

Name:

Roll No:

## PRACTICAL 9

### To implement Word Count problem using Pig

#### Apache Pig

- **Apache Pig** is a platform for analyzing large data sets that consists of a high-level language for expressing data analysis programs, coupled with infrastructure for evaluating these programs. The salient property of Pig programs is that their structure is amenable to substantial parallelization, which in turns enables them to handle very large data sets.
  - The language used for Pig is Pig Latin. The Pig scripts get internally converted to Map Reduce jobs and get executed on data stored in HDFS.
  - Apart from that, Pig can also execute its job in Apache Tez or Apache Spark.
  - Pig can handle any type of data, i.e., structured, semi-structured or unstructured and stores the corresponding results into Hadoop Data File System.
  - Every task which can be achieved using PIG can also be achieved using java used in MapReduce.
- Pig's language layer currently consists of a textual language called Pig Latin, which has the following key properties:
- **Ease of programming.** It is trivial to achieve parallel execution of simple, "embarrassingly parallel" data analysis tasks. Complex tasks comprised of multiple interrelated data transformations are explicitly encoded as data flow sequences, making them easy to write, understand, and maintain.
  - **Optimization opportunities.** The way in which tasks are encoded permits the system to optimize their execution automatically, allowing the user to focus on semantics rather than efficiency.
  - **Extensibility.** Users can create their own functions to do special-purpose processing.

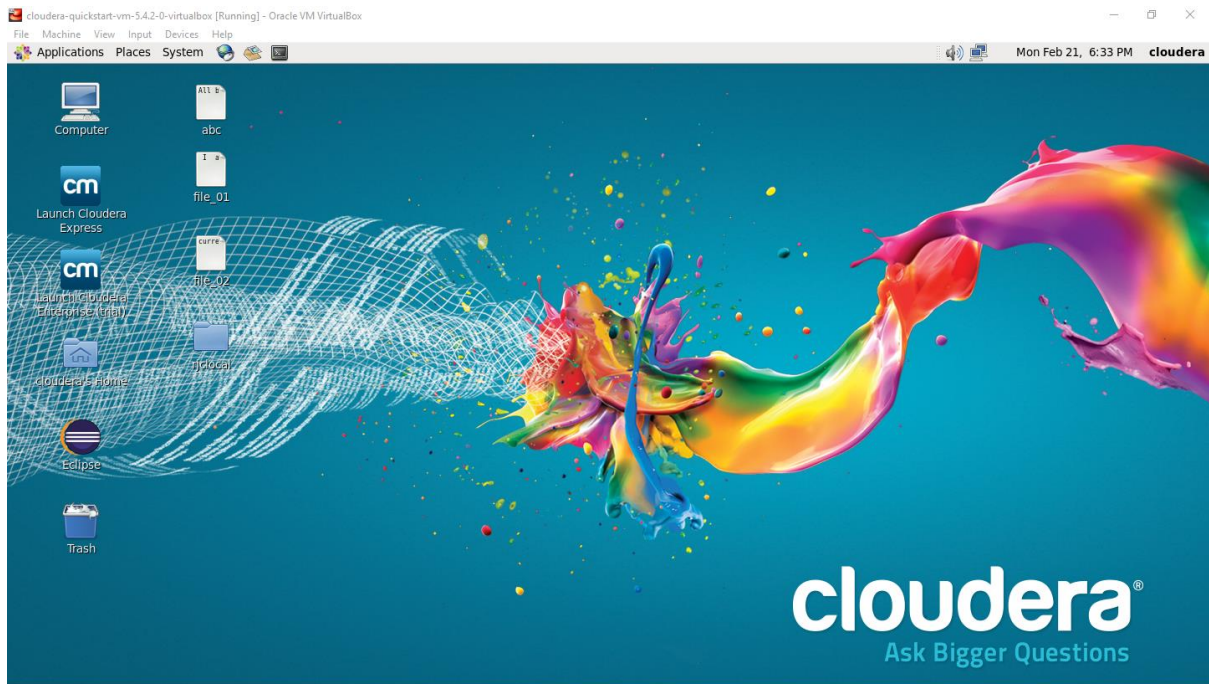
Name:

Roll No:

## **To implement Word Count problem using Pig**

### **Steps:**

1) Start the cloudera.



2) Open the browser. And then open Hue and login.

Name:

Roll No:



Sign in to continue to Hue

3) In Hue Go to file browser and Now open the directory /user/cloudera

quickstart.cloudera:8888/filebrowser/

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started Pract 9 - Pig

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

File Browser

Search for file name Actions Move to trash Upload New

Home / user / cloudera History Trash

	Name	Size	User	Group	Permissions	Date
	↓		hdfs	supergroup	drwxr-xr-x	March 10, 2022 08:19 PM
	.		cloudera	cloudera	drwxr-xr-x	March 10, 2022 08:42 PM
	Desktop		cloudera	cloudera	drwxr-xr-x	February 14, 2022 07:56 PM

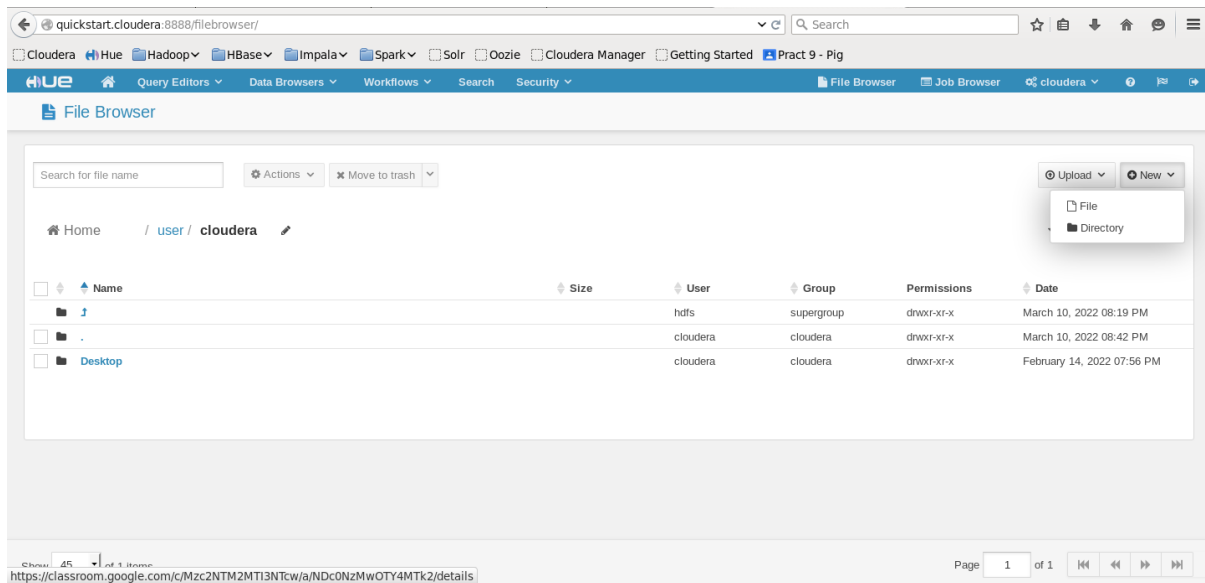
Show 45 of 1 Items Page 1 of 1

Hue - File Browser - ... cloudera@quickstar... [Experiment 11 - Sc... cloudera@quickstar... /root@quickstart/h... cloudera@quickstar... Implementation of ...

4) Now we are creating the directory as Training inside /user/cloudera

Name:

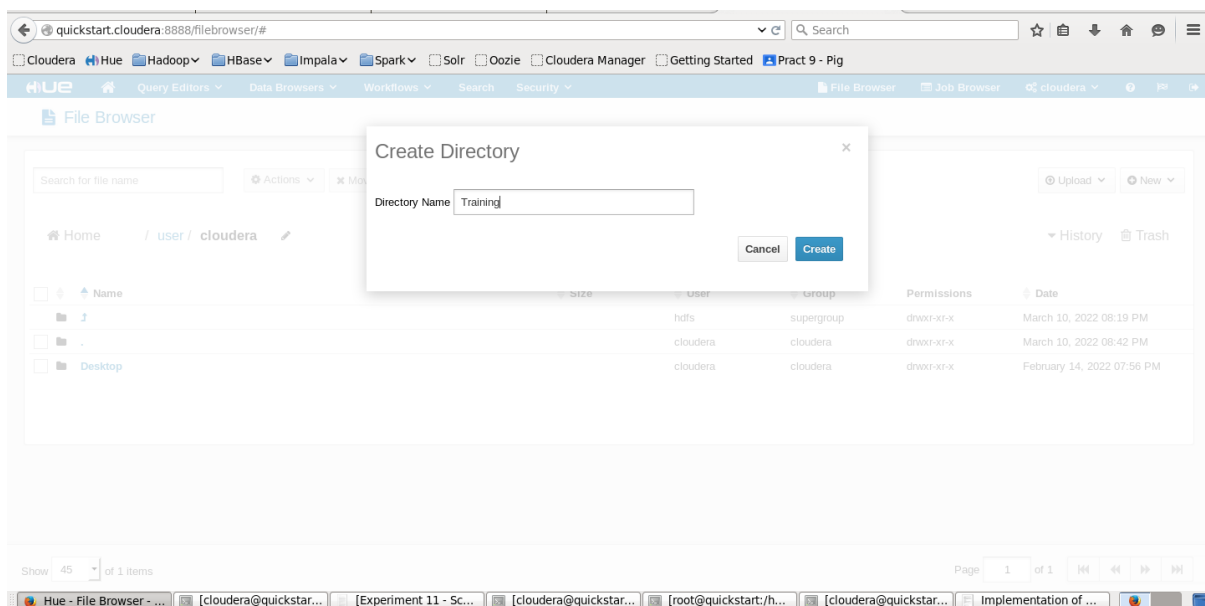
Roll No:



In File Browser we have New option in right corner

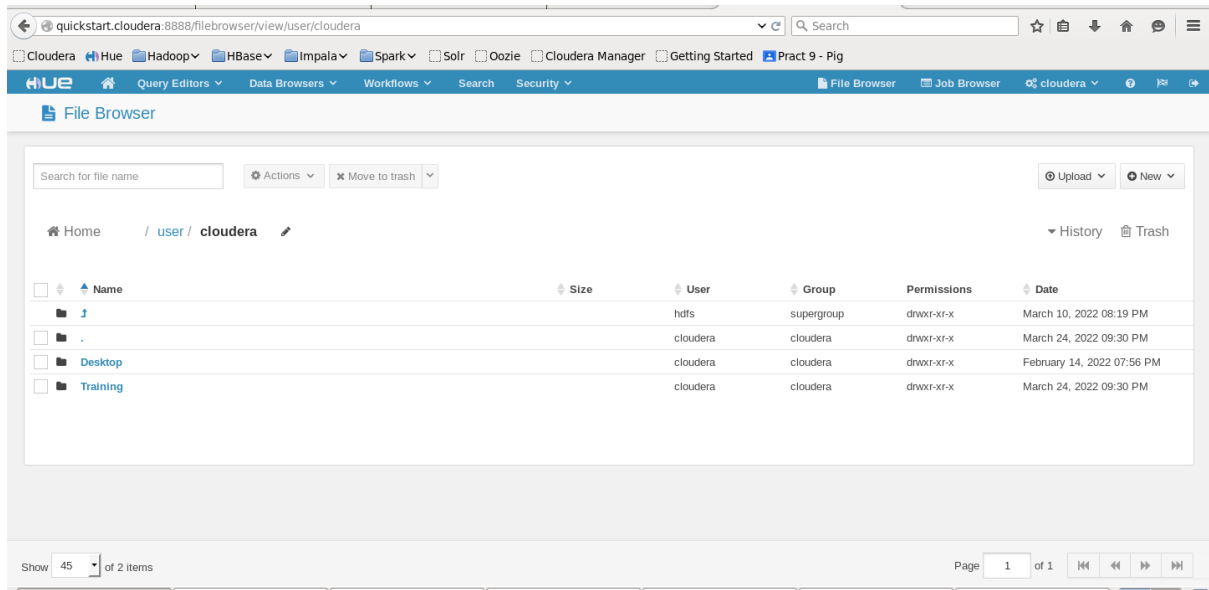
Click on New → Directory

Give the directory name And click on Create

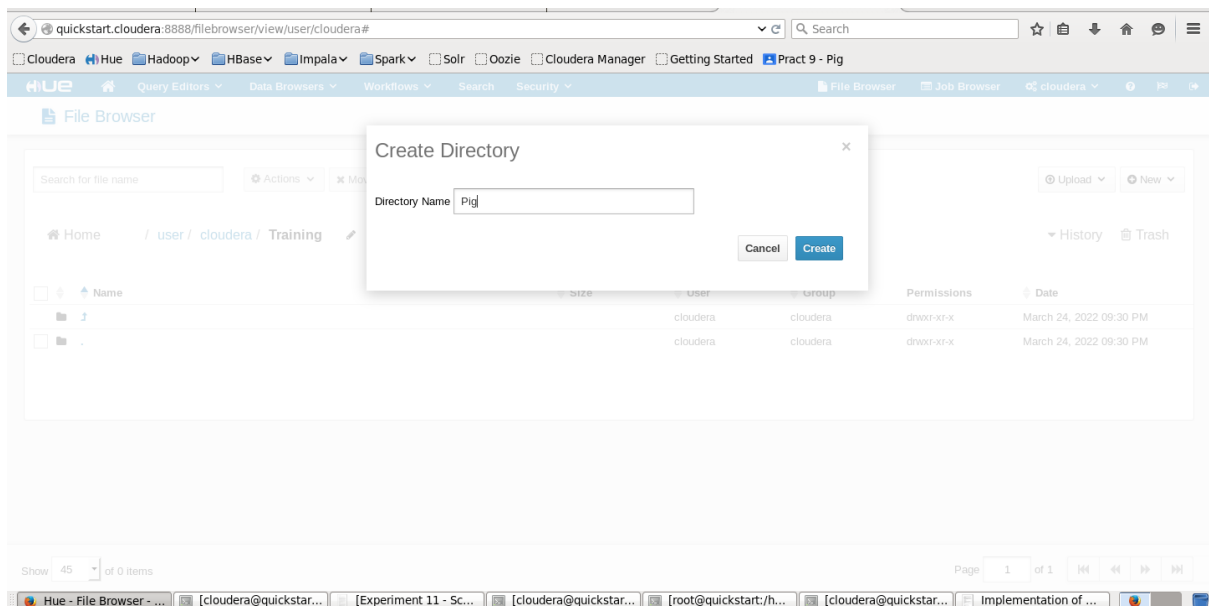


Name:

Roll No:



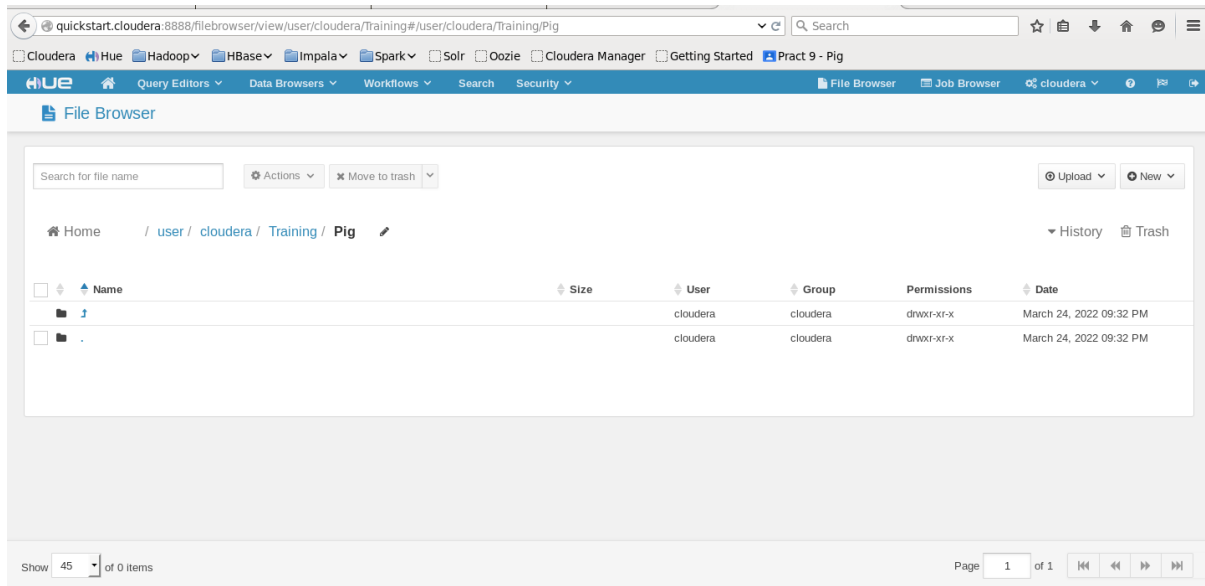
5) After creating Training directory now creating the Pig directory inside Training.



6) Pig directory has been created inside /user/cloudera/Training

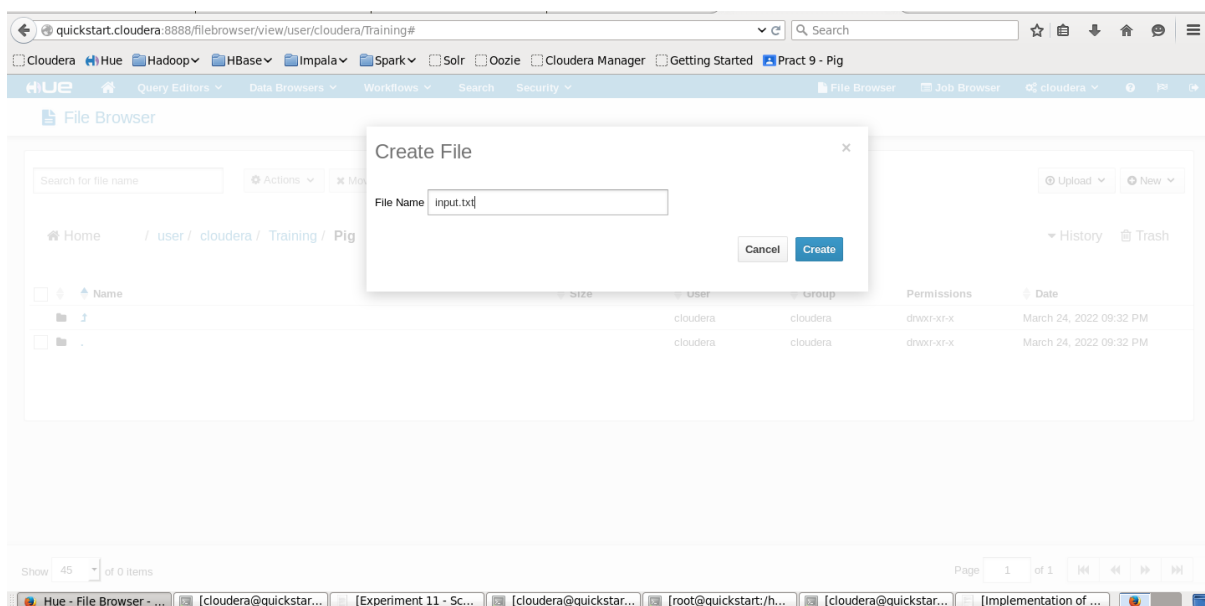
Name:

Roll No:



7) Creating input.txt file inside /usr/cloudera/Training/Pig directory

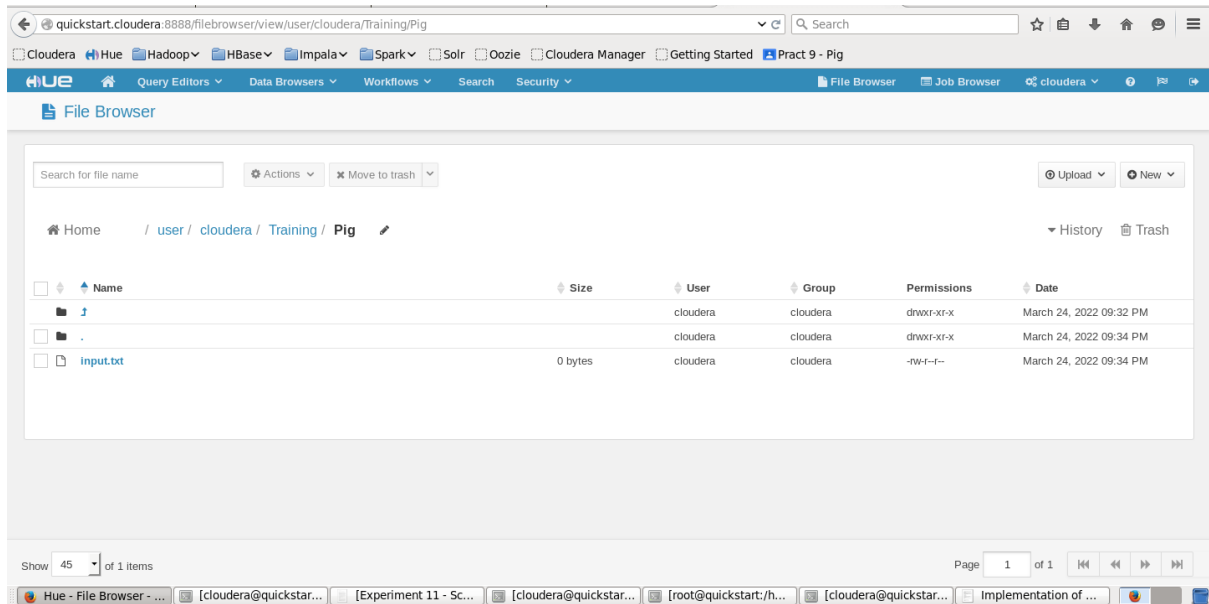
Again inside the Pig directory click on New and create file as 'input.txt'



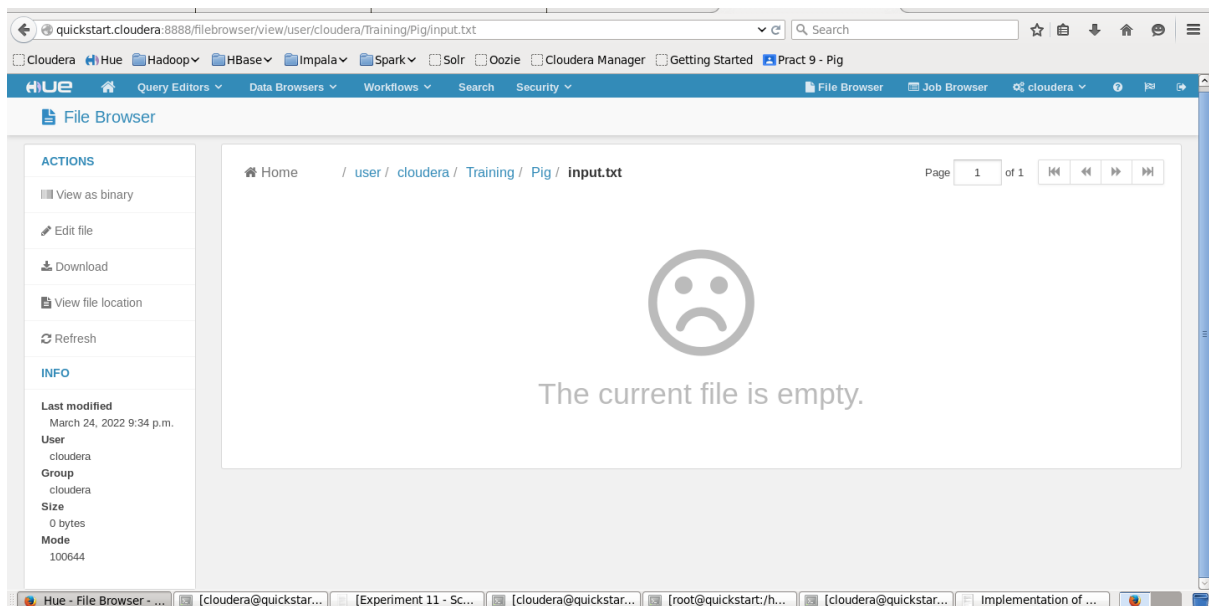
Once the file has been created click on 'input.txt' to add the content in it

Name:

Roll No:



8) Adding some contents to this input.txt file.

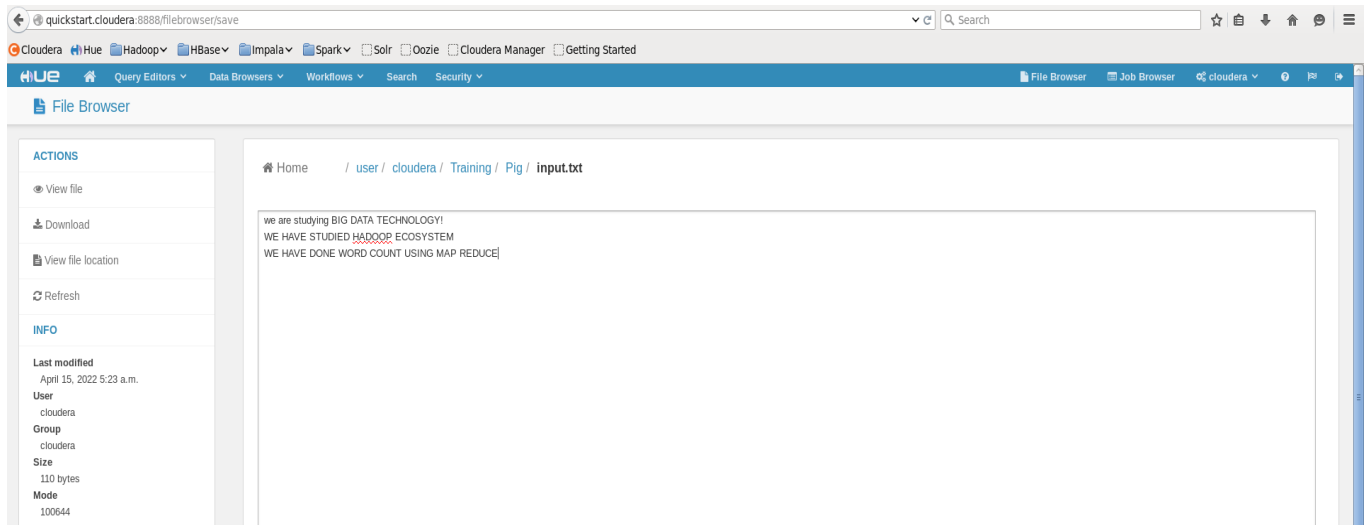


For adding content in the input file, Click on 'Edit file' option then add the content.

Save the input.txt file

Name:

Roll No:



9) Now Open the terminal. And start Pig by typing pig on terminal.



Name:

Roll No:

```
[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.
2022-03-24 21:27:30,534 [main] INFO org.apache.pig.Main - Apache Pig version 0.12.0-cdh5.4.2 (reexported) compiled May 19 2015, 17:03:41
2022-03-24 21:27:30,534 [main] INFO org.apache.pig.Main - Logging error messages to: /home/cloudera/pig_1648182450505.log
2022-03-24 21:27:30,564 [main] INFO org.apache.pig.impl.util.Utils - Default bootstrap file /home/cloudera/.pigbootstrap not found
2022-03-24 21:27:31,410 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-24 21:27:31,411 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-24 21:27:31,411 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://quickstart.cloudera:8020
2022-03-24 21:27:33,409 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-24 21:27:33,409 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to map-reduce job tracker at: localhost:8021
2022-03-24 21:27:33,415 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-24 21:27:33,466 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-24 21:27:33,468 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-24 21:27:33,517 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-24 21:27:33,517 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-24 21:27:33,571 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-24 21:27:33,571 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-24 21:27:33,670 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
```

Roll No:

10) Now we have to load that input file where ever it is stored. By typing the command

```
Input1 = LOAD '/usr/cloudera/Training/pig/input.txt' AS (f1:chararray);
grunt> Input1 = LOAD '/usr/cloudera/Training/pig/input.txt' AS (f1:chararray);
grunt>
```

DUMP input1;

```

2022-03-24 21:34:49.873 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: UNKNOWN
2022-03-24 21:34:49.874 [main] INFO org.apache.hadoop.mapreduce.lib.input.TextInputFormat - FILEFS_ENABLED=false, ColumnsMatchPrune, DuplicatesForEachColumnHeader, GroupConcatParallelSetter, ImplicitSplitter, Limit
2022-03-24 21:34:49.873 [main] INFO org.apache.hadoop.mapreduce.lib.output.TextOutputFormat - MergeFilter, MergeFilter, NewPartitionFilter, PushDownForEachJoin, PushDownFilter, SplitFilter, StreamTypeCastInsertion, FILES_DISK_READ=FilterIgnoreExpressionsImplicit, PartitionFilter, Limit
2022-03-24 21:34:49.875 [main] INFO org.apache.hadoop.mapreduce.executionengine.mapreducelayer.MRCCompiler - File compilation threshold: 100 optimistic? false
2022-03-24 21:34:49.876 [main] INFO org.apache.hadoop.mapreduce.executionengine.mapreducelayer.MultiQueryOptimizer - MR plan size before optimization: 1
2022-03-24 21:34:49.876 [main] INFO org.apache.hadoop.mapreduce.executionengine.mapreducelayer.MultiQueryOptimizer - MR plan size after optimization: 1
2022-03-24 21:34:49.900 [main] INFO org.apache.hadoop.yarn.client.Proxy - Connecting to ResourceManager at /0.0.0.0:8032
2022-03-24 21:34:49.900 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
2022-03-24 21:34:49.905 [main] INFO org.apache.hadoop.mapreduce.executionengine.mapreducelayer.JobControlCompiler - mpared, job.reduce.markreset.buflen=99999999 is not set, set to default 0.3
2022-03-24 21:34:49.958 [main] INFO org.apache.hadoop.mapreduce.executionengine.mapreducelayer.JobControlCompiler - creating jar file Job0515548170405991125.jar
2022-03-24 21:34:49.958 [main] INFO org.apache.hadoop.mapreduce.executionengine.mapreducelayer.JobControlCompiler - jar file Job0515548170405991125.jar created
2022-03-24 21:34:50.060 [main] INFO org.apache.hadoop.mapreduce.executionengine.mapreducelayer.JobControlCompiler - Setting up single store job
2022-03-24 21:34:53.060 [main] INFO org.apache.pig.data.SchemaTempFrontend - Key [pig:schematemp] is false, will not generate code.
2022-03-24 21:34:53.060 [main] INFO org.apache.pig.data.SchemaTempFrontend - Starting process to move generated code to distributed cache
2022-03-24 21:34:53.060 [main] INFO org.apache.pig.data.SchemaTempFrontend - Setting key [pig:schematemp-classes] with classes to deserialize []
2022-03-24 21:34:53.067 [main] INFO org.apache.hadoop.mapreduce.executionengine.mapreducelayer.Launcher - 1 map reduce job(s) waiting for submission.
2022-03-24 21:34:53.067 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mpared, job.tracker is deprecated. Instead, use mapreduce.job.tracker.address
2022-03-24 21:34:53.067 [main] INFO org.apache.hadoop.yarn.client.Proxy - Connecting to ResourceManager at /0.0.0.0:8032
2022-03-24 21:34:53.074 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-24 21:34:53.270 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2022-03-24 21:34:53.271 [JobControl] INFO org.apache.pig.hadoop.executionengine.mapreducelayer.JobControl - Total input paths to process : 1
2022-03-24 21:34:53.278 [JobControl] INFO org.apache.pig.hadoop.executionengine.mapreducelayer.JobControl - Total input paths (combined) to process : 1
2022-03-24 21:34:53.331 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - number of splits: 1
2022-03-24 21:34:53.799 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job_1645434315262_0028
2022-03-24 21:34:53.800 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting application_1645434315262_0028

```

Roll No:

12) Here we are counting the words in each line for that we are using the following command

```

grunts> wordsInEachLine = FOREACH Input1 GENERATE flatten(TOKENIZE(f1)) as word;
2022-03-25 22:09:29,072 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-25 22:09:29,072 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
grunts>

```

```

2022-03-25 22:09:29,072 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-25 22:09:29,072 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-25 22:09:29,072 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-25 22:10:40,971 [main] INFO org.apache.pig.tools.pigstats.ScriptStats - Pig features used in the script: UNKNOWN
2022-03-25 22:10:40,971 [main] INFO org.apache.pig.mapred.logical.optimizer.LocalPlanOptimizer - RULES ENABLED:AdaptFilter,ColumnMajorGroupPrune, DuplicateForEachColumnRewrite, GroupConcatParallelSetter, ImplicitSplitInserter, LimitBy
2022-03-25 22:10:40,971 [main] INFO org.apache.pig.mapred.logical.optimizer.LocalPlanOptimizer - RULES DISABLED:FilterConcatExpressionsInserter, PartitionFilterOptimizer,
2022-03-25 22:10:40,971 [main] INFO org.apache.pig.mapred.logical.optimizer.LocalPlanOptimizer - Filter MergeRefactor, NewTableOptimizer, PushDownJoinInserter, PushDownJoinSplitter, StreamConcatInserter,
2022-03-25 22:10:41,014 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MRCompiler - File concatenation threshold: 160 optimistic: false
2022-03-25 22:10:41,016 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2022-03-25 22:10:41,016 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1
2022-03-25 22:10:41,062 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2022-03-25 22:10:41,062 [main] INFO org.apache.pig.tools.pigstats.ScriptStats - Pig script settings are added to the job
2022-03-25 22:10:41,084 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.JobControlCompiler - mapred.job.reduce.markrest.buffer.percent is not set, set to default 0.3
2022-03-25 22:10:42,118 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.JobControlCompiler - creating jar file job065080738905845951.jar
2022-03-25 22:10:42,118 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.JobControlCompiler - jar file job065080738905845951.jar created
2022-03-25 22:10:43,336 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.JobControlCompiler - Setting up single store job
2022-03-25 22:10:45,537 [main] INFO org.apache.pig.data.SchemaTableFrontend - Key [pig.schematuple] is false, will not generate code.
2022-03-25 22:10:45,537 [main] INFO org.apache.pig.data.SchemaTableFrontend - Starting process to move generated code to distributed cache
2022-03-25 22:10:45,537 [main] INFO org.apache.pig.data.SchemaTableFrontend - Setting key [pig.schematuple.classes] with classes to deserialize []
2022-03-25 22:10:45,546 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MapReduceLauncher - 1 map-reduce jobs (waiting for submission.
2022-03-25 22:10:45,572 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-25 22:10:45,572 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-25 22:10:45,589 [JobControl] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2022-03-25 22:10:45,589 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-25 22:10:45,589 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-25 22:10:45,776 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
2022-03-25 22:10:45,778 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths (combined) to process : 1
2022-03-25 22:10:45,780 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Submitting tokens for job: job 1644548343526 0029
2022-03-25 22:10:45,928 [JobControl] INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl - Submitted application application 1644548343526 0029
2022-03-25 22:10:46,009 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MapReduceLauncher - http://quickstart:8030/proxy/application 1644548343526 0029/
2022-03-25 22:10:46,046 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MapReduceLauncher - http://localhost:8030/job/ 1644548343526 0029/
2022-03-25 22:10:46,046 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MapReduceLauncher - Develloped aliases Input,wordsInEachLine
2022-03-25 22:10:46,046 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MapReduceLauncher - Percolated locations: M Input[1],E wordsInEachLine[1]-1,C M:
2022-03-25 22:10:46,046 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MapReduceLauncher - More information at: http://localhost:8080/jobdetails_1644548343526 0029
2022-03-25 22:10:46,100 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MapReduceLauncher - 0% complete
2022-03-25 22:10:46,100 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MapReduceLauncher - 50% complete
2022-03-25 22:10:46,100 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MapReduceLauncher - 100% complete
2022-03-25 22:11:11,459 [main] INFO org.apache.pig.tools.pigstats.SimplePigStats - Script Statistics:
HadoopVersion PigVersion UserStd StartedAt FinishedAt Features
2.6.0-cdh.4.2 0.12.0-cdh.4.2 2 cloudera 2022-03-25 22:10:41 2022-03-25 22:11:11 UNKNOWN

```



Name:

Roll No:

```
Job Stats (time in seconds):
JobId  Maps  Reduces  MaxMapTime  MinMapTime  AvgMapTime  MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReduceTime  Alias  Feature Outputs
job_1644548343526_0029  1  0  5  5  5  5  n/a  n/a  n/a  n/a  Input1,wordsInEachLine  MAP_ONLY  hdfs://quickstart.cloudera:8020/tmp/temp-669075149/tmp-341911088,

Input(s):
Successfully read 3 records (497 bytes) from: "/user/cloudera/Training/pig/input.txt"

Output(s):
Successfully stored 19 records (225 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-669075149/tmp-341911088"

Counters:
Total records written : 19
Total bytes written : 225
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1644548343526_0029

2022-03-25 22:12:11.537 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2022-03-25 22:12:11.537 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-25 22:12:11.537 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-25 22:12:11.538 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] was not set... will not generate code.
2022-03-25 22:12:11.543 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2022-03-25 22:12:11.543 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
(are)
(studying)
(BTG)
(DATA)
(TECHNOLOGY)
(WE)
(HAVE)
(STUDIED)
(HADDOOP)
(ECOSYSTEM)
(WE)
(HAVE)
(DONE)
(WORD)
(COUNT)
(USING)
(WAS)
(REDUCE)
grunt>
```

- 14) Now grouping the words present in each line.  
groupedWords = group wordsInEachLine by word;

```
grunt> groupedWords = group wordsInEachLine by word;
grunt>
```

And then dumping the data by the following command.  
dump groupedWords;

```
grunt> dump groupedWords;
2022-03-25 22:12:59.694 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP BY
2022-03-25 22:12:59.695 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - RULES_ENABLED=[AddForEach, ColumnKeyPrune, DuplicateForEachColumnWrite, GroupByConstParallelSetter, ImplicitSplitInserter, LimitFilter, LoadTypeCaster, MergeFilter, MergeForEach, NewPartitionFilterOptimizer, PushDownForEachFlatten, PushJoinFilter, SplitFilter, StreamTypeCaster], RULES_DISABLED=[FilterLogicExpressionsSimplifier, PartitionFilterOptimizer]
2022-03-25 22:12:59.702 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic? false
2022-03-25 22:12:59.706 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2022-03-25 22:12:59.706 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1
2022-03-25 22:12:59.736 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2022-03-25 22:12:59.738 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig script settings are added to the job
2022-03-25 22:12:59.750 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2022-03-25 22:12:59.750 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reduce phase detected, estimating # of required reducers.
2022-03-25 22:12:59.752 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Using reducer estimator: org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputSizeReducerEstimator
2022-03-25 22:12:59.757 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputSizeReducerEstimator - BytesPerReducer=1000000000 maxReducers=999 totalInputFileSize=110
2022-03-25 22:12:59.757 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting Parallelism to 1
2022-03-25 22:12:59.760 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - creating jar file Job016439686515978146.jar
2022-03-25 22:12:59.767 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - jar file Job016439686515978146.jar created
2022-03-25 22:13:03.573 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
2022-03-25 22:13:03.574 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2022-03-25 22:13:03.574 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
2022-03-25 22:13:03.574 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserialize []
2022-03-25 22:13:03.676 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
2022-03-25 22:13:03.676 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-25 22:13:03.677 [JobControl] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2022-03-25 22:13:03.687 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-25 22:13:03.825 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2022-03-25 22:13:03.825 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths to process : 1
2022-03-25 22:13:03.829 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths (combined) to process : 1
2022-03-25 22:13:03.858 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapReduceUtil - Total input paths (combined) to process : 1
2022-03-25 22:13:03.985 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job_1644548343526_0030
2022-03-25 22:13:03.991 [JobControl] INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl - Submitted application application_1644548343526_0030
2022-03-25 22:13:03.993 [JobControl] INFO org.apache.hadoop.mapreduce.Job - The url to track the job: http://quickstart.cloudera:8080/proxy/application_1644548343526_0030/
2022-03-25 22:13:04.176 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - HadoopJobId: job_1644548343526_0030
2022-03-25 22:13:04.176 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases Input1,groupedWords,wordsInEachLine
2022-03-25 22:13:04.176 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - detailed locations: W: Input[1:9],wordsInEachLine[1-11],groupedWords[5,15] C: R:
2022-03-25 22:13:04.176 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - More information at: http://localhost:50030/jobdetails.jsp?jobid=job_1644548343526_0030
2022-03-25 22:13:04.246 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 0% complete
2022-03-25 22:13:20.489 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 50% complete
2022-03-25 22:13:34.143 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2022-03-25 22:13:34.144 [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  2022-03-25 22:12:59  2022-03-25 22:13:34  GROUP BY

Success!

Job Stats (time in seconds):
JobId  Maps  Reduces  MaxMapTime  MinMapTime  AvgMapTime  MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReduceTime  Alias  Feature Outputs
job_1644548343526_0030  1  1  5  5  5  5  6  6  6  6  Input1,groupedWords,wordsInEachLine  GROUP BY  hdfs://quickstart.cloudera:8020/tmp/temp-669075149/tmp325198341,
```

Name:

Roll No:

```
Success!
Job Stats (time in seconds):
JobId  Maps  Reduces  MaxMapTime  MinMapTime  AvgMapTime  MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReduceTime  Alias  Feature Outputs
job_1644548343526_0030  1  1  5  5  5  5  6  6  6  6  Input1.groupedWords,wordsInEachLine  GROUP_BY  hdf://quickstart.cloudera:8020/tmp/temp-669075149/tmp325198341.

Input(s):
Successfully read 3 records (497 bytes) from: "/user/cloudera/Training/pig/input.txt"

Output(s):
Successfully stored 17 records (407 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-669075149/tmp325198341"

Counters:
Total records written : 17
Total bytes written : 407
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
Job 1644548343526_0030

2022-03-25 22:13:34,226 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2022-03-25 22:13:34,226 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-25 22:13:34,226 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-25 22:13:34,227 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set. - Will not generate code.
2022-03-25 22:13:34,235 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2022-03-25 22:13:34,235 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(we, (we), (we))
(we, (we))
(BIG, (BIG))
(WMP, (WMP))
(are, (are))
(DATA, (DATA))
(DONE, (DONE))
(HAVE, (HAVE), (HAVE))
(WORD, (WORD))
(COUNT, (COUNT))
(USING, (USING))
(HADOOP, (HADOOP))
(REDUCE, (REDUCE))
(STUDED, (STUDED))
(studying, (studying))
(ECOSYSTEM, (ECOSYSTEM))
(TECHNOLOGY, (TECHNOLOGY))
grunt>
```

15) Now we count those words. For each group we count words in each line.

countedWords = foreach groupedWords generate group, COUNT(wordsInEachLine);

```
grunt> countedWords = foreach groupedWords generate group, COUNT(wordsInEachLine);
grunt>
```

16) After every counting of words commands, we are dumping the data dump countedWords;

Now the Final Output we are getting as word count for every word.

```
grunt> countedWords = foreach groupedWords generate group, COUNT(wordsInEachLine);
grunt> dump countedWords;
2022-03-25 22:17:07,065 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP BY
2022-03-25 22:17:07,066 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2022-03-25 22:17:07,070 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCompiler - File concatenation threshold: 100 optimistic? false
2022-03-25 22:17:07,075 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.CombinerOptimizer - Choosing to move algebraic foreach to combiner
2022-03-25 22:17:07,080 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2022-03-25 22:17:07,110 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - BytesPerReducer=1000000000 maxReducers=999 totalInputFileSize=110
2022-03-25 22:17:07,110 [main] INFO org.apache.hadoop.yarn.client.NMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2022-03-25 22:17:07,113 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting Parallelism to 1
2022-03-25 22:17:07,132 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Using reducer estimator: org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputSizeReducerEstimator
2022-03-25 22:17:07,133 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reduce phase detected, estimating # of required reducers.
2022-03-25 22:17:07,136 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting Parallelism to 1
2022-03-25 22:17:07,136 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Creating jar file Job7925997203331980059.jar
2022-03-25 22:17:11,021 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Jar file Job7925997203331980059.jar created
2022-03-25 22:17:11,036 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
2022-03-25 22:17:11,037 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2022-03-25 22:17:11,037 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserialize []
2022-03-25 22:17:11,047 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
2022-03-25 22:17:11,077 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-25 22:17:11,078 [JobControl] INFO org.apache.hadoop.yarn.client.NMProxy - Connecting to ResourceManager at /0.0.0.0:8032
2022-03-25 22:17:11,087 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-25 22:17:11,215 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2022-03-25 22:17:11,215 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
2022-03-25 22:17:11,222 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths (combined) to process : 1
2022-03-25 22:17:11,990 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - number of splits:1
2022-03-25 22:17:12,055 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job_1644548343526_0031
2022-03-25 22:17:12,059 [JobControl] INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl - Submitted application application_1644548343526_0031
2022-03-25 22:17:12,104 [JobControl] INFO org.apache.hadoop.mapreduce.Job - The url to track the job: http://quickstart.cloudera:8088/proxy/application_1644548343526_0031/
2022-03-25 22:17:12,106 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - HadoopJobId: job_1644548343526_0031
2022-03-25 22:17:12,106 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases Input1.countedWords,groupedWords,wordsInEachLine
2022-03-25 22:17:12,106 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - detailed locations: M: Input1[3.9],wordsInEachLine[1-1,1],countedWords[6.15],groupedWords[6.15]; R: countedWords[6.15]; groupedWords[6.15] R: countedWords[6.15]
2022-03-25 22:17:12,141 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - More information at: http://localhost:50830/jobdetails.jsp?jobid=job_1644548343526_0031
2022-03-25 22:17:12,141 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 0% complete
2022-03-25 22:17:20,042 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 50% complete
```

Name:

Roll No:

```
2022-03-25 22:17:42,771 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2022-03-25 22:17:42,771 [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 2022-03-25 22:17:07 2022-03-25 22:17:42 GROUP_BY
Success!
Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Feature Outputs
job_1644548343526_0031 1 1 5 5 5 5 6 6 6 6 Input1,CountedWords,groupedWords,wordsInEachLine GROUP_BY,COMBINER hdfs://quickstart.cloudera:8020/tmp/temp-669075149/tmp1191620340
1191620340
Input(s):
Successfully read 3 records (497 bytes) from: "/user/cloudera/Training/pig/input.txt"
Output(s):
Successfully stored 17 records (224 bytes) in: "hdfs://quickstart.cloudera:8020/tmp/temp-669075149/tmp1191620340"
Counters:
Total records written : 17
Total bytes written : 224
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_1644548343526_0031
2022-03-25 22:17:42,858 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2022-03-25 22:17:42,858 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2022-03-25 22:17:42,858 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2022-03-25 22:17:42,858 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... Will not generate code.
2022-03-25 22:17:42,866 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2022-03-25 22:17:42,866 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer - Total input paths to process : 1
```

```
(WE,2)
(we,1)
(BIG,1)
(MAP,1)
(are,1)
(DATA,1)
(DONE,1)
(HAVE,2)
(WORD,1)
(COUNT,1)
(USING,1)
(HADOOP,1)
(REDUCE,1)
(STUDED,1)
(studying,1)
(ECOSYSTEM,1)
(TWECHNOLOGY!,1)
grunt> █
```

As we can see from above image the Word “a” occurred twice, word “for, data” start with small w occurred twice, word “I” occurred once, and so on.

17) Now Exit from the grunt shell using quit command.

```
grunt> quit
[cloudera@quickstart ~]$ █
```