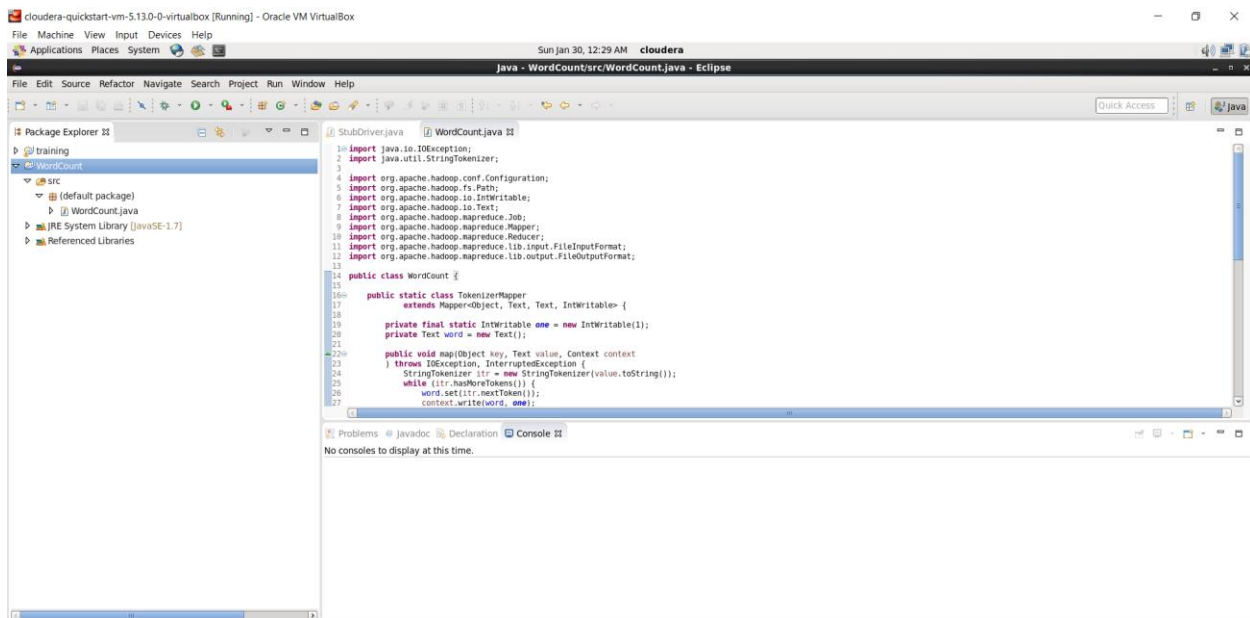# CSCE 5300 INTRODUCTION TO BIG DATA
# AND DATA SCIENCE

**NAME :** Eswara Reddy Thimmapuram

## Use case Mapper Reduce word count process

- Firstly, the input is given in a specified size and then the input is split into different parts. Then the next job is to convert input block into string. So the mapper takes starting address and ending address as input.
- Then the mapper converts input into string and tokenizes into words. The mapper will then append 1 to each word for first mapper output and so on. The output of mapper would be key-value pairs. After mapper process the next jobs are sorting and shuffling.
- In next phase, all words come to their respective places in sorting phase and in shuffling phase all similar words go to similar reducer. Then the reducer will sum up values by using keys. The output of reducer will be the final result which is the sum of words that are present in the input file.

## WordCount Program



- Firstly we create new java project WordCount and import all header files and will add all required external jar files. Then we create new class named WordCount .
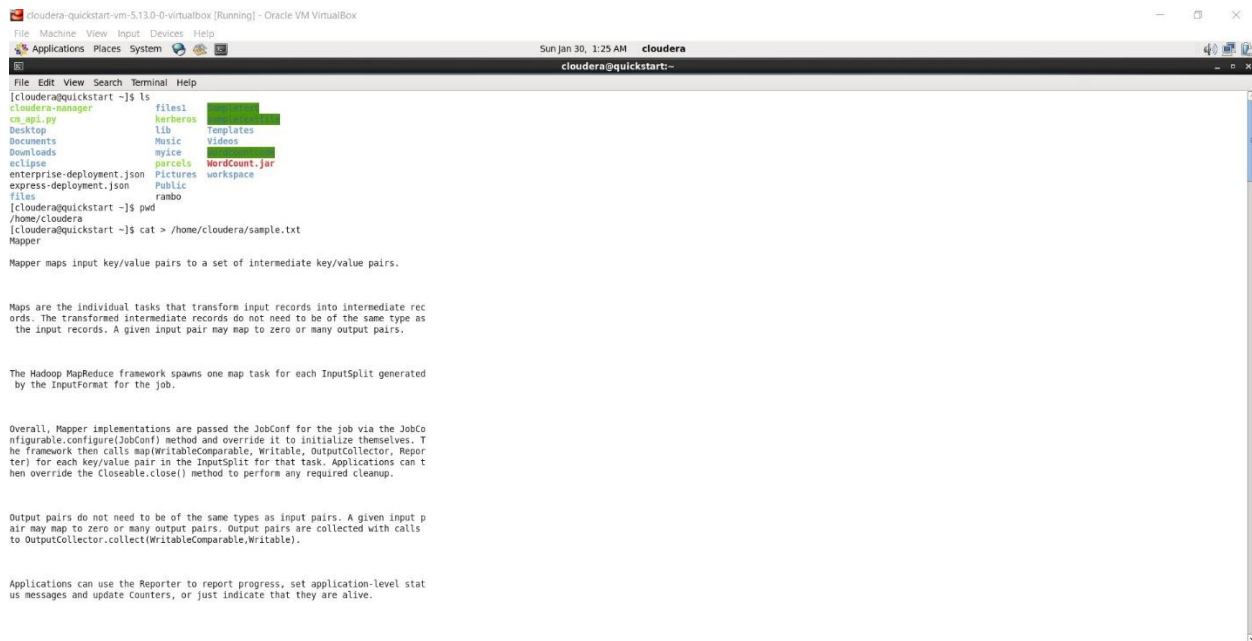
- Then we extend the mapper class by creating TokenizerMapper and it takes arguments like key-in,value-in,key-out and value-out. Then we create word variable to store word and one variable to store value.
- Then we create map function which takes arguments like key-in as key, textual content as value variable and context to emit result to next stage.
- Then we StringTokenizer which divides string into tokens and store each of them in itr variable.
- Then we define reducer class as IntSumReducer which takes four inputs. Then we create result variable to store how many number of times the key is appearing.
- Then we write reduce function which takes word as key and value as value and if the word repeats then it will be added to the sum value.
- Finally the reducer gives the final output of the program which is the total sum of words that are present in input file.



- Then we create Sampletext file to store the input we have to give to the program.

- Then we use cat command to store the input into sampletext file.



- Then we create input1 file and use put command to copy file from local system to Hadoop system.

Applications can use the Reporter to report progress, set application-level status messages and update Counters, or just indicate that they are alive.

The output of the Reducer is not sorted.
[cloudera@quickstart ~]$ hadoop jar /home/cloudera/WordCount.jar WordCount input1/sample.txt output1
22/01/30 01:33:21 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
22/01/30 01:33:24 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
22/01/30 01:33:25 INFO input.FileInputFormat: Total input paths to process : 1
22/01/30 01:33:25 INFO mapreduce.JobSubmitter: number of splits:1
22/01/30 01:33:27 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1643520845524_0001
22/01/30 01:33:31 INFO impl.YarnClientImpl: Submitted application application_1643520845524_0001
22/01/30 01:33:32 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8088/proxy/application_1643520845524_0001/
22/01/30 01:33:32 INFO mapreduce.Job: Running job: job_1643520845524_0001
22/01/30 01:34:11 INFO mapreduce.Job: Job job_1643520845524_0001 running in uber mode : false
22/01/30 01:34:11 INFO mapreduce.Job:  map 0% reduce 0%
22/01/30 01:34:59 INFO mapreduce.Job:  map 100% reduce 0%
22/01/30 01:39:55 INFO mapreduce.Job:  map 100% reduce 100%
22/01/30 01:39:55 INFO mapreduce.Job: Job job_1643520845524_0001 completed successfully
22/01/30 01:39:55 INFO mapreduce.Job: Counters: 49
        File System Counters
                FILE: Number of bytes read=4204
                FILE: Number of bytes written=295149
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=4924
                HDFS: Number of bytes written=3035
                HDFS: Number of read operations=6
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
        Job Counters
                Launched map tasks=1
                Launched reduce tasks=1
                Data-local map tasks=1
                Total time spent by all maps in occupied slots (ms)=43587
                Total time spent by all reduces in occupied slots (ms)=107215
                Total time spent by all map tasks (ms)=43587
                Total time spent by all reduce tasks (ms)=107215
                Total vcore-milliseconds taken by all map tasks=43587
                Total vcore-milliseconds taken by all reduce tasks=107215
                Total megabyte-milliseconds taken by all map tasks=44633088
                Total megabyte-milliseconds taken by all reduce tasks=109788160
        Map-Reduce Framework
                Map input records=111
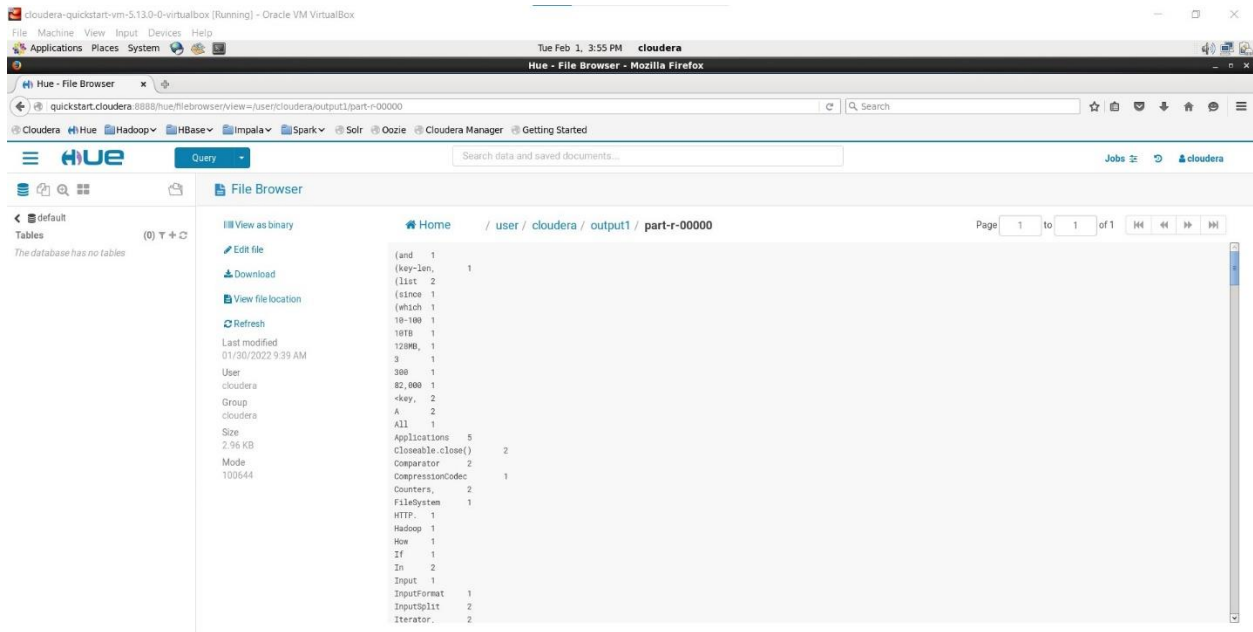                Map output records=719

- Then we run the mapreduce job command and then the execution continues.



                WRONG_REDUCE=0
        File Input Format Counters
                Bytes Read=4796
        File Output Format Counters
                Bytes Written=3035
[cloudera@quickstart ~]$ hadoop fs -ls output1
Found 2 items
-rw-r--r--   1 cloudera cloudera          0 2022-01-30 01:39 output1/_SUCCESS
-rw-r--r--   1 cloudera cloudera       3035 2022-01-30 01:39 output1/part-r-00000
[cloudera@quickstart ~]$ hadoop fs -cat output1/part-r-00000
(and     1
(key-len,     1
(list    2
(since   1
(which   1
10-100   1
10TB     1
128MB,   1
3        1
300      1
82,000   1
<key,    2
A        2
All      1
Applications     5
Closeable.close()       2
Comparator       2
CompressionCodec        1
Counters,        2
FileSystem       1
HTTP.    1
Hadoop   1
How      1
If       1
In       2
Input    1
InputFormat      1
InputSplit       2
Iterator,        2
JobConf 2
JobConf.         1
JobConf.setCombinerClass(Class),        1
JobConf.setNumReduceTasks(int). 1
JobConf.setOutputKeyComparatorClass(Class)      1
JobConf.setOutputKeyComparatorClass(Class).     1
JobConf.setOutputValueGroupingComparator(Class).     1
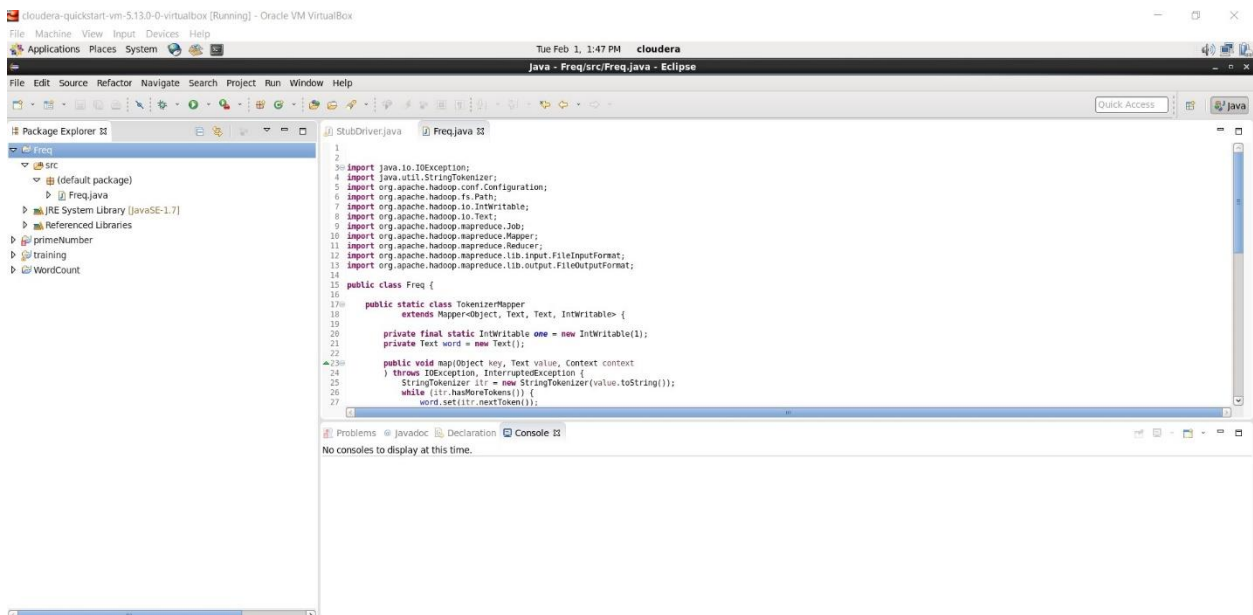JobConfigurable.configure(JobConf)      2
Many     1
MapReduce        1

- Finally, we list the files in the output and use cat command to display the output which is the number of times the words are repeated in given file.

- Here is the visualized version of wordcount program

## Frequency of words with letter 'b'



- Firstly we create a java project named Freq and then import all the required header files and add all external jar libs to the program.

- Then we extend the mapper class by creating TokenizerMapper and it takes arguments like key-in,value-in,key-out and value-out. Then we create word variable to store word and one variable to store value.
- Then we create map function which takes arguments like key-in as key, textual content as value variable and context to emit result to next stage.
- Then we StringTokenizer which divides string into tokens and store each of them in itr variable.
- Then we define reducer class as IntSumReducer which takes four inputs. Then we create result variable to store how many number of times the key is appearing.
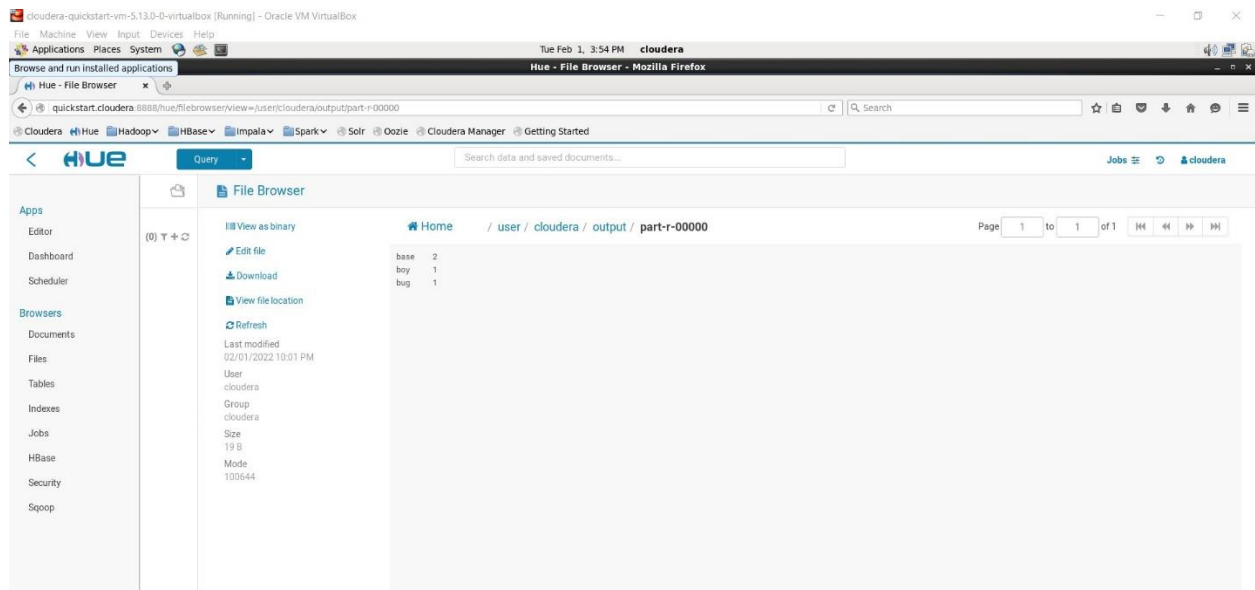- Then we write logic which adds words which starts with letter b and gives output to the reducer.



- Then we create a file and add required input to it using cat command. Then after we import files from local to Hadoop using put command.
- Then we display input using cat command. Then we execute map reduce job and execution continues.

- Finally, we list the files in the output and use cat command to display the output which is the number of times the words with letter 'b' are repeated in given file.



- Here is visualized version of frequency of b words file

# Prime Number Program



- Firstly we create a java project named prime and then import all the required header files and add all external jar libs to the program.
- Then we extend the mapper class by creating TokenizerMapper and it takes arguments like key-in,value-in,key-out and value-out. Then we create word variable to store word and one variable to store value.
- Then we create map function which takes arguments like key-in as key, textual content as value variable and context to emit result to next stage.
- Then we StringTokenizer which divides string into tokens and store each of them in itr variable.
- Then we define reducer class as IntSumReducer which takes four inputs. Then we create result variable to store how many number of times the key is appearing.
- Then we write logic which adds words which starts with letter b and gives output to the reducer.

Java - Prime/src/Prime.java - Eclipse          _  □  ✕

File  Edit  Source  Refactor  Navigate  Search  Project  Run  Window  Help

Quick Access          🔲  | 🟦 Java  🔲 CVS Repository Exploring

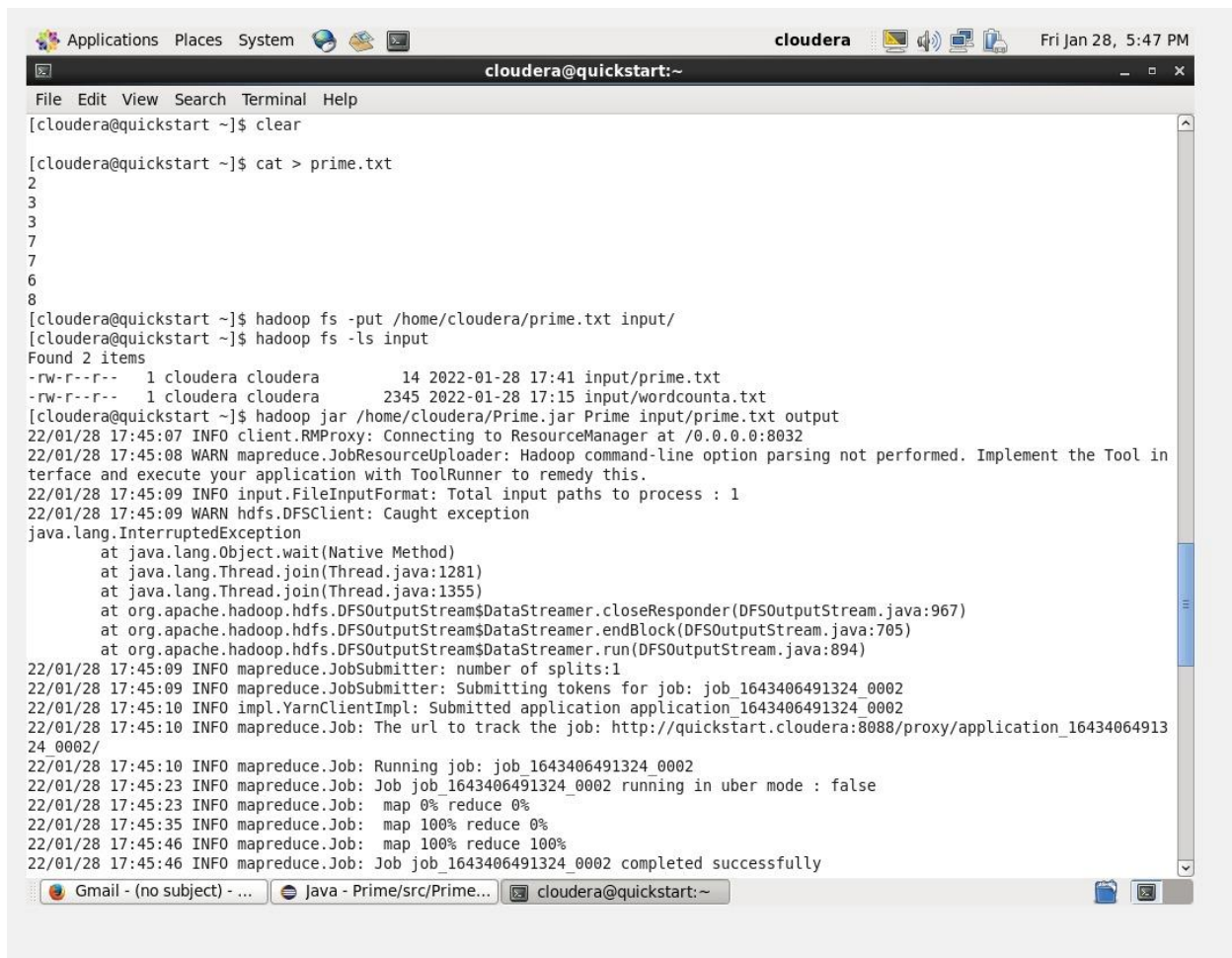| 📄 WordCounta.java | 📄 WordCo.java | 📄 Prime.java ⊠ |

```
32        extends Reducer<Text, IntWritable, Text, IntWritable> {
33            private IntWritable result = new IntWritable();
34            public void reduce(Text key, Iterable<IntWritable> values,
35                             Context context
36                             ) throws IOException, InterruptedException {
37                int output = 0;
38                int num = Integer.parseInt(key.toString());
39
40                for (int i = 2; i < num; i++) {
41                        if(num % i == 0) {
42                        output = 1;
43                        }
44                }
45                result.set(output);
46                context.write(key, result);
47            }
48        }
49
50        public static void main(String[] args) throws Exception {
51            Configuration conf = new Configuration();
52            Job job = Job.getInstance(conf, "word count");
53            job.setJarByClass(Prime.class);
54            job.setMapperClass(TokenizerMapper.class);
55            job.setCombinerClass(IntSumReducer.class);
56            job.setReducerClass(IntSumReducer.class);
57            job.setOutputKeyClass(Text.class);
58            job.setOutputValueClass(IntWritable.class);
59            FileInputFormat.addInputPath(job, new Path(args[0]));
60            FileOutputFormat.setOutputPath(job, new Path(args[1]));
61            System.exit(job.waitForCompletion(true) ? 0 : 1);
62        }
63
64
65 }
```

|  | Writable | Smart Insert | 1 : 1 |

🔴 [Cloudera Live : Welco...    ⬤ Java - Prime/src/Prime...          📋  ⬤

- Then I have added input to file using cat command and listed the files using ls command.

- Then I have run mapreduce job and exceuted the code.

cloudera@quickstart:~

File  Edit  View  Search  Terminal  Help

```
                Total megabyte-milliseconds taken by all map tasks=10501120
                Total megabyte-milliseconds taken by all reduce tasks=7277568
        Map-Reduce Framework
                Map input records=7
                Map output records=7
                Map output bytes=42
                Map output materialized bytes=46
                Input split bytes=126
                Combine input records=7
                Combine output records=5
                Reduce input groups=5
                Reduce shuffle bytes=46
                Reduce input records=5
                Reduce output records=5
                Spilled Records=10
                Shuffled Maps =1
                Failed Shuffles=0
                Merged Map outputs=1
                GC time elapsed (ms)=254
                CPU time spent (ms)=1670
                Physical memory (bytes) snapshot=335568896
                Virtual memory (bytes) snapshot=3015303168
                Total committed heap usage (bytes)=226365440
        Shuffle Errors
                BAD_ID=0
                CONNECTION=0
                IO_ERROR=0
                WRONG_LENGTH=0
                WRONG_MAP=0
                WRONG_REDUCE=0
        File Input Format Counters
                Bytes Read=14
        File Output Format Counters
                Bytes Written=20
[cloudera@quickstart ~]$ hadoop fs -ls output
Found 2 items
-rw-r--r--   1 cloudera cloudera          0 2022-01-28 17:45 output/_SUCCESS
-rw-r--r--   1 cloudera cloudera         20 2022-01-28 17:45 output/part-r-00000
[cloudera@quickstart ~]$
```

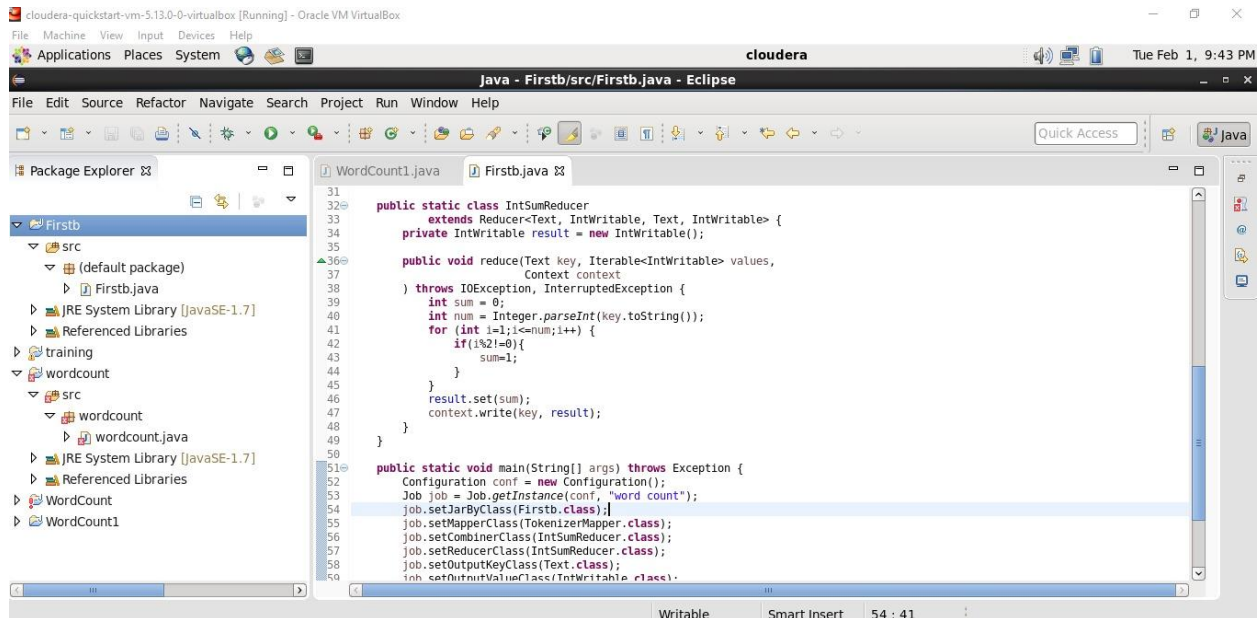Gmail - (no subject) - ...    Java - Prime/src/Prime...    cloudera@quickstart:~

Then I displayed the output using ls command and it gives prime numbers as 1 and others as 0.

Here is visualized version of prime number program

# Odd number Program



```java
31
32    public static class IntSumReducer
33            extends Reducer<Text, IntWritable, Text, IntWritable> {
34        private IntWritable result = new IntWritable();
35
36        public void reduce(Text key, Iterable<IntWritable> values,
37                            Context context
38        ) throws IOException, InterruptedException {
39            int sum = 0;
40            int num = Integer.parseInt(key.toString());
41            for (int i=1;i<=num;i++) {
42                if(i%2!=0){
43                    sum=1;
44                }
45            }
46            result.set(sum);
47            context.write(key, result);
48        }
49    }
50
51    public static void main(String[] args) throws Exception {
52        Configuration conf = new Configuration();
53        Job job = Job.getInstance(conf, "word count");
54        job.setJarByClass(Firstb.class);
55        job.setMapperClass(TokenizerMapper.class);
56        job.setCombinerClass(IntSumReducer.class);
57        job.setReducerClass(IntSumReducer.class);
58        job.setOutputKeyClass(Text.class);
59        job.setOutputValueClass(IntWritable.class);
```



```
            Bytes Written=12
[cloudera@quickstart ~]$ hadoop fs -ls out
Found 2 items
-rw-r--r--   1 cloudera cloudera          0 2022-02-01 21:27 out/_SUCCESS
-rw-r--r--   1 cloudera cloudera         12 2022-02-01 21:27 out/part-r-00000
[cloudera@quickstart ~]$ hadoop fs -cat out/part-r-00000
1       1
3       1
5       1
[cloudera@quickstart ~]$
```