

## Experiment : 06

### Implement word count problem using pig

The image shows two screenshots of a Cloudera Live QuickStart VM environment. The top screenshot displays the 'Welcome to Your Cloudera QuickStart VM!' page, which includes a table of cluster nodes and sections for getting started, analyzing data, and managing the cluster. The bottom screenshot shows the Hue File Browser interface, displaying a file listing for the 'cloudera' user.

**Welcome to Your Cloudera QuickStart VM!**

Your Cluster

Node	Address
Manager Node	10.0.2.15
Worker Node 1	10.0.2.15

**Get Started**

The tutorial below guides you through some analytic use cases, using the most popular open source tools included with CDH (including Cloudera Impala, Cloudera Search, and Hue).

[Start Tutorial](#)

**Analyze Your Data**

Hue is the open source web interface for Hadoop that lets you analyze your data. Simply load in your data and then easily begin to analyze, search, and visualize it. In the QuickStart VM, the administrative username for Hue is 'cloudera' and the password is 'cloudera'.

[Launch Hue UI](#)

**Manage Your Cluster**

**Hue - File Browser - Mozilla Firefox**

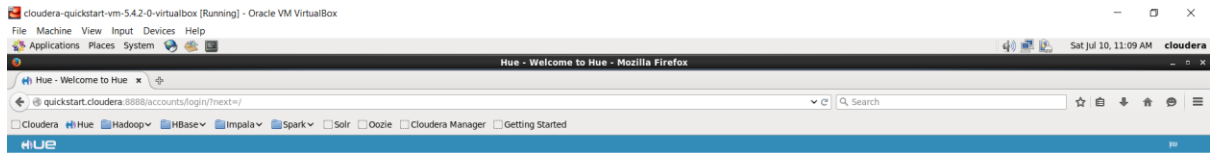
Search for file name

Actions: Move to trash

Home / user / cloudera

Name	Size	User	Group	Permissions	Date
2		hdts	supergroup	drwxr-xr-x	June 09, 2015 03:38 AM
.		cloudera	cloudera	drwxr-xr-x	June 09, 2015 03:37 AM

Page 1 of 1

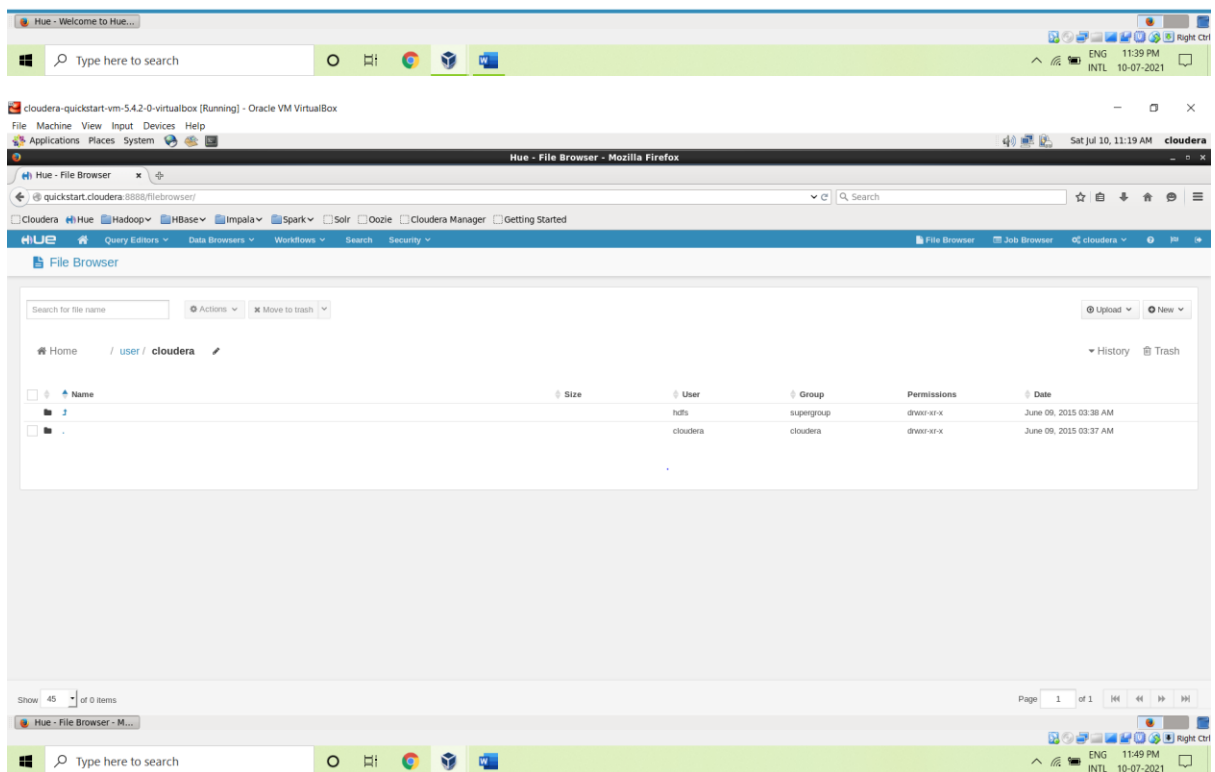


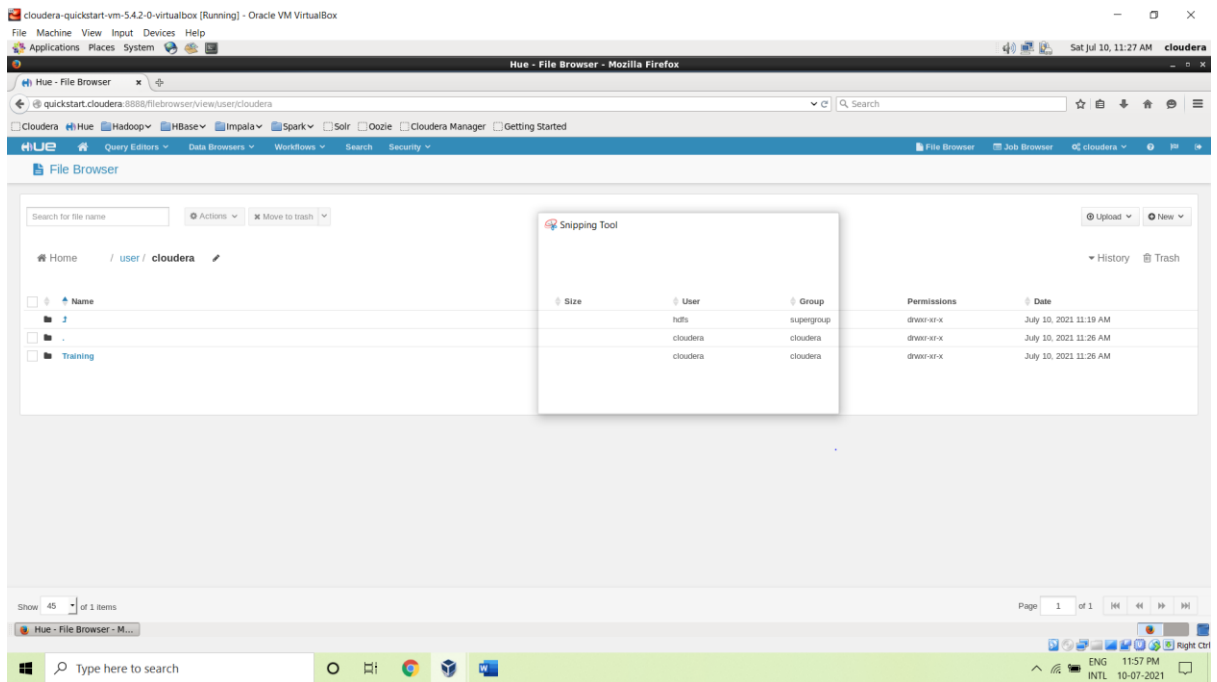
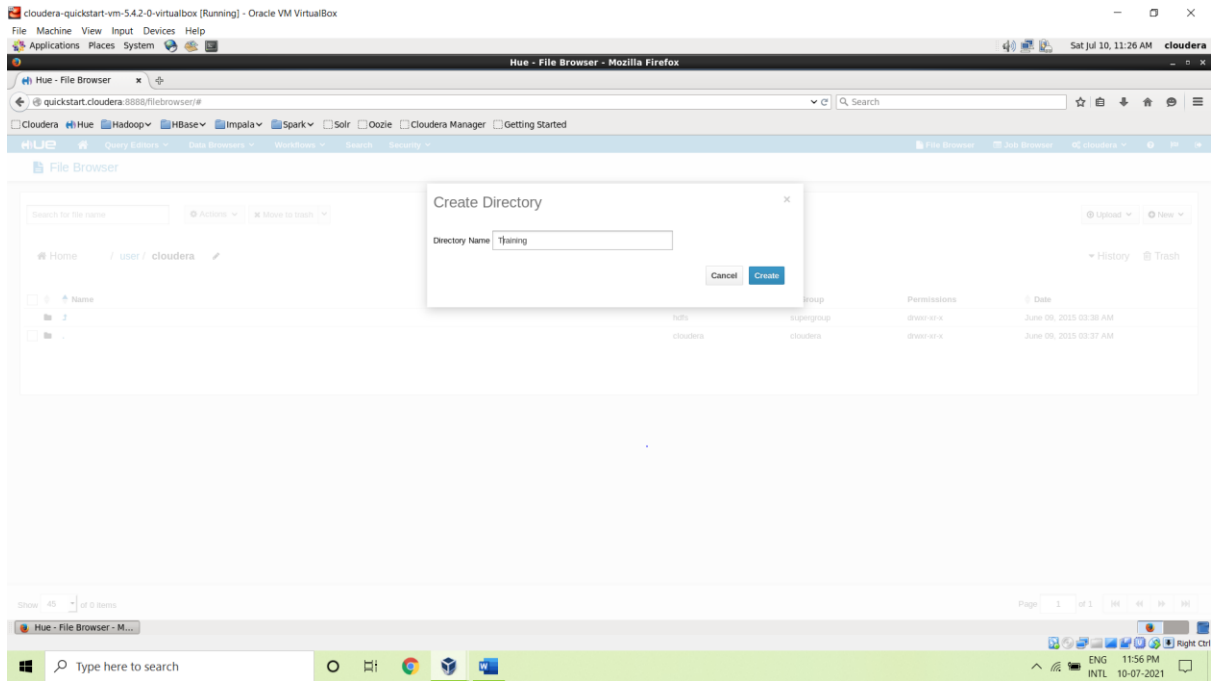
Sign in to continue to Hue

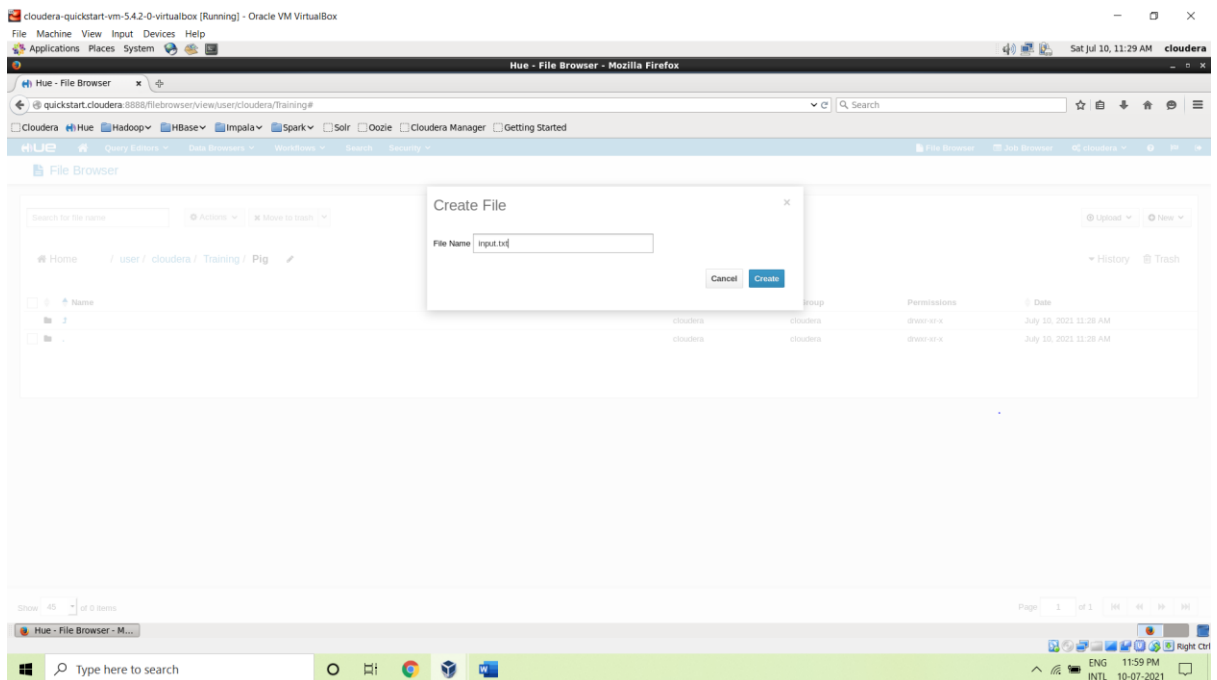
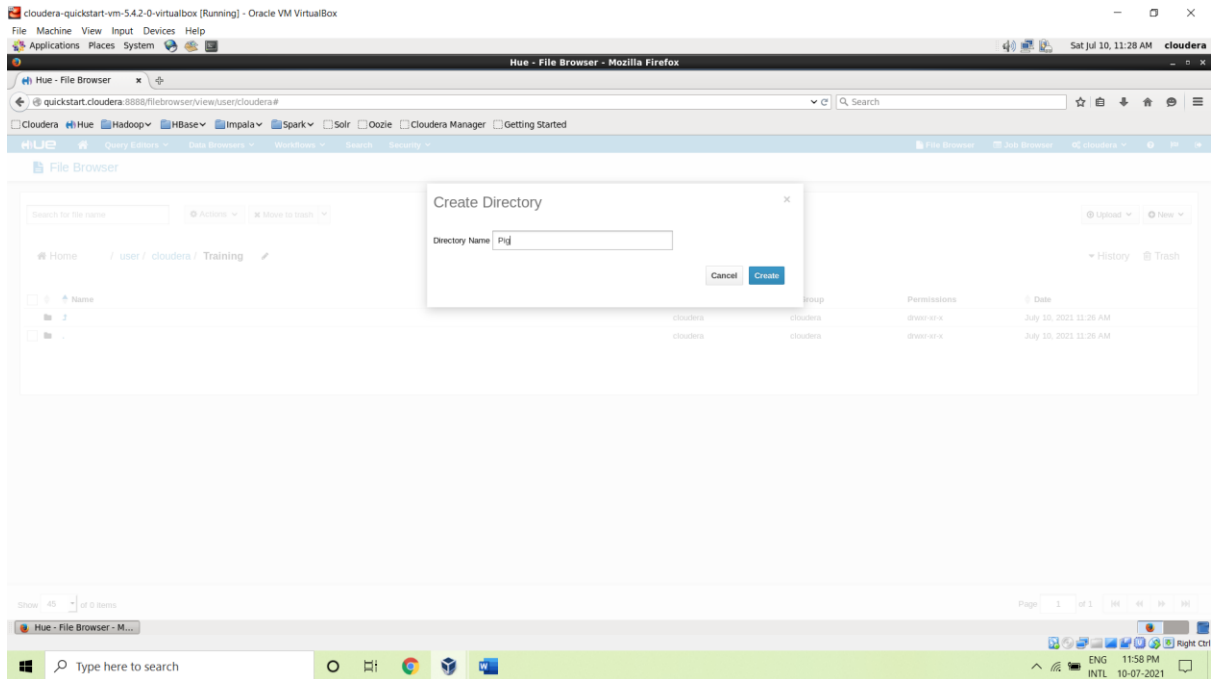
cloudera

\*\*\*\*\*

Sign in







cloudera-quickstart-vm-5.42-0-virtualbox [Running] - Oracle VM VirtualBox

File Machine View Input Devices Help

Applications Places System

Hue - File Browser - input.txt - File Viewer - Mozilla Firefox

Hue - File Browser - I... x

@ quickstart.cloudera:8888/filebrowser/view/user/cloudera/Training/Pig/input.txt

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

HUE Query Editors Data Browsers Workflows Search Security File Browser Job Browser cloudera

File Browser

ACTIONS

- View as binary
- Edit file
- Download
- View file location
- Refresh

INFO

Last modified  
July 10, 2021 11:30 a.m.

User  
cloudera

Group  
cloudera

Size  
0 bytes

Mode  
100644

Home / user / cloudera / Training / Pig / input.txt

Page 1 of 1

The current file is empty.

Hue - File Browser - in...

Type here to search

ENG 12:00 AM  
INTL 11-07-2021

cloudera-quickstart-vm-5.42-0-virtualbox [Running] - Oracle VM VirtualBox

File Machine View Input Devices Help

Applications Places System

Hue - File Browser - input.txt - File Viewer - Mozilla Firefox

Hue - File Browser - I... x

@ quickstart.cloudera:8888/filebrowser/edit/user/cloudera/Training/Pig/input.txt

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

HUE Query Editors Data Browsers Workflows Search Security File Browser Job Browser cloudera

File Browser

ACTIONS

- View file
- Download
- View file location
- Refresh

INFO

Last modified  
July 10, 2021 11:30 a.m.

User  
cloudera

Group  
cloudera

Size  
0 bytes

Mode  
100644

Home / user / cloudera / Training / Pig / input.txt

Big Data is a term used to describe a collection of data that is huge in size and yet growing exponentially with time. But its not the amount of data that's important , Big data can be analyzed for insight that lead to better decisions and strategic business moves.

definitions of Big Data as the V's

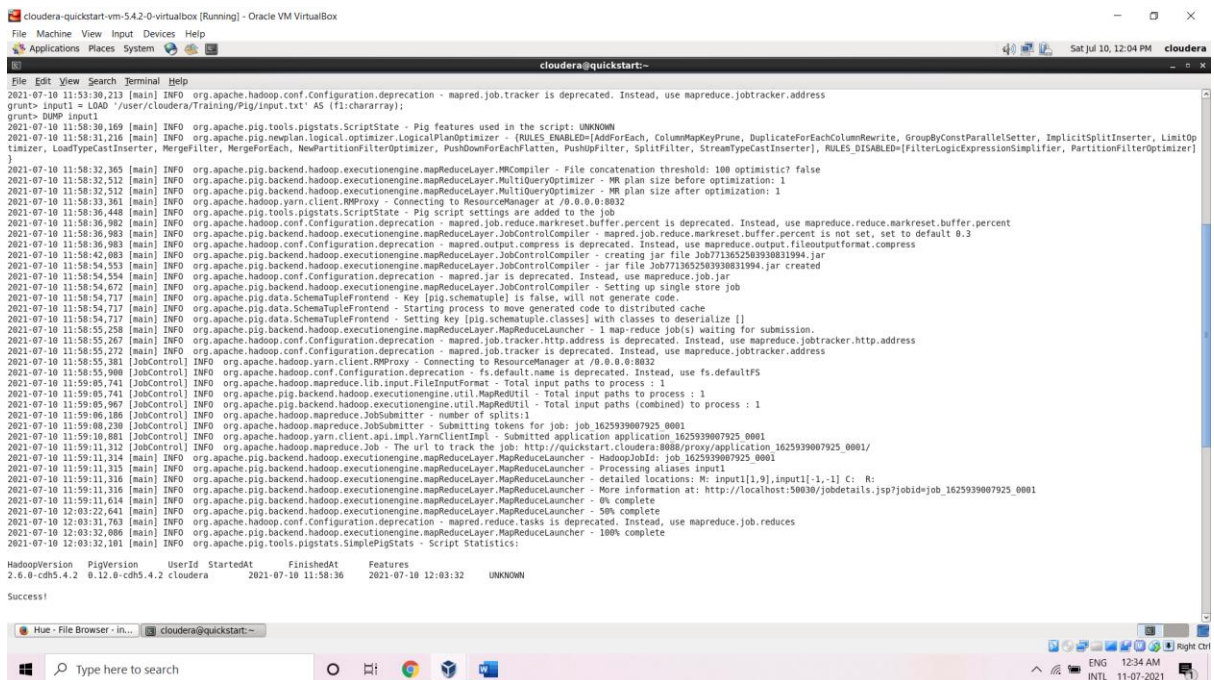
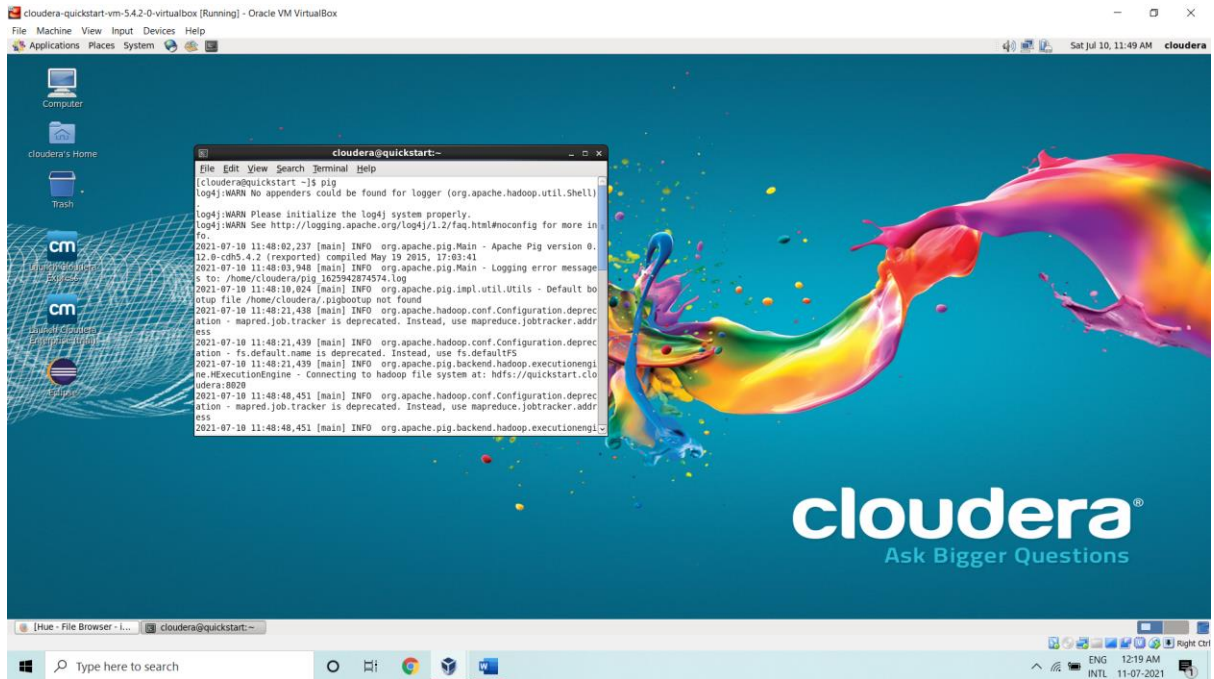
1. Volume
2. Velocity
3. Variety
4. Veracity

Save Save as

Hue - File Browser - in...

Type here to search

ENG 12:11 AM  
INTL 11-07-2021



```
Job Stats (time in seconds):
JobId  Maps  Reduces MaxMapTime  MinMapTime  AvgMapTime  MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReduceTime  Alias  Feature Outputs
job_1625939007925_0001  1  0  70  70  70  70  n/a  n/a  n/a  n/a  input1  MAP_ONLY  hdfs://quickstart.cloudera:8020/tmp/temp-91315653/tmp-1510987010,
```

```
Input(s):
Successfully read 8 records (756 bytes) from: "/user/cloudera/Training/Pig/input.txt"
```

```
Output(s):
Successfully stored 8 records (416 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-91315653/tmp-1510987010"
```

```
Counters:
Total records written : 8
Total bytes written : 416
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
```

```
Job DAG:
job_1625939007925_0001
```

```
2021-07-10 12:03:32,310 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-07-10 12:03:32,449 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-07-10 12:03:32,449 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-07-10 12:03:32,473 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-07-10 12:03:32,842 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-07-10 12:03:32,842 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Big Data is a term used to describe a collection of data that is huge in size and yet growing exponentially with time. )
(But its not the amount of data that's important , Big data can be analyzed for insight that lead to better decisions and strategic business moves.)
()
(definitions of Big Data as the V's )
( 1. Volume)
( 2. Velocity)
( 3. Variety)
( 4. Veracity )
grunt>
```

cloudera-quickstart-vm-5.42-0-virtualbox [Running] - Oracle VM VirtualBox

File Machine View Input Devices Help

Applications Places System

cloudera@quickstart:~

Sat Jul 10, 12:15

File Edit View Search Terminal Help

( 4. Veracity )

grunt: wordInEachLine = FOREACH input1 GENERATE flatten(TOKENIZE(f1)) as word;

2021-07-10 12:11:39,259 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS

2021-07-10 12:11:39,259 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address

grunt: dump wordInEachLine

2021-07-10 12:12:25,531 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: UNKNOWN

2021-07-10 12:12:26,264 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES\_ENABLED=[AddForEach, ColumnMapKeyPrune, DuplicateForEachColumnRewrite, GroupByConstParallelSetter, ImplicitSplitInsertLimiter, LoadTypeCaster, MergeFilter, MergeForEach, NewPartitionFilterOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCasterInsert], RULES\_DISABLED=[FilterLogicExpressionsSimplifier, PartitionFilter]}

2021-07-10 12:12:26,310 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic? false

2021-07-10 12:12:26,320 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1

2021-07-10 12:12:26,321 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1

2021-07-10 12:12:26,388 [main] INFO org.apache.hadoop.yarn.client.NMProxy - Connecting to ResourceManager at /0.0.0.0:8032

2021-07-10 12:12:26,531 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job

2021-07-10 12:12:26,605 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3

2021-07-10 12:12:28,612 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - creating jar file Job6736055733475564116.jar

2021-07-10 12:12:38,327 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - jar file Job6736055733475564116.jar created

2021-07-10 12:12:38,409 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job

2021-07-10 12:12:38,648 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.

2021-07-10 12:12:38,648 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache

2021-07-10 12:12:38,648 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserialize []

2021-07-10 12:12:38,695 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.

2021-07-10 12:12:38,695 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address

2021-07-10 12:12:38,790 [JobControl] INFO org.apache.hadoop.yarn.client.NMProxy - Connecting to ResourceManager at /0.0.0.0:8032

2021-07-10 12:12:38,716 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS

2021-07-10 12:12:39,651 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1

2021-07-10 12:12:39,651 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1

2021-07-10 12:12:39,657 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths (combined) to process : 1

2021-07-10 12:12:39,764 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - number of splits:1

2021-07-10 12:12:40,093 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job\_1625939007925\_0002

2021-07-10 12:12:40,381 [JobControl] INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl - Submitted application application\_1625939007925\_0002

2021-07-10 12:12:40,388 [JobControl] INFO org.apache.hadoop.mapreduce.Job - The url to track the job: http://quickstart.cloudera:8080/proxy/application\_1625939007925\_0002/

2021-07-10 12:12:40,388 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - HadoopJobId: job\_1625939007925\_0002

2021-07-10 12:12:40,388 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases input1,wordInEachLine

2021-07-10 12:12:40,388 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - detailed locations: W: input[1,9],wordInEachLine[1..1] C: R:

2021-07-10 12:12:40,388 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - More information at: http://localhost:50038/jobdetails.jsp?jobid=job\_1625939007925\_0002

2021-07-10 12:12:40,473 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 0% complete

2021-07-10 12:14:35,118 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 50% complete

2021-07-10 12:14:41,874 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete

2021-07-10 12:14:41,875 [main] INFO org.apache.pig.tools.pigstats.SimplePigStats - Script Statistics:

HadoopVersion 2.6.0-cdh5.4.2 PigVersion 0.12.0-cdh5.4.2 UserId cloudera StartedAt 2021-07-10 12:12:26 FinishedAt 2021-07-10 12:14:41 Features UNKNOWN

Success!

Job Stats (time in seconds):

JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Feature Outputs

job\_1625939007925\_0002 1 0 54 54 54 54 n/a n/a n/a n/a input1,wordInEachLine MAP\_ONLY hdfs://quickstart.cloudera:8020/tmp/temp-91315653/tmp217140104,

Hue - File Browser - in-... cloudera@quickstart:~

Type here to search

ENG 12:45



```
cloudera-quickstart-vm-542-0-virtualbox [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
Applications Places System
cloudera@quickstart:~
File Edit View Search Terminal Help

Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Feature Outputs
job_1625939007925_0002 1 0 54 54 54 54 n/a n/a n/a n/a input,wordInEachLine MAP_ONLY hdfs://quickstart.cloudera:8020/tmp/temp-91315653/tmp217140104,

Input(s):
Successfully read 8 records (756 bytes) from: "/user/cloudera/Training/Pig/input.txt"

Output(s):
Successfully stored 64 records (727 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-91315653/tmp217140104"

Counters:
Total records written : 64
Total bytes written : 727
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1625939007925_0002

2021-07-10 12:14:42.328 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-07-10 12:14:42.329 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-07-10 12:14:42.329 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-07-10 12:14:42.337 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2021-07-10 12:14:42.353 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-07-10 12:14:42.353 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Data)
(is)
(a)
(term)
(used)
(to)
(describe)
(a)
(collection)
(of)
(data)
(that)
(is)
(huge)
(in)
(size)
(and)
(yet)
(growing)
(exponentially)

Hue - File Browser - in... cloudera@quickstart:~
Type here to search
```

```
cloudera-quickstart-vm-542-0-virtualbox [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
Applications Places System
cloudera@quickstart:~
File Edit View Search Terminal Help

(Veracity)
grunt> groupedWords = group wordInEachLine by word;
grunt> dump groupedWords;
2021-07-10 12:19:43.274 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP BY
2021-07-10 12:19:43.276 [main] INFO org.apache.pig.nplan.logical.optimizer.LogicalPlanOptimizer - [RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, DuplicateForEachColumnRewrite, GroupByConstParallelSetter, ImplicitSplitInserter, LimitOp
timizer, LoadTypeCastInserter, MergeFilter, MergeForEach, NewPartitionFilterOptimizer, PushDownForEachLatten, PushUpFilter, SplitFilter, StreamTypeCastInserter], RULES_DISABLED=[FilterLogicExpressionSimplifier, PartitionFilterOptimizer]
]
2021-07-10 12:19:43.583 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic? false
2021-07-10 12:19:43.585 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2021-07-10 12:19:43.585 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - Reduce phase detected, estimating # of required reducers.
2021-07-10 12:19:43.658 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at 0.0.0.0:8032
2021-07-10 12:19:43.655 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
2021-07-10 12:19:43.781 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2021-07-10 12:19:43.781 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - BytesPerReducer=1000000000 maxReducers=999 totalInputFileSize=369
2021-07-10 12:19:43.780 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Using reducer estimator: org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputSizeReducerEstimator
2021-07-10 12:19:43.780 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting Parallelism to 1
2021-07-10 12:19:43.784 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
2021-07-10 12:19:56.156 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - creating jar file Job5100193755170024685.jar
2021-07-10 12:19:56.156 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - jar file Job5100193755170024685.jar created
2021-07-10 12:19:56.156 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
2021-07-10 12:19:56.200 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2021-07-10 12:19:56.200 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
2021-07-10 12:19:56.200 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserialize []
2021-07-10 12:19:56.473 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
2021-07-10 12:19:56.477 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-07-10 12:19:56.492 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2021-07-10 12:19:56.954 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-07-10 12:19:56.954 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2021-07-10 12:19:56.959 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths (combined) to process : 1
2021-07-10 12:19:57.265 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - number of splits: 1
2021-07-10 12:19:57.466 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job_1625939007925_0003
2021-07-10 12:19:57.561 [JobControl] INFO org.apache.hadoop.yarn.client.impl.YarnClientImpl - Submitted application application_1625939007925_0003
2021-07-10 12:19:57.570 [JobControl] INFO org.apache.hadoop.mapreduce.Job - The url to track the job: http://quickstart.cloudera:8080/proxy/application_1625939007925_0003/
2021-07-10 12:19:57.570 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - HadoopJobId: job_1625939007925_0003
2021-07-10 12:19:57.571 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases groupedWords,input,wordInEachLine
2021-07-10 12:19:57.571 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Detailed locations: N: input[1,9],wordInEachLine[1..1],groupedWords[3,15] C: R:
2021-07-10 12:19:57.673 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - More information at: http://localhost:50030/jobdetails.jsp?jobid=job_1625939007925_0003
2021-07-10 12:21:15.216 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 0% complete
2021-07-10 12:21:43.889 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 50% complete
2021-07-10 12:21:43.890 [main] INFO org.apache.pig.tools.pigstats.SimplePigStats - Script Statistics:

HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.6.0-cdh5.4.2 0.12.0-cdh5.4.2 cloudera 2021-07-10 12:19:43 2021-07-10 12:21:43 GROUP BY

Success!

Job Stats (time in seconds):

Hue - File Browser - in... cloudera@quickstart:~
Type here to search
```



```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
(3.,{(3.)})  
(4.,{(4.)})  
(as,{(as)})  
(be,{(be)})  
(in,{(in)})  
(is,{(is),(is)})  
(of,{(of),(of),(of)})  
(to,{(to),(to)})  
(Big,{(Big),(Big),(Big)})  
(But,{(But)})  
(V's,{(V's)})  
(and,{(and),(and)})  
(can,{(can)})  
(for,{(for)})  
(its,{(its)})  
(not,{(not)})  
(the,{(the),(the)})  
(yet,{(yet)})  
(Data,{(Data),(Data)})  
(data,{(data),(data),(data)})  
(huge,{(huge)})  
(lead,{(lead)})  
(size,{(size)})  
(term,{(term)})  
(that,{(that),(that)})  
(used,{(used)})  
(with,{(with)})  
(time.,{(time.)})  
(Volume,{(Volume)})  
(amount,{(amount)})  
(better,{(better)})  
(moves.,{(moves.)})  
(that's,{(that's)})  
(Variety,{(Variety)})  
(growing,{(growing)})  
(insight,{(insight)})  
(Velocity,{(Velocity)})  
(Veracity,{(Veracity)})  
(analyzed,{(analyzed)})  
(business,{(business)})  
(describe,{(describe)})  
(decisions,{(decisions)})  
(important,{(important)})  
(strategic,{(strategic)})  
(collection,{(collection)})  
(definitions,{(definitions)})  
(exponentially,{(exponentially)})  
(,{{{}}})  
grunt>
```

```
cloudera-quickstart-vm-5.4.2-0-virtualbox [Running] - Oracle VM VirtualBox  
File Machine View Input Devices Help  
Applications Places System  
cloudera@quickstart:~  
File Edit View Search Terminal Help  
(exponentially,{(exponentially)})  
(,{{{}}})  
grunt> countedWords = foreach groupedWords generate group, COUNT(wordInEachLine);  
grunt> dump countedWords;  
2021-07-10 12:25:44,631 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP_BY  
2021-07-10 12:25:45,635 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - {RULES_ENABLED={AddForEach, ColumnMapKeyPrune, DuplicateForEachColumnRewrite, GroupByConstParallelSetter, Impl  
timizer, LoadTypeCastInserter, MergeFilter, MergeForEach, NewPartitionFilterOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter}, RULES_DISABLED={FilterLogicExpressionSimplifier,  
}  
2021-07-10 12:25:44,653 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic? false  
2021-07-10 12:25:44,700 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.CombinerOptimizer - Choosing to move algebraic foreach to combiner  
2021-07-10 12:25:45,058 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1  
2021-07-10 12:25:45,058 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1  
2021-07-10 12:25:45,124 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032  
2021-07-10 12:25:45,130 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job  
2021-07-10 12:25:45,179 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3  
2021-07-10 12:25:45,180 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reduce phase detected, estimating # of required reducers.  
2021-07-10 12:25:45,180 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Using reducer estimator: org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputS  
2021-07-10 12:25:45,191 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputSizeReducerEstimator - BytesPerReducer=1000000000 maxReducers=999 totalInputFileSize=369  
2021-07-10 12:25:45,192 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting Parallelism to 1  
2021-07-10 12:25:46,704 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - creating jar file Job427095258085334023.jar  
2021-07-10 12:25:56,315 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - jar file Job427095258085334023.jar created  
2021-07-10 12:25:56,360 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job  
2021-07-10 12:25:56,370 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.  
2021-07-10 12:25:56,370 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache  
2021-07-10 12:25:56,370 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserialize []  
2021-07-10 12:25:56,483 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.  
2021-07-10 12:25:56,483 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address  
2021-07-10 12:25:56,487 [JobControl] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032  
2021-07-10 12:25:56,499 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS  
2021-07-10 12:25:56,916 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1  
2021-07-10 12:25:56,916 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1  
2021-07-10 12:25:56,922 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths (combined) to process : 1  
2021-07-10 12:25:57,022 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - number of splits:1  
2021-07-10 12:25:57,318 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job_1625939007925_0004  
2021-07-10 12:25:57,467 [JobControl] INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl - Submitted application application_1625939007925_0004  
2021-07-10 12:25:57,483 [JobControl] INFO org.apache.hadoop.mapreduce.Job - The url to track the job: http://quickstart.cloudera:8080/proxy/application_1625939007925_0004/  
2021-07-10 12:25:57,483 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - HadoopJobId: job_1625939007925_0004  
2021-07-10 12:25:57,483 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases countedWords,groupedWords,input1,wordInEachLine  
2021-07-10 12:25:57,483 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - detailed locations: R: input1[1,9],wordInEachLine[-1,-1],countedWords[4,15],groupedWords[3,15]  
2021-07-10 12:25:57,483 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - countedWords[3,15] R: countedWords[4,15]  
2021-07-10 12:25:57,483 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - More information at: http://localhost:50030/jobdetails.jsp?jobid=job_1625939007925_0004  
2021-07-10 12:25:57,636 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 0% complete  
2021-07-10 12:27:09,118 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 50% complete  
2021-07-10 12:27:44,238 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete  
2021-07-10 12:27:44,239 [main] INFO org.apache.pig.tools.pigstats.SimplePigStats - Script Statistics:  
  
HadoopVersion PigVersion UserId StartedAt FinishedAt Features  
2.6.0-cdh5.4.2 0.12.0-cdh5.4.2 cloudera 2021-07-10 12:25:45 2021-07-10 12:27:44 GROUP_BY
```

cloudera-quickstart-vm-5.4.2-0-virtualbox [Running] - Oracle VM VirtualBox

File Machine View Input Devices Help

Applications Places System

cloudera@quickstart:~

File Edit View Search Terminal Help

```
(3.,1)
(4.,1)
(as,1)
(be,1)
(in,1)
(is,2)
(of,3)
(to,2)
(Big,3)
(But,1)
(V's,1)
(and,2)
(can,1)
(for,1)
(its,1)
(not,1)
(the,2)
(yet,1)
(Data,2)
(data,3)
(huge,1)
(lead,1)
(size,1)
(term,1)
(that,2)
(used,1)
(with,1)
(time.,1)
(Volume,1)
(amount,1)
(better,1)
(moves.,1)
(that's,1)
(Variety,1)
(growing,1)
(insight,1)
(Velocity,1)
(Veracity,1)
(analyzed,1)
(business,1)
(describe,1)
(decisions,1)
(important,1)
(strategic,1)
(collection,1)
(definitions,1)
(exponentially,1)
(,0)
grunt>
```

Hue - File Browser - in...

cloudera@quickstart:~