logs_clean.clickstream_logs_final
ok ONLY /user/cloudera/clickstream/logs_noheader

Input(s):
Successfully read 1147 records (321710 bytes) from: "hdfs://quickstart.cloudera:8820/user/cloudera/clickstream/logs/*"

```

```
[cloudera@quickstart- Job Stats (time in seconds):
JobId  Maps Reduces MaxMapTime     MinMapTime      AvgMapTime      MedianMapTime   MaxReduceTime  MinReduceTime  AvgReduceTime
job_1628829248144_0015 1          0            3             3              3             n/a           n/a           n/a           clickstream_logs,clickstream_logs_clean,clickstream_logs_final
AP_ONLY /user/cloudera/clickstream/logs_noheader,
Input(s):
Successfully read 1147 records (321710 bytes) from: "hdfs://quickstart.cloudera:8020/user/cloudera/clickstream/logs/*"
Output(s):
Successfully stored 1147 records (227447 bytes) in: "/user/cloudera/clickstream/logs_noheader"
Counters:
Total records written : 1147
Total bytes written : 227447
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_1628829248144_0015
grunt> 1554514 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
```

```
cloudera@quickstart:~  
  
grunt> fs -ls clickstream  
Found 4 items  
drwxr-xr-x  - cloudera cloudera          0 2021-08-08 17:44 clickstream/iplookup  
drwxr-xr-x  - cloudera cloudera          0 2021-08-13 21:19 clickstream/iplookup_noheader  
drwxr-xr-x  - cloudera cloudera          0 2021-08-08 17:45 clickstream/logs  
drwxr-xr-x  - cloudera cloudera          0 2021-08-13 23:01 clickstream/logs_noheader  
grunt> fs -ls clickstream/logs_noheader  
Found 2 items  
-rw-r--r--  1 cloudera cloudera          0 2021-08-13 23:01 clickstream/logs_noheader/_SUCCESS  
-rw-r--r--  1 cloudera cloudera  227447 2021-08-13 23:01 clickstream/logs_noheader/part-m-00000  
grunt> ■
```

769-Win10Docker-srajendi

```
cloudera@quickstart:~$ grunt quit
[cloudera@quickstart ~]$ hdfs dfs -ls
Found 7 items
drwxr-xr-x  - cloudera cloudera          0 2021-08-13 23:01 clickstream
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:32 fudgemart-clothing
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:23 fudgemart-products-by-clothing
drwxr-xr-x  - cloudera cloudera          0 2021-08-09 23:40 sotu2016
drwxr-xr-x  - cloudera cloudera          0 2021-08-13 00:23 testing
drwxr-xr-x  - cloudera cloudera          0 2021-08-06 06:34 text
drwxr-xr-x  - cloudera cloudera          0 2021-08-10 00:47 tweets
[cloudera@quickstart ~]$ hdfs dfs -ls clickstream
Found 4 items
drwxr-xr-x  - cloudera cloudera          0 2021-08-08 17:44 clickstream/iplookup
drwxr-xr-x  - cloudera cloudera          0 2021-08-13 21:19 clickstream/iplookup_noheader
drwxr-xr-x  - cloudera cloudera          0 2021-08-08 17:45 clickstream/logs
drwxr-xr-x  - cloudera cloudera          0 2021-08-13 23:01 clickstream/logs_noheader
[cloudera@quickstart ~]$ hdfs dfs -ls clickstream/*header
Found 2 items
-rw-r--r--  1 cloudera cloudera          0 2021-08-13 21:19 clickstream/iplookup_noheader/_SUCCESS
-rw-r--r--  1 cloudera cloudera        1185 2021-08-13 21:19 clickstream/iplookup_noheader/part-r-00000
Found 2 items
-rw-r--r--  1 cloudera cloudera          0 2021-08-13 23:01 clickstream/logs_noheader/_SUCCESS
-rw-r--r--  1 cloudera cloudera      227447 2021-08-13 23:01 clickstream/logs_noheader/part-m-00000
[cloudera@quickstart ~]$
```

769-Win10Docker-srajendi

```
cloudera@quickstart:~$ beeline -u jdbc:hive2://localhost:10000/default -u cloudera -p cloudera
2021-08-13 23:26:29,123 WARN [main] mapreduce.TableMapReduceUtil: The hbase-prefix-tree module jar containing PrefixTreeCodec is not present. Continuing without it.
scan complete in 3ms
Connecting to jdbc:hive2://localhost:10000/default
Connected to: Apache Hive (version 1.1.0-cdh5.7.0)
Driver: Hive JDBC (version 1.1.0-cdh5.7.0)
Transaction isolation: TRANSACTION_REPEATABLE_READ
Beeline version 1.1.0-cdh5.7.0 by Apache Hive
0: jdbc:hive2://localhost:10000/default> show databases;
INFO : Compiling command(queryId:hive_20210813231515_4b33b72a-dc97-4428-974e-a04502dff72): show databases
INFO : Semantic Analysis Completed
INFO : Returning Hive schema:Schema(fieldsSchemas:[FieldSchema(name:database_name, type:string, comment:from deserializer)], properties:null)
INFO : Completed compiling command(queryId:hive_20210813231515_4b33b72a-dc97-4428-974e-a04502dff72); Time taken: 1.298 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId:hive_20210813231515_4b33b72a-dc97-4428-974e-a04502dff72): show databases
INFO : Starting task [Stage-0:DDL] in serial mode
INFO : Completed executing command(queryId:hive_20210813231515_4b33b72a-dc97-4428-974e-a04502dff72); Time taken: 0.419 seconds
INFO : OK
+-----+-----+
| database_name |
+-----+-----+
| clickstream   |
| default       |
+-----+-----+
2 rows selected (2.162 seconds)
```

```
cloudera@quickstart:~$ beeline -u jdbc:hive2://localhost:10000/default -u cloudera -p cloudera
2021-08-13 23:26:29,123 WARN [main] mapreduce.TableMapReduceUtil: The hbase-prefix-tree module jar containing PrefixTreeCodec is not present. Continuing without it.
scan complete in 3ms
Connecting to jdbc:hive2://localhost:10000/default
Connected to: Apache Hive (version 1.1.0-cdh5.7.0)
Driver: Hive JDBC (version 1.1.0-cdh5.7.0)
Transaction isolation: TRANSACTION_REPEATABLE_READ
Beeline version 1.1.0-cdh5.7.0 by Apache Hive
0: jdbc:hive2://localhost:10000/default> show databases;
INFO : Compiling command(queryId:hive_20210813232626_92d6ab46-8595-48f2-96c5-ce0c5b27efd0): show databases
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fieldsSchemas:[FieldSchema(name:database_name, type:string, comment:from deserializer)], properties:null)
INFO : Completed compiling command(queryId:hive_20210813232626_92d6ab46-8595-48f2-96c5-ce0c5b27efd0); Time taken: 0.031 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(queryId:hive_20210813232626_92d6ab46-8595-48f2-96c5-ce0c5b27efd0): show databases
INFO : Starting task [Stage-0:DDL] in serial mode
INFO : Completed executing command(queryId:hive_20210813232626_92d6ab46-8595-48f2-96c5-ce0c5b27efd0); Time taken: 0.049 seconds
INFO : OK
+-----+-----+
| database_name |
+-----+-----+
| clickstream   |
| default       |
+-----+-----+
2 rows selected (0.224 seconds)
```

769-Win10Docker-srajendi

```
cloudera@quickstart:~  
0: jdbc:hive2://localhost:10000/default> DROP DATABASE IF EXISTS clickstream CASCADE;  
INFO : Compiling command(queryId=hive_20210813232727_3d960416-98b1-460b-92c1-f86f496fcbe8): DROP DATABASE IF EXISTS clickstream CASCADE  
INFO : Semantic Analysis Completed  
INFO : Returning Hive schema: Schema(fieldSchemas:null, properties:null)  
INFO : Completed compiling command(queryId=hive_20210813232727_3d960416-98b1-460b-92c1-f86f496fcbe8); Time taken: 0.665 seconds  
INFO : Concurrency mode is disabled, not creating a lock manager  
INFO : Executing command(queryId=hive_20210813232727_3d960416-98b1-460b-92c1-f86f496fcbe8): DROP DATABASE IF EXISTS clickstream CASCADE  
INFO : Starting task [Stage-0:DDL] in serial mode  
INFO : Completed executing command(queryId=hive_20210813232727_3d960416-98b1-460b-92c1-f86f496fcbe8); Time taken: 7.948 seconds  
INFO : OK  
No rows affected (8.656 seconds)  
0: jdbc:hive2://localhost:10000/default> show databases;  
INFO : Compiling command(queryId=hive_20210813232727_8c103cad-6b1e-4dbf-8a82-66746a8cf1b2): show databases  
INFO : Semantic Analysis Completed  
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:database_name, type:string, comment:from deserializer)], properties:null)  
INFO : Completed compiling command(queryId=hive_20210813232727_8c103cad-6b1e-4dbf-8a82-66746a8cf1b2); Time taken: 0.007 seconds  
INFO : Concurrency mode is disabled, not creating a lock manager  
INFO : Executing command(queryId=hive_20210813232727_8c103cad-6b1e-4dbf-8a82-66746a8cf1b2): show databases  
INFO : Starting task [Stage-0:DDL] in serial mode  
INFO : Completed executing command(queryId=hive_20210813232727_8c103cad-6b1e-4dbf-8a82-66746a8cf1b2); Time taken: 0.019 seconds  
INFO : OK  
+-----+  
| database_name |  
+-----+  
| default |  
+-----+  
1 row selected (0.001 seconds)  
0: jdbc:hive2://localhost:10000/default>
```

769-Win10Docker-srajendi

```
cloudera@quickstart:~  
0: jdbc:hive2://localhost:10000/default> CREATE DATABASE clickstream;  
INFO : Compiling command(queryId=hive_20210813232828_04815a86-f0fa-49f4-8d31-214de74d3e6d): CREATE DATABASE clickstream  
INFO : Semantic Analysis Completed  
INFO : Returning Hive schema: Schema(fieldSchemas:null, properties:null)  
INFO : Completed compiling command(queryId=hive_20210813232828_04815a86-f0fa-49f4-8d31-214de74d3e6d); Time taken: 0.007 seconds  
INFO : Concurrency mode is disabled, not creating a lock manager  
INFO : Executing command(queryId=hive_20210813232828_04815a86-f0fa-49f4-8d31-214de74d3e6d): CREATE DATABASE clickstream  
INFO : Starting task [Stage-0:DDL] in serial mode  
INFO : Completed executing command(queryId=hive_20210813232828_04815a86-f0fa-49f4-8d31-214de74d3e6d); Time taken: 0.212 seconds  
INFO : OK  
No rows affected (0.243 seconds)  
0: jdbc:hive2://localhost:10000/default> show databases;  
INFO : Compiling command(queryId=hive_20210813232828_666a979d-f78a-4c96-ab33-f216dac5f4c4): show databases  
INFO : Semantic Analysis Completed  
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:database_name, type:string, comment:from deserializer)], properties:null)  
INFO : Completed compiling command(queryId=hive_20210813232828_666a979d-f78a-4c96-ab33-f216dac5f4c4); Time taken: 0.008 seconds  
INFO : Concurrency mode is disabled, not creating a lock manager  
INFO : Executing command(queryId=hive_20210813232828_666a979d-f78a-4c96-ab33-f216dac5f4c4): show databases  
INFO : Starting task [Stage-0:DDL] in serial mode  
INFO : Completed executing command(queryId=hive_20210813232828_666a979d-f78a-4c96-ab33-f216dac5f4c4); Time taken: 0.006 seconds  
INFO : OK  
+-----+  
| database_name |  
+-----+  
| clickstream |  
| default |  
+-----+  
2 rows selected (0.055 seconds)  
0: jdbc:hive2://localhost:10000/default> show tables;  
INFO : Compiling command(queryId=hive_20210813232929_506dd7cd-f3eb-49be-9578-7d2262b007c2): show tablesream  
INFO : Semantic Analysis Completed  
INFO : Returning Hive schema: Schema(fieldSchemas:[FieldsSchema(name:tab_name, type:string, comment:from deserializer)], properties:null)  
INFO : Completed compiling command(queryId=hive_20210813232929_506dd7cd-f3eb-49be-9578-7d2262b007c2); Time taken: 0.007 seconds  
INFO : Concurrency mode is disabled, not creating a lock manager  
INFO : Executing command(queryId=hive_20210813232929_506dd7cd-f3eb-49be-9578-7d2262b007c2): show tablesream  
INFO : Starting task [Stage-0:DDL] in serial mode  
INFO : Completed executing command(queryId=hive_20210813232929_506dd7cd-f3eb-49be-9578-7d2262b007c2); Time taken: 0.012 seconds  
INFO : OK  
+-----+ (0.059 seconds)  
| tab_name |://localhost:10000/default> swo  
+-----+  
+-----+  
No rows selected (0.055 seconds)  
0: jdbc:hive2://localhost:10000/default> -
```

769-Win10Docker-srajendi

```
cloudera@quickstart:~  
0: jdbc:hive2://localhost:10000/default> show tables;  
INFO : Compiling command(queryId=hive_20210813232929_506dd7cd-f3eb-49be-9578-7d2262b007c2): show tablesream  
INFO : Semantic Analysis Completed  
INFO : Returning Hive schema: Schema(fieldsSchemas:[fieldsSchema(name:tab_name, type:string, comment:from deserializer)], properties:null)  
INFO : Completed compiling command(queryId=hive_20210813232929_506dd7cd-f3eb-49be-9578-7d2262b007c2); Time taken: 0.007 seconds  
INFO : Concurrency mode is disabled, not creating a lock manager  
INFO : Executing command(queryId=hive_20210813232929_506dd7cd-f3eb-49be-9578-7d2262b007c2): show tablesream  
INFO : Starting task [Stage-0:DDL] in serial mode  
INFO : Completed executing command(queryId=hive_20210813232929_506dd7cd-f3eb-49be-9578-7d2262b007c2); Time taken: 0.012 seconds  
INFO : OK  
+-----+--- (0.059 seconds)  
| tab_name |:/localhost:10000/default> swo  
+-----+  
+-----+  
No rows selected (0.059 seconds)  
0: jdbc:hive2://localhost:10000/default> !sh hdfs dfs -ls clickstream/*header  
Found 2 items  
-rw-r--r-- 1 cloudera cloudera 0 2021-08-13 21:19 clickstream/iplookup_noheader/_SUCCESS  
-rw-r--r-- 1 cloudera cloudera 1185 2021-08-13 21:19 clickstream/iplookup_noheader/part-r-00000  
Found 2 items  
-rw-r--r-- 1 cloudera cloudera 0 2021-08-13 23:01 clickstream/logs_noheader/_SUCCESS  
-rw-r--r-- 1 cloudera cloudera 227447 2021-08-13 23:01 clickstream/logs_noheader/part-m-00000  
0: jdbc:hive2://localhost:10000/default>
```

769-Win10Docker-srajendi

```
cloudera@quickstart:~  
0: jdbc:hive2://localhost:10000/default> create external table iplookup (   
0: jdbc:hive2://localhost:10000/default> ip string,country string,state string,city string,lat double,lan double)  
0: jdbc:hive2://localhost:10000/default> row format delimited  
0: jdbc:hive2://localhost:10000/default> fields terminated by ','  
0: jdbc:hive2://localhost:10000/default> location '/user/cloudera/clickstream/iplookup_noheader/'  
INFO : Compiling command(queryId=hive_20210813234848_d3ab2589-aba9-4f97-8a56-701a5b8c27f8): create external table iplookup (   
ip string,country string,state string,city string,lat double,lan double)  
row format delimited  
fields terminated by ','  
location '/user/cloudera/clickstream/iplookup_noheader/'  
INFO : Semantic Analysis Completed  
INFO : Returning Hive schema: Schema(fieldsSchemas:null, properties:null)  
INFO : Completed compiling command(queryId=hive_20210813234848_d3ab2589-aba9-4f97-8a56-701a5b8c27f8); Time taken: 0.027 seconds  
INFO : Concurrency mode is disabled, not creating a lock manager  
INFO : Executing command(queryId=hive_20210813234848_d3ab2589-aba9-4f97-8a56-701a5b8c27f8): create external table iplookup (   
ip string,country string,state string,city string,lat double,lan double)  
row format delimited  
fields terminated by ','  
location '/user/cloudera/clickstream/iplookup_noheader/'  
INFO : Starting task [Stage-0:DDL] in serial mode  
INFO : Completed executing command(queryId=hive_20210813234848_d3ab2589-aba9-4f97-8a56-701a5b8c27f8); Time taken: 0.16 seconds  
INFO : OK  
No rows affected (0.204 seconds)  
0: jdbc:hive2://localhost:10000/default>
```

769-Win10Docker-srajendi

```
Enforce US Keyboard  
cloudera@quickstart:~  
0: jdbc:hive2://localhost:10000/default> create external table weblogs (   
0: jdbc:hive2://localhost:10000/default> redate date,reqtime timestamp,method string,uri string,ipaddress string,useragent string)  
0: jdbc:hive2://localhost:10000/default> row format delimited  
0: jdbc:hive2://localhost:10000/default> fields terminated by '\n'  
0: jdbc:hive2://localhost:10000/default> location '/user/cloudera/clickstream/logs_noheader/'  
Error: Failed while compiling statement: FAILED: ParseException line 2:56 mismatched input 'ipaddress' expecting ) near 'string' in create table statement (state=42000,code=40000)  
0: jdbc:hive2://localhost:10000/default> create external table weblogs (   
0: jdbc:hive2://localhost:10000/default> redate date,reqtime timestamp,method string,uri string,ipaddress string,useragent string)  
0: jdbc:hive2://localhost:10000/default> row format delimited  
0: jdbc:hive2://localhost:10000/default> fields terminated by '\n'  
0: jdbc:hive2://localhost:10000/default> location '/user/cloudera/clickstream/logs_noheader/'  
INFO : Compiling command(queryId=hive_20210813235353_a49da797-ee6b-40d8-96eb-41a9ed5b1a7): create external table weblogs (   
redate date,reqtime timestamp,method string,uri string,ipaddress string,useragent string)  
row format delimited  
fields terminated by '\n'  
location '/user/cloudera/clickstream/logs_noheader/'  
INFO : Semantic Analysis Completed  
INFO : Returning Hive schema: Schema(fieldsSchemas:null, properties:null)  
INFO : Completed compiling command(queryId=hive_20210813235353_a49da797-ee6b-40d8-96eb-41a9ed5b1a7); Time taken: 0.027 seconds  
INFO : Concurrency mode is disabled, not creating a lock manager  
INFO : Executing command(queryId=hive_20210813235353_a49da797-ee6b-40d8-96eb-41a9ed5b1a7): create external table weblogs (   
redate date,reqtime timestamp,method string,uri string,ipaddress string,useragent string)  
row format delimited  
fields terminated by '\n'  
location '/user/cloudera/clickstream/logs_noheader/'  
INFO : Starting task [Stage-0:DDL] in serial mode  
INFO : Completed executing command(queryId=hive_20210813235353_a49da797-ee6b-40d8-96eb-41a9ed5b1a7); Time taken: 0.001 seconds  
INFO : OK  
No rows affected (0.107 seconds)  
0: jdbc:hive2://localhost:10000/default>
```

```

768-Win10Docker-srajendi
[cloudera@quickstart: ~]
: jdbc:hive2://localhost:10000/default> select * from iplookup limit 5
INFO : Compiling command(queryId=hive_20210813235454_fcc6-a9d3-b5e4-b246f47e948); select * from iplookup limit 5
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Struct<fieldSchemas:[FieldSchema(name:iplookup_ip, type:string, comment:null), FieldSchema(name:iplookup_country, type:string, comment:null), FieldSchema(name:iplookup_state, type:string, comment:null), FieldSchema(name:iplookup_city, type:string, comment:null), FieldSchema(name:iplookup_lat, type:double, comment:null), FieldSchema(name:iplookup_lng, type:double, comment:null)], properties=null>
INFO : Compiled compiling command(queryId=hive_20210813235454_fcc6-a9d3-b5e4-b246f47e948); Time taken: 0.350 seconds
INFO : Executing command(queryId=hive_20210813235454_fcc6-a9d3-b5e4-b246f47e948); select * from iplookup limit 5
INFO : Compiled executing command(queryId=hive_20210813235454_fcc6-a9d3-b5e4-b246f47e948); Time taken: 0.001 seconds
INFO : OK
+-----+
| iplookup_ip | iplookup_country | iplookup_state | iplookup_city | iplookup_lat | iplookup_lng |
+-----+
| 128.122.140.238 | USA | NY | New York | 40.712784 | -74.005041 |
| 131.159.205.152 | USA | CA | San Jose | 37.331459 | -122.051529 |
| 155.100.169.152 | USA | UT | Salt Lake City | 40.768779 | -111.891047 |
| 172.189.252.8 | USA | VA | Dulles | 38.955855 | -77.447819 |
| 215.82.209.10 | USA | OH | Columbus | 39.961376 | 82.998794 |
+-----+
5 rows selected (0.02 seconds)
: jdbc:hive2://localhost:10000/default> 

```



```

[cloudera@quickstart: ~]
: jdbc:hive2://localhost:10000/default> describe iplookup;
INFO : Describing table iplookup in database default
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Struct<fieldSchemas:[FieldSchema(name:col_name, type:date, comment:null), FieldSchema(name:data_type, type:string, comment:from deserializer), FieldSchemas(name:comment, type:string, comment:from deserializer), properties:null>
INFO : Compiled compiling command(queryId=hive_20210813235456_f722ff601-3466-4595-93de-f6142008f317); describe iplookup
INFO : Executing command(queryId=hive_20210813235456_f722ff601-3466-4595-93de-f6142008f317); describe iplookup
INFO : Starting task [Stage=0+0.001] in serial mode
INFO : Completed successfully executing command(queryId=hive_20210813235456_f722ff601-3466-4595-93de-f6142008f317); Time taken: 0.024 seconds
INFO : OK
+-----+
| col_name | data_type | comment |
+-----+
| ip | string |          |
| country | string |          |
| state | string |          |
| city | string |          |
| lat | double |          |
| lng | double |          |
+-----+
6 rows selected (0.144 seconds)
: jdbc:hive2://localhost:10000/default> describe weblogs;
INFO : Compiling command(queryId=hive_20210813235566_f2816851-5e3b-4011-ba99-bea7c85b3e00); describe weblogs
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Struct<fieldSchemas:[FieldSchema(name:col_name, type:string, comment:from deserializer), FieldSchema(name:data_type, type:string, comment:from deserializer), FieldSchema(name:comment, type:string, comment:from deserializer), properties:null>
INFO : Compiled compiling command(queryId=hive_20210813235566_f2816851-5e3b-4011-ba99-bea7c85b3e00); Time taken: 0.078 seconds
INFO : Executing command(queryId=hive_20210813235566_f2816851-5e3b-4011-ba99-bea7c85b3e00); describe weblogs
INFO : Starting task [Stage=0+0.001] in serial mode
INFO : Completed successfully executing command(queryId=hive_20210813235566_f2816851-5e3b-4011-ba99-bea7c85b3e00); Time taken: 0.024 seconds
INFO : OK
+-----+
| col_name | data_type | comment |
+-----+
| regdate | date |          |
| regtime | timestamp |          |
| user | string |          |
| url | string |          |
| ipaddress | string |          |
| useragent | string |          |
+-----+
6 rows selected (0.148 seconds)
: jdbc:hive2://localhost:10000/default> 

```

769-Win10Docker-srajendi

```

[cloudera@quickstart: ~]
: jdbc:hive2://localhost:10000/default> select
0: jdbc:hive2://localhost:10000/default>   i.city, count(w.method) as http_requests
0: jdbc:hive2://localhost:10000/default>   from iplookup as i join weblogs as w
0: jdbc:hive2://localhost:10000/default>   on i.ip = w.ipaddress
0: jdbc:hive2://localhost:10000/default>   group by i.city
0: jdbc:hive2://localhost:10000/default>   order by http_requests desc;

```

```

cloudera@quickstart:~$ cloudera@quickstart:~$ 
9: jdbc:hive2://localhost:10000/default> select
9: jdbc:hive2://localhost:10000/default> i.city.count(w.method) as http_requests
9: jdbc:hive2://localhost:10000/default> from weblogs w inner join weblogs as w
9: jdbc:hive2://localhost:10000/default> on i.ip = w.ipaddress
9: jdbc:hive2://localhost:10000/default> group by i.city
9: jdbc:hive2://localhost:10000/default> order by http_requests desc;
INFO : Compiling command(queryId=hive_20210814000404_239f1460-9b8f-4cbc-a59e-337aa1ebd637): select
i.city,count(w.method) as http_requests
From iplookup as i join weblogs as w
on i.ip = w.ipaddress
group by i.city
order by http_requests desc;
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema[FieldSchemas:[FieldSchema(name:i.city, type:string, comment:null), FieldSchema(name:http_requests, type:bigint, comment:null)], properties:null]
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing query(queryId=hive_20210814000404_239f1460-9b8f-4cbc-a59e-337aa1ebd637): select
i.city,count(w.method) as http_requests
From iplookup as i join weblogs as w
on i.ip = w.ipaddress
group by i.city
order by http_requests desc;
INFO : Total Job = 2
INFO : Starting task [Stage-0:MAPRED] in serial mode
INFO : Mapred task completed successfully
INFO : Mapredlocal task succeeded
INFO : Launching Job 1 out of 2
INFO : Starting task [Stage-2:MAPRED] in serial mode
INFO : Number of reduce tasks determined at compile time: 1
INFO : In order to change the average load for a reducer (in bytes):
INFO :   set hive.exec.reducers.bytes.per.reducer=<number>
INFO : In order to limit the maximum number of reducers:
INFO :   set hive.exec.reducers.max=<number>
INFO : In order to set a constant number of reducers:
INFO :   set mapreduce.job.reduces=<number>
INFO : Starting Job : job_1628829248144_0016, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1628829248144_0016/
INFO : Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1628829248144_0016
INFO : MapReduce job information for Stage-0: number of mappers: 1; number of reducers: 1
INFO : 2021-08-14 00:05:03,089 Stage-0 map = 100%, reduce = 0%, Cumulative CPU 2.15 sec
INFO : 2021-08-14 00:05:11,509 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 3.62 sec
INFO : MapReduce Total cumulative CPU time: 3 seconds 620 msec
INFO : Ended Job = job_1628829248144_0016
INFO : Launching Job 2 out of 2
INFO : Starting task [Stage-3:MAPRED] in serial mode
INFO : Number of reduce tasks determined at compile time: 1
INFO : In order to change the average load for a reducer (in bytes):
INFO :   set hive.exec.reducers.bytes.per.reducer=<number>
INFO : In order to limit the maximum number of reducers:
INFO :   set hive.exec.reducers.max=<number>
INFO : In order to set a constant number of reducers:
INFO :   set mapreduce.job.reduces=<number>
INFO : Starting Job : job_1628829248144_0017, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1628829248144_0017/
INFO : Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1628829248144_0017
INFO : Hadoop job information for Stage-3: number of mappers: 1; number of reducers: 1
INFO : 2021-08-14 00:05:18,582 Stage-3 map = 0%, reduce = 0%
INFO : 2021-08-14 00:05:23,856 Stage-3 map = 100%, reduce = 0%, Cumulative CPU 0.95 sec
INFO : 2021-08-14 00:05:30,175 Stage-3 map = 100%, reduce = 100%, Cumulative CPU 2.45 sec
INFO : MapReduce Total cumulative CPU time: 2 seconds 450 msec
INFO : Ended Job = job_1628829248144_0017
INFO : MapReduce Jobs Launched:
INFO : Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 3.62 sec HDFS Read: 237563 HDFS Write: 456 SUCCESS
INFO : Stage-Stage-3: Map: 1 Reduce: 1 Cumulative CPU: 2.45 sec HDFS Read: 5093 HDFS Write: 164 SUCCESS
INFO : Total MapReduce CPU Time Spent: 6 seconds 70 msec
INFO : Completed executing command(queryId=hive_20210814000404_239f1460-9b8f-4cbc-a59e-337aa1ebd637); Time taken: 47.451 seconds
INFO : OK
+-----+-----+
| i.city | http_requests |
+-----+-----+
| Syracuse | 272 |
| Columbus | 152 |
| Dulles | 91 |
| Jersey City | 86 |
| Dallas | 85 |
| Freeport | 78 |
| Arlington | 72 |
| New York | 68 |
| Raleigh | 54 |
| Tampa | 54 |
| Salt Lake City | 51 |
| Cleveland | 35 |
| Los Angeles | 24 |
+-----+-----+
13 rows selected (47.911 seconds)
9: jdbc:hive2://localhost:10000/default>

```

```

769-Win10Docker-srajendi:~$ cloudera@quickstart:~$ 
cloudera@quickstart:~$ 
INFO : Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1628829248144_0016
INFO : Hadoop job information for Stage-2: number of mappers: 1; number of reducers: 1
INFO : 2021-08-14 00:04:56,675 Stage-2 map = 0%, reduce = 0%
INFO : 2021-08-14 00:05:03,089 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 2.15 sec
INFO : 2021-08-14 00:05:11,509 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 3.62 sec
INFO : MapReduce Total cumulative CPU time: 3 seconds 620 msec
INFO : Ended Job = job_1628829248144_0016
INFO : Launching Job 2 out of 2
INFO : Starting task [Stage-3:MAPRED] in serial mode
INFO : Number of reduce tasks determined at compile time: 1
INFO : In order to change the average load for a reducer (in bytes):
INFO :   set hive.exec.reducers.bytes.per.reducer=<number>
INFO : In order to limit the maximum number of reducers:
INFO :   set hive.exec.reducers.max=<number>
INFO : In order to set a constant number of reducers:
INFO :   set mapreduce.job.reduces=<number>
INFO : Starting Job : job_1628829248144_0017, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1628829248144_0017/
INFO : Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1628829248144_0017
INFO : Hadoop job information for Stage-3: number of mappers: 1; number of reducers: 1
INFO : 2021-08-14 00:04:55,675 Stage-3 map = 0%, reduce = 0%
INFO : 2021-08-14 00:05:03,089 Stage-3 map = 100%, reduce = 100%, Cumulative CPU 0.95 sec
INFO : 2021-08-14 00:05:30,175 Stage-3 map = 100%, reduce = 100%, Cumulative CPU 2.45 sec
INFO : MapReduce Total cumulative CPU time: 2 seconds 450 msec
INFO : Ended Job = job_1628829248144_0017
INFO : MapReduce Jobs Launched:
INFO : Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 3.62 sec HDFS Read: 237563 HDFS Write: 456 SUCCESS
INFO : Stage-Stage-3: Map: 1 Reduce: 1 Cumulative CPU: 2.45 sec HDFS Read: 5093 HDFS Write: 164 SUCCESS
INFO : Total MapReduce CPU Time Spent: 6 seconds 70 msec
INFO : Completed executing command(queryId=hive_20210814000404_239f1460-9b8f-4cbc-a59e-337aa1ebd637); Time taken: 47.451 seconds
INFO : OK
+-----+-----+
| i.city | http_requests |
+-----+-----+
| Syracuse | 272 |
| Columbus | 152 |
| Dulles | 91 |
| Jersey City | 86 |
| Dallas | 85 |
| Freeport | 78 |
| Arlington | 72 |
| New York | 68 |
| Raleigh | 54 |
| Tampa | 54 |
| Salt Lake City | 51 |
| Cleveland | 35 |
| Los Angeles | 24 |
+-----+-----+
13 rows selected (47.911 seconds)
9: jdbc:hive2://localhost:10000/default>

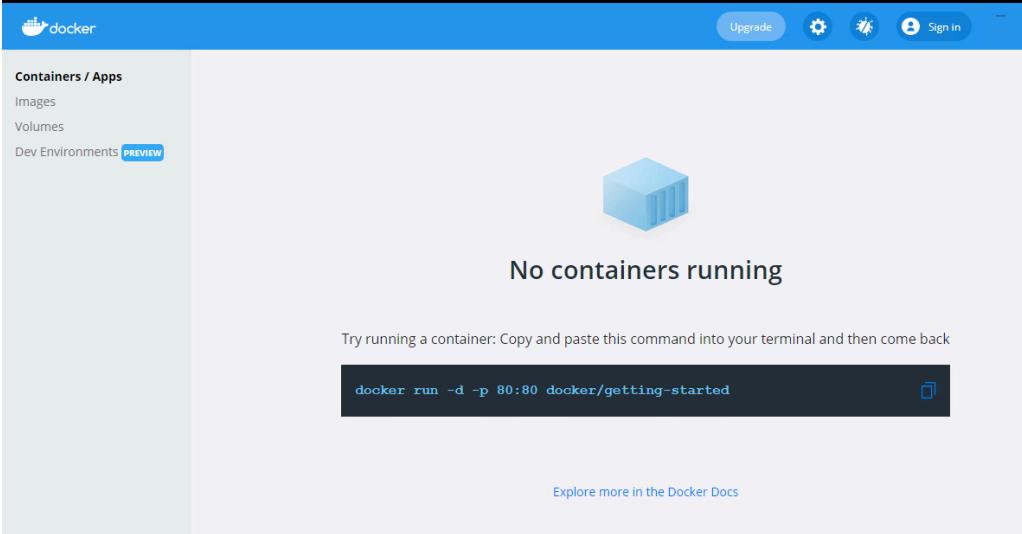
```

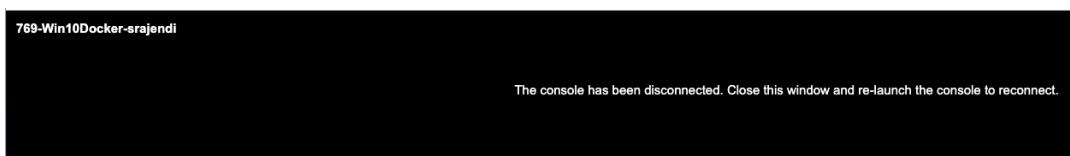
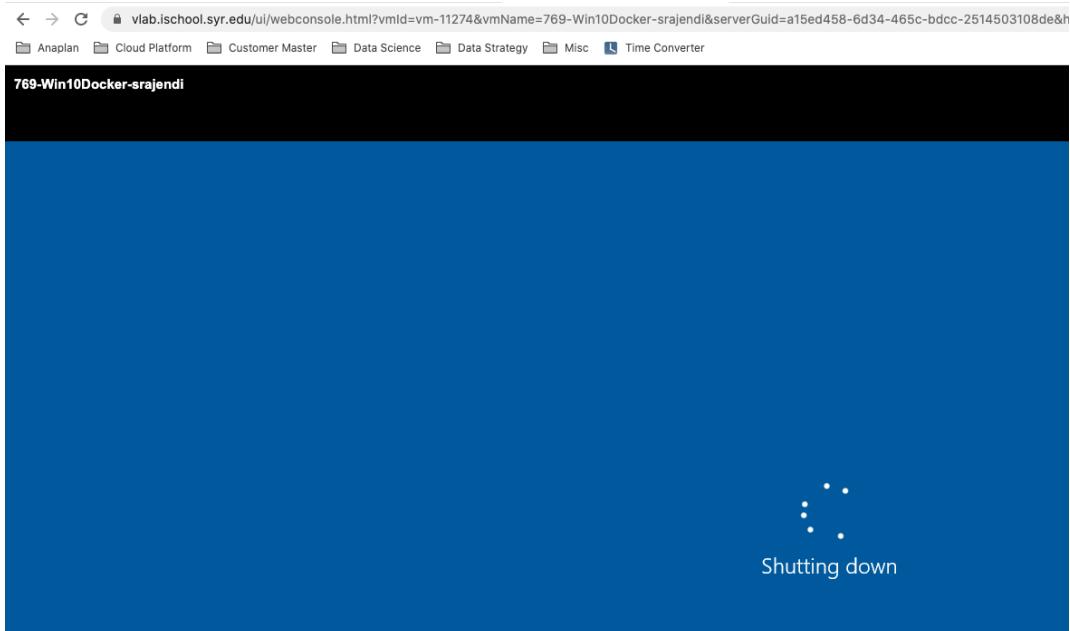
769-Win10Docker-srajendi

```
c:\ Administrator: Command Prompt
+-----+-----+
| i.city | http_requests |
+-----+-----+
| Syracuse | 272
| Columbus | 152
| Dulles | 91
| Jersey City | 86
| Dallas | 85
| Freeport | 78
| Arlington | 72
| New York | 68
| Raleigh | 54
| Tampa | 54
| Salt Lake City | 51
| Cleveland | 35
| Los Angeles | 24
+-----+-----+
13 rows selected (47.911 seconds)
0: jdbc:hive2://localhost:10000/default> !exit
Closing: 0: jdbc:hive2://localhost:10000/default
```

```
[cloudera@quickstart ~]$ exit
logout
C:\Users\LocalAdmin\srajendi\adv-db-labs\hadoop>docker-compose ps
Name          Command           State          Ports
cloudera     /usr/bin/docker-quickstart Up      0.0.0.0:7180->7180/tcp,:::7180->7180/tcp, 0.0.0.0:8080->80/tcp,:::8080->80/tcp, 0.0.0.0:8888->8888/tcp,:::8888->8888/tcp
C:\Users\LocalAdmin\srajendi\adv-db-labs\hadoop>docker-compose down
Stopping cloudera ... done
Removing cloudera ... done
Removing network hadoop_default
C:\Users\LocalAdmin\srajendi\adv-db-labs\hadoop>docker-compose ps
Name          Command           State
C:\Users\LocalAdmin\srajendi\adv-db-labs\hadoop>
```

769-Win10Docker-srajendi





A screenshot of the vSphere Client interface. The top navigation bar includes "vm", "vSphere Client", "Menu", and a search bar. The left sidebar shows a tree structure of vSphere resources under "vlab.ischool.syr.edu": "vLab-Students" contains "Classes-Students" and "IST769-M400-Block", which in turn contains "srajendi" and the specific VM "769-Win10Docker-srajendi". The right pane is titled "769-Win10Docker-srajendi" and shows the "Summary" tab selected. The summary details indicate the VM is "Powered Off". Other tabs include "Monitor", "Configure", "Permissions", "Datastores", and "Networks". Below the summary are sections for "Guest OS", "Compatibility", "VMware Tools", "DNS Name", "IP Addresses", and "Host". Buttons for "Launch Web Console" and "Launch Remote Console" are also present.