

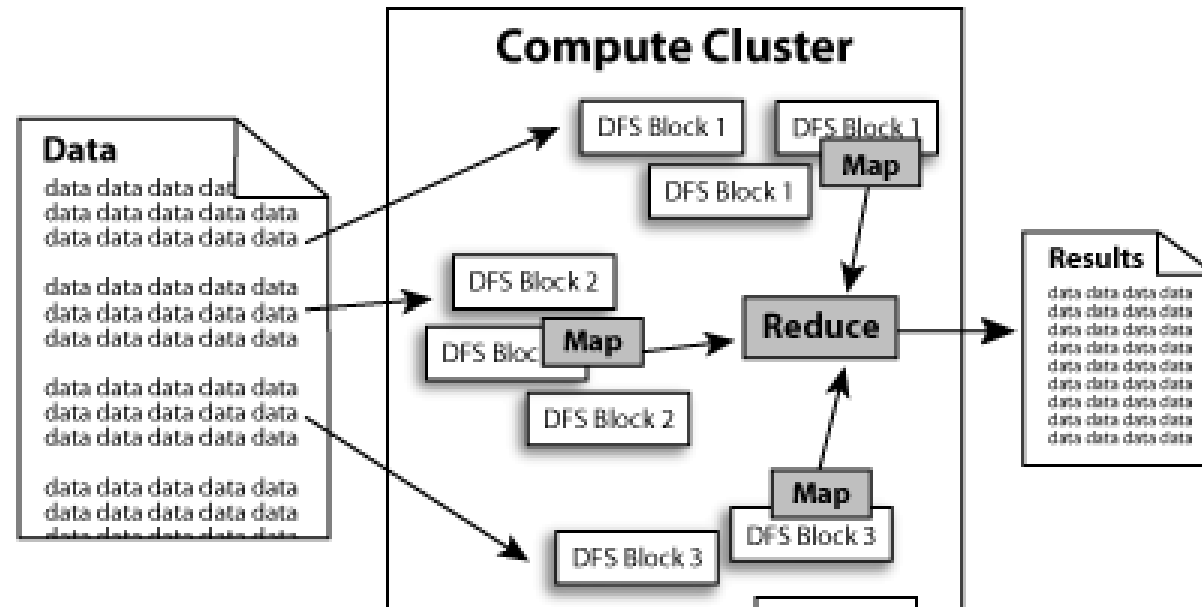
Word Counting Using MapReduce

Sorayut Glomglome

Content

1. Explain MapReduce for word count
2. Hands-on practice for MapReduce

MapReduce on Hadoop



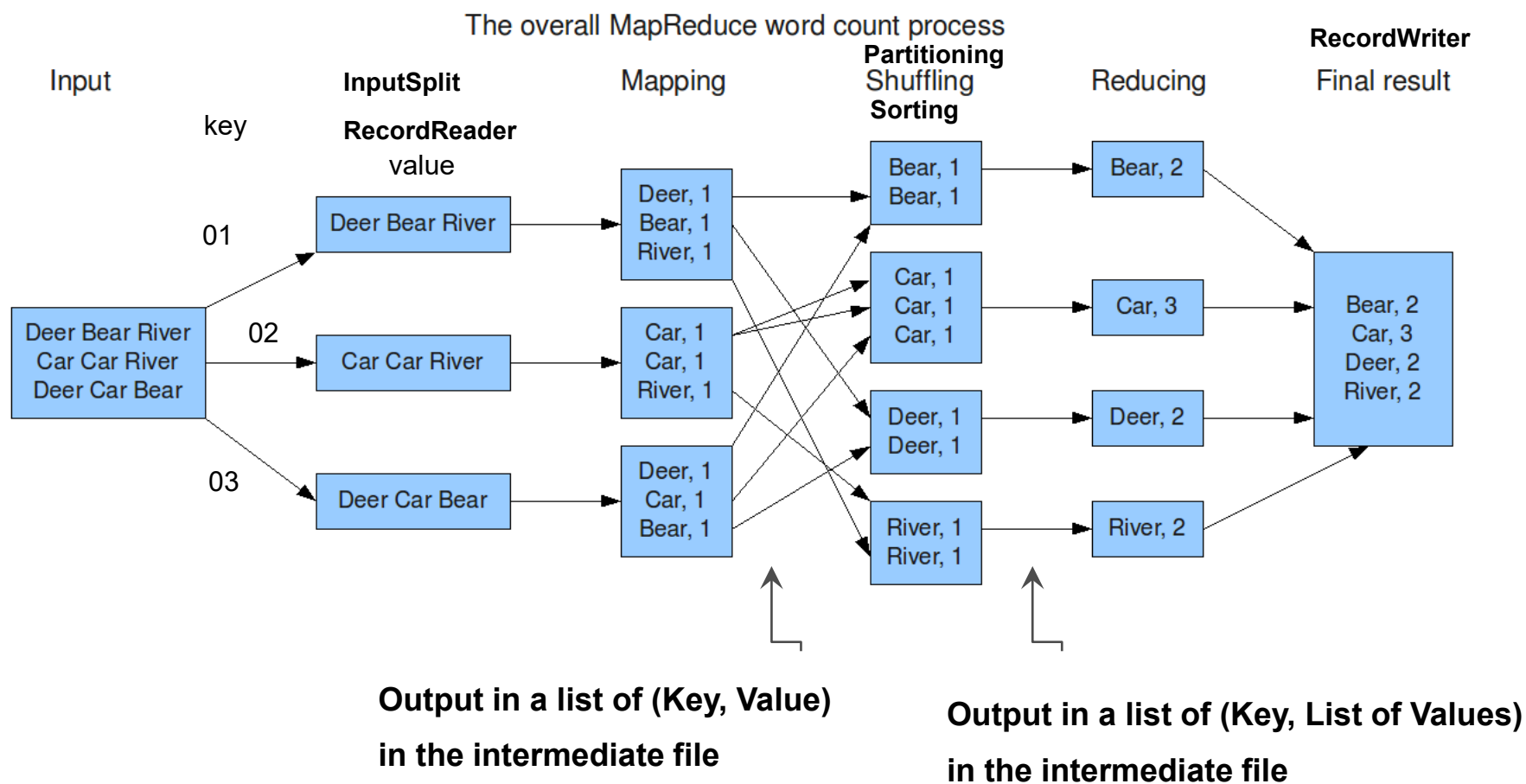
map: (K1, V1) -> list(K2, V2)

reduce: (K2, list(V2)) -> list(K3, V3)

MapReduce - Word Count

map: (K1, V1) -> list(K2, V2))

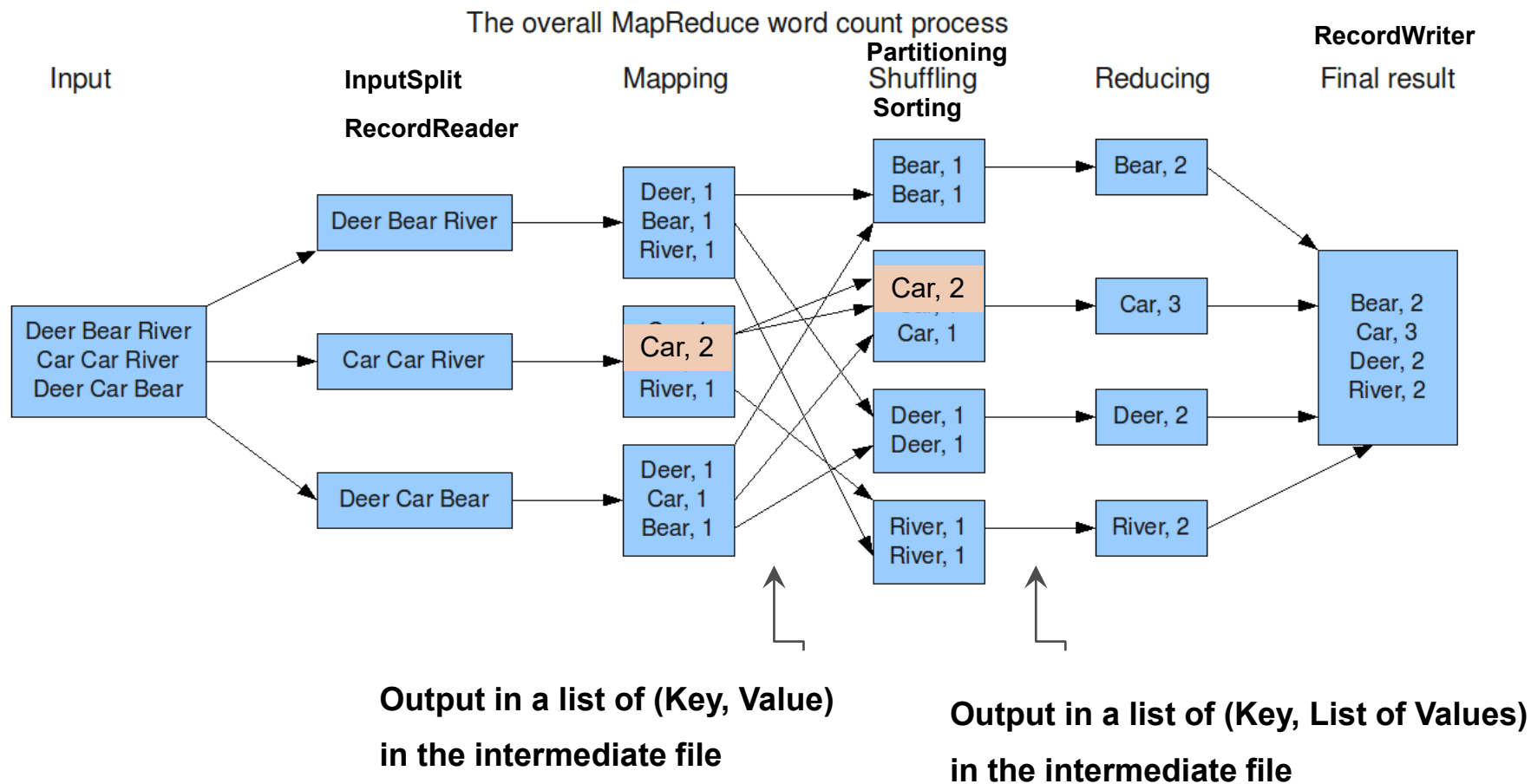
reduce: (K2, list(V2)) -> list(K3, V3)



MapReduce - Word Count (accelerating)

map: (K1, V1) -> list(K2, V2))

reduce: (K2, list(V2)) -> list(K3, V3)



MapReduce Processing – The Data flow

1. InputFormat, InputSplits, RecordReader
2. Mapper - *your focus is here*
3. Partition, Shuffle & Sort
4. Reducer - *your focus is here*
5. OutputFormat, RecordWriter

InputFormat

InputFormat:	Description:	Key:	Value:
TextInputFormat	Default format; reads lines of text files	The byte offset of the line	The line contents
KeyValueInputFormat	Parses lines into key, val pairs	Everything up to the first tab character	The remainder of the line
SequenceFileInputFormat	A Hadoop-specific high-performance binary format	user-defined	user-defined

InputSplit

An InputSplit describes a unit of work that comprises a single *map task*.

InputSplit presents a byte-oriented view of the input.

You can control this value by setting the **mapred.min.split.size** parameter in **core-site.xml**, or by overriding the parameter in the JobConf object used to submit a particular MapReduce job.

RecordReader

RecordReader reads <key, value> pairs from an InputSplit.

Typically the RecordReader converts the byte-oriented view of the input, provided by the InputSplit, and presents a record-oriented to the Mapper

Mapper

```
map: (K1, V1) -> list(K2, V2)
reduce: (K2, list(V2)) -> list(K3, V3)
```

The **Mapper** performs the **user-defined** logic to the input a key, value and emits **(key, value) pair(s)** which are forwarded to the **Reducers**.

Partition, Shuffle & Sort

```
map: (K1, V1) -> list(K2, V2)  
reduce: (K2, list(V2)) -> list(K3, V3)
```

After the first map tasks have completed, the nodes may still be performing several more map tasks each. But they also begin exchanging the intermediate outputs from the map tasks to where they are required by the reducers.

Partitioner controls the partitioning of map-outputs to assign to reduce task . The total number of partitions is the same as the number of reduce tasks for the job

The set of intermediate keys on a single node is **automatically sorted** by internal Hadoop before they are presented to the Reducer

This process of moving map outputs to the reducers is known as **shuffling**.

Reducer

```
map: (K1, V1) -> list(K2, V2)
reduce: (K2, list(V2)) -> list(K3, V3)
```

This is an instance of **user-provided code** that performs read each key, iterator of values in the partition assigned. The **OutputCollector** object in **Reducer** phase has a method named **collect()** which will collect a (key, value) output.

OutputFormat, Record Writer

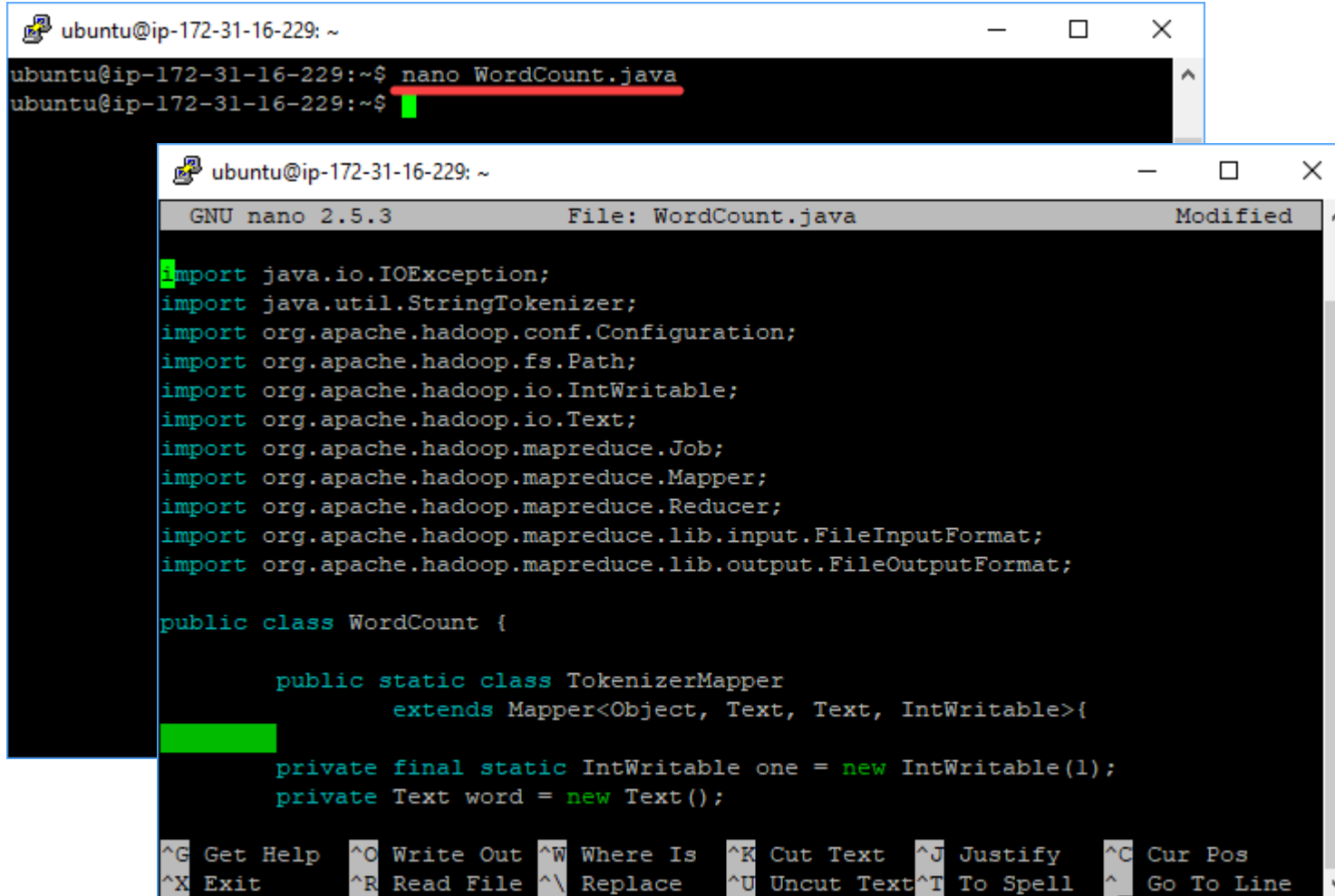
OutputFormat governs the writing format in **OutputCollector** and **RecordWriter** writes output into HDFS.

TextOutputFormat	Default; writes lines in "key \t value" form
SequenceFileOutputFormat	Writes binary files suitable for reading into subsequent MapReduce jobs
NullOutputFormat	generates no output files

Hands-on practice for MapReduce (Word Count)

Create a Java file

\$nano WordCount.java

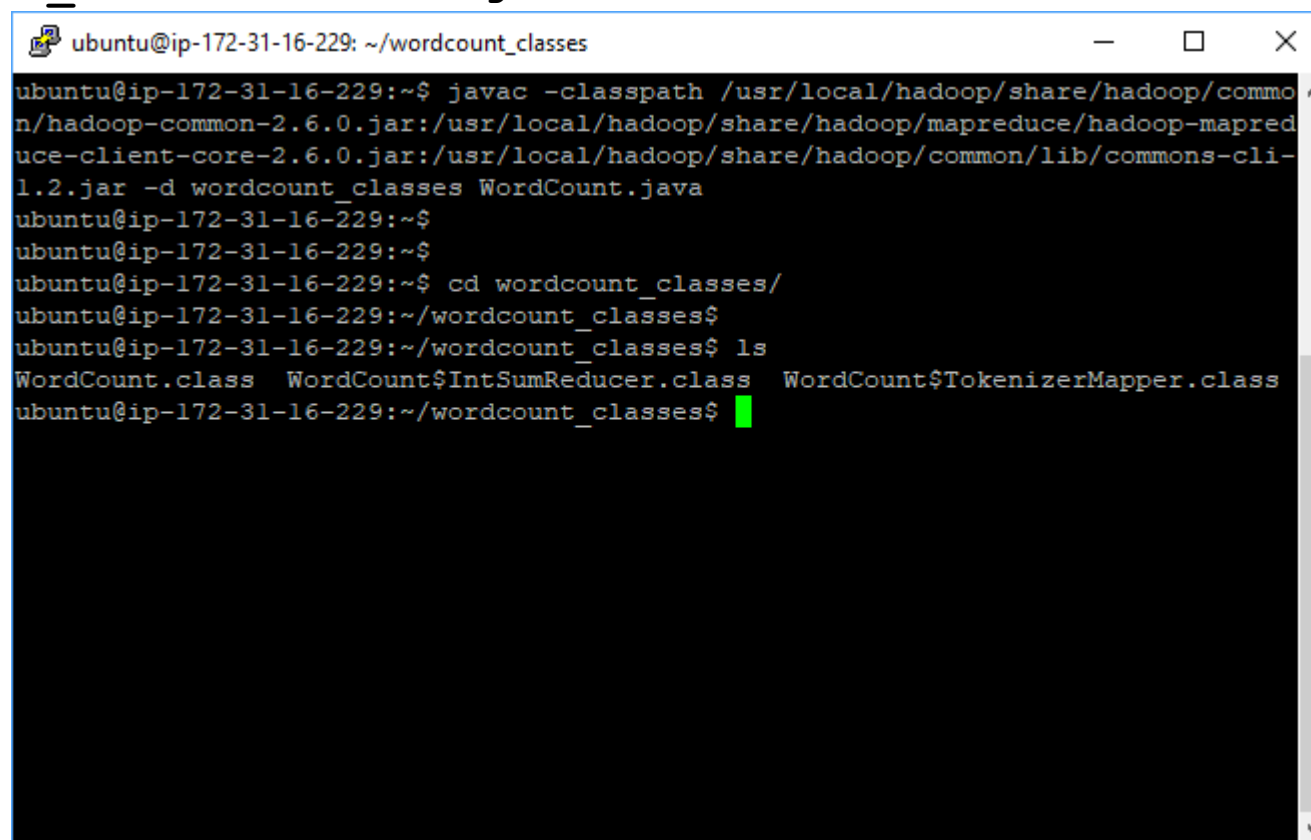


```
ubuntu@ip-172-31-16-229: ~  
ubuntu@ip-172-31-16-229:~$ nano WordCount.java  
ubuntu@ip-172-31-16-229:~$  
GNU nano 2.5.3 File: WordCount.java Modified  
import java.io.IOException;  
import java.util.StringTokenizer;  
import org.apache.hadoop.conf.Configuration;  
import org.apache.hadoop.fs.Path;  
import org.apache.hadoop.io.IntWritable;  
import org.apache.hadoop.io.Text;  
import org.apache.hadoop.mapreduce.Job;  
import org.apache.hadoop.mapreduce.Mapper;  
import org.apache.hadoop.mapreduce.Reducer;  
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;  
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;  
  
public class WordCount {  
  
    public static class TokenizerMapper  
        extends Mapper<Object, Text, Text, IntWritable>{  
  
        private final static IntWritable one = new IntWritable(1);  
        private Text word = new Text();  
  
^G Get Help ^O Write Out ^W Where Is ^K Cut Text ^J Justify ^C Cur Pos  
^X Exit ^R Read File ^\ Replace ^U Uncut Text ^T To Spell ^_ Go To Line
```

Create a directory for compiled code and compile

```
$mkdir wordcount_classes
```

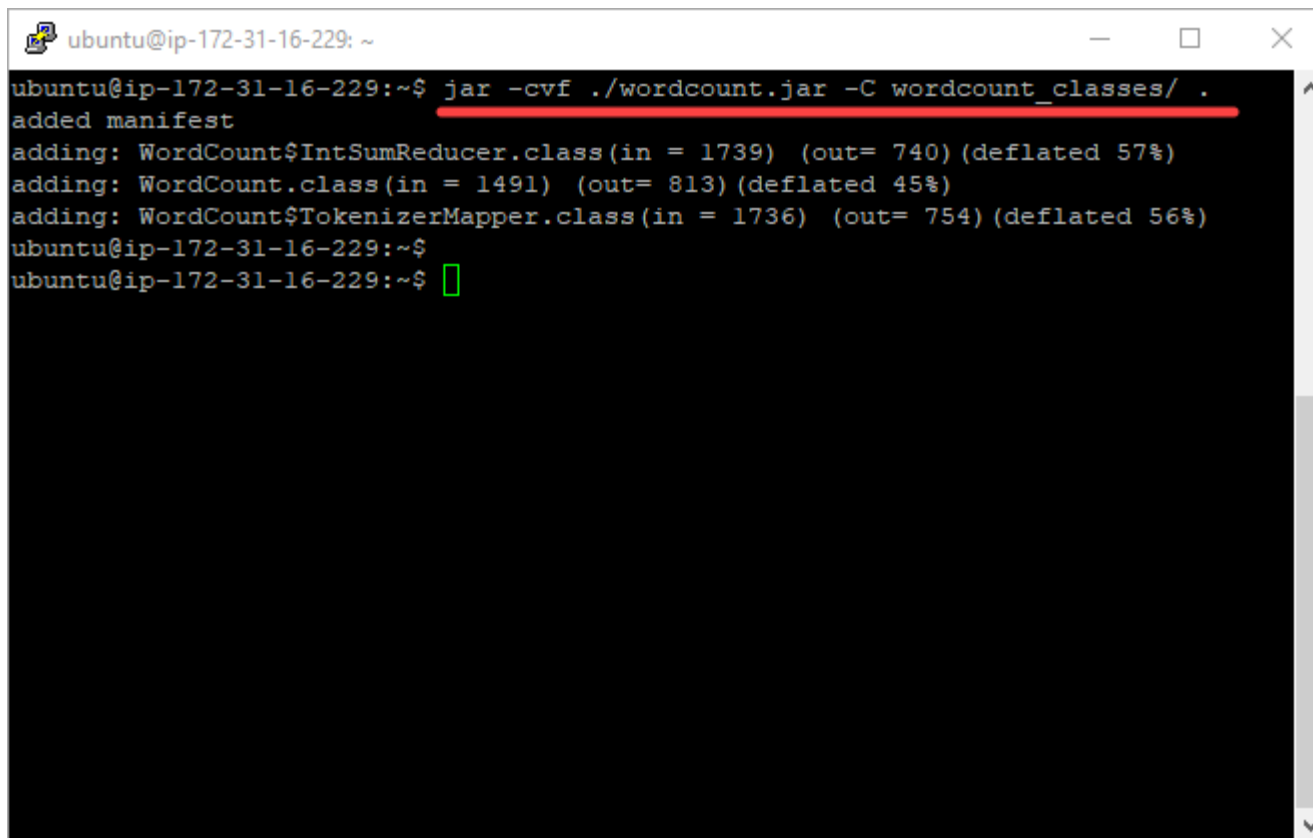
```
$javac -classpath /usr/local/hadoop/share/hadoop/common/hadoop-common-  
2.6.0.jar:/usr/local/hadoop/share/hadoop/mapreduce/hadoop-mapreduce-client-core-  
2.6.0.jar:/usr/local/hadoop/share/hadoop/common/lib/commons-cli-1.2.jar -d  
wordcount_classes WordCount.java
```



```
ubuntu@ip-172-31-16-229: ~/wordcount_classes  
ubuntu@ip-172-31-16-229:~$ javac -classpath /usr/local/hadoop/share/hadoop/commo  
n/hadoop-common-2.6.0.jar:/usr/local/hadoop/share/hadoop/mapreduce/hadoop-mapred  
uce-client-core-2.6.0.jar:/usr/local/hadoop/share/hadoop/common/lib/commons-cli-  
1.2.jar -d wordcount_classes WordCount.java  
ubuntu@ip-172-31-16-229:~$  
ubuntu@ip-172-31-16-229:~$  
ubuntu@ip-172-31-16-229:~$ cd wordcount_classes/  
ubuntu@ip-172-31-16-229:~/wordcount_classes$  
ubuntu@ip-172-31-16-229:~/wordcount_classes$ ls  
WordCount.class  WordCount$IntSumReducer.class  WordCount$TokenizerMapper.class  
ubuntu@ip-172-31-16-229:~/wordcount_classes$
```


Create a Jar file

```
$jar -cvf ./wordcount.jar -C wordcount_classes/ .
```

A terminal window titled 'ubuntu@ip-172-31-16-229: ~' with standard window controls. The command 'jar -cvf ./wordcount.jar -C wordcount_classes/ .' is entered and executed. The output shows 'added manifest' followed by three lines of class addition statistics: 'adding: WordCount\$IntSumReducer.class(in = 1739) (out= 740) (deflated 57%)', 'adding: WordCount.class(in = 1491) (out= 813) (deflated 45%)', and 'adding: WordCount\$TokenizerMapper.class(in = 1736) (out= 754) (deflated 56%)'. The prompt returns to the shell.

```
ubuntu@ip-172-31-16-229:~$ jar -cvf ./wordcount.jar -C wordcount_classes/ .
added manifest
adding: WordCount$IntSumReducer.class(in = 1739) (out= 740) (deflated 57%)
adding: WordCount.class(in = 1491) (out= 813) (deflated 45%)
adding: WordCount$TokenizerMapper.class(in = 1736) (out= 754) (deflated 56%)
ubuntu@ip-172-31-16-229:~$
```

Don't forget "."!!!

Add a MapReduce Job to YARN

```
$yarn jar ./wordcount.jar WordCount /inputs/* /outputs/wordcount_output_dir01
```

Run all files in directory.

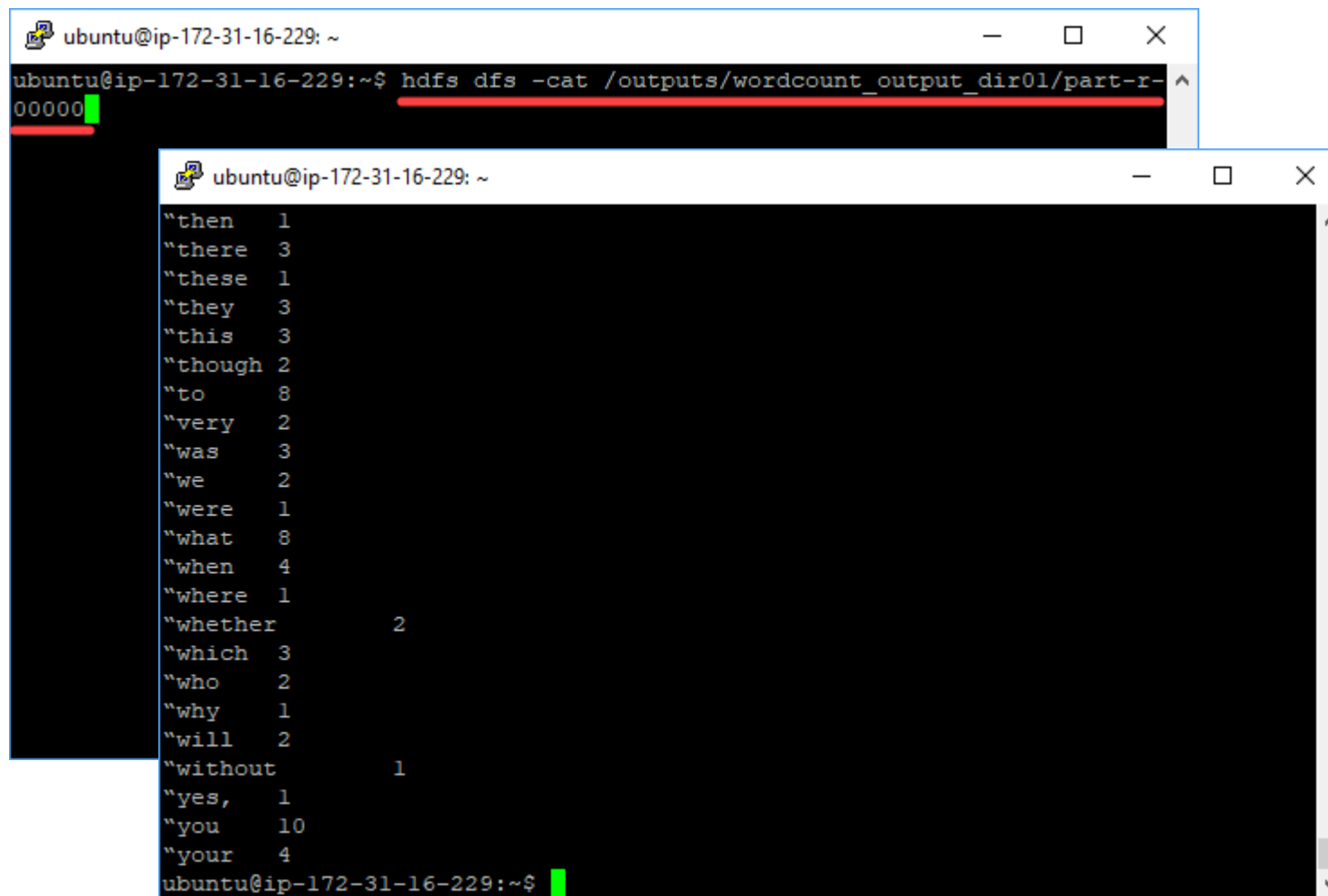
```
ubuntu@ip-172-31-16-229: ~  
ubuntu@ip-172-31-16-229:~$ jar -cvf ./wordcount.jar -C wordcount_classes/ .  
added manifest  
adding: WordCount$IntSumReducer.class(in = 1739) (out= 740) (deflated 57%)  
adding: WordCount.class(in = 1491) (out= 813) (deflated 45%)  
adding: WordCount$TokenizerMapper.class(in = 1736) (out= 754) (deflated 56%)  
ubuntu@ip-172-31-16-229:~$  
ubuntu@ip-172-31-16-229:~$ yarn jar ./wordcount.jar WordCount /inputs/* /outputs  
/wordcount output dir01  
19/02/20 13:21:53 INFO Configuration.deprecation: session.id is deprecated. Inst  
ead, use dfs.metrics.session-id  
19/02/20 13:21:53 INFO jvm.JvmMetrics: Initializing JVM Metrics with processName  
=JobTracker, sessionId=  
19/02/20 13:21:53 WARN mapreduce.JobSubmitter: Hadoop command-line option parsin  
g not performed. Implement the Tool interface and execute your application with  
ToolRunner to remedy this.  
19/02/20 13:21:53 INFO input.FileInputFormat: Total input paths to process : 1  
19/02/20 13:21:53 INFO mapreduce.JobSubmitter: number of splits:1  
19/02/20 13:21:53 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_lo  
cal581792515_0001  
19/02/20 13:21:53 INFO mapreduce.Job: The url to track the job: http://localhost  
:8080/  
19/02/20 13:21:53 INFO mapreduce.Job: Running job: job_local581792515_0001  
19/02/20 13:21:53 INFO mapred.LocalJobRunner: OutputCommitter set in config null  
19/02/20 13:21:53 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hado
```

Can't overwrite existing directory.

Use the new one each time.

View the result using command line

```
$hdfs dfs -cat /outputs/wordcount_output_dir01/part-r-00000
```



The image shows two terminal windows. The top window displays the command `hdfs dfs -cat /outputs/wordcount_output_dir01/part-r-00000` being executed. The bottom window shows the output of the command, which is a list of words and their counts, formatted as `"word" count`.

```
ubuntu@ip-172-31-16-229: ~  
ubuntu@ip-172-31-16-229:~$ hdfs dfs -cat /outputs/wordcount_output_dir01/part-r-00000  
"then" 1  
"there" 3  
"these" 1  
"they" 3  
"this" 3  
"though" 2  
"to" 8  
"very" 2  
"was" 3  
"we" 2  
"were" 1  
"what" 8  
"when" 4  
"where" 1  
"whether" 2  
"which" 3  
"who" 2  
"why" 1  
"will" 2  
"without" 1  
"yes," 1  
"you" 10  
"your" 4  
ubuntu@ip-172-31-16-229:~$
```

View the result using web console

<http://52.26.15.54:50070>

Hadoop Overview Datanodes Snapshot Startup Progress Utilities

Browse Directory

/outputs/wordcount_output_dir01 Go!

Permission	Owner	Group	Size	Replication	Block Size	Name
-rw-r--r--	ubuntu	supergroup	0 B	1	128 MB	_SUCCESS
-rw-r--r--	ubuntu	supergroup	147.67 KB	1	128 MB	part-r-00000

Hadoop, 2014.

1	#1342]	1	
2	\$5,000)	1	
3	'AS-IS'	1	
4	'Ah!	1	
5	'Bingley,	1	
6	'Having	1	
7	'I	1	
8	'Keep	1	
9	'Lady	1	
10	'Lately,	1	
11	'Lydia	1	
12	'Mr.	1	
13	'Oh!	1	
14	'This	1	
15	'Tis	1	
16	'Yes,'	1	
17	'_She	1	
18	'had	1	
19	'violently	1	
20	'you	1	

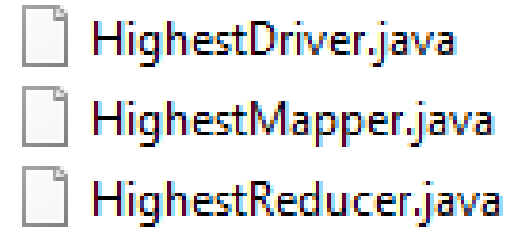
13619	"they	3	
13620	"this	3	
13621	"though	2	
13622	"to	8	
13623	"very	2	
13624	"was	3	
13625	"we	2	
13626	"were	1	
13627	"what	8	
13628	"when	4	
13629	"where	1	
13630	"whether	2	
13631	"which	3	
13632	"who	2	
13633	"why	1	
13634	"will	2	
13635	"without	1	
13636	"yes,	1	
13637	"you	10	
13638	"your	4	

Practice for MapReduce (Weather Data)

<ftp://ftp.ncdc.noaa.gov/pub/data/noaa/2018>

Instruction

1. Import all files to EC2
2. Compile three java files
3. Pack class files to a jar file
4. Input weather data files to hdfs directory
5. Run MapReduce



HighestDriver.java
HighestMapper.java
HighestReducer.java

Don't Forget to
STOP or TERMINATE
EC2 Instance !!!!!