

Lab 7 TopN

TopN Driver Code:

```
package samples.topn;
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.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.GenericOptionsParser;

public class TopN {
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        String[] otherArgs = (new GenericOptionsParser(conf, args)).getRemainingArgs();
        if (otherArgs.length != 2) {
            System.err.println("Usage: TopN <in> <out>");
            System.exit(2);
        }
        Job job = Job.getInstance(conf);
        job.setJobName("Top N");
        job.setJarByClass(TopN.class);
        job.setMapperClass(TopNMapper.class);
        job.setReducerClass(TopNReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        FileInputFormat.addInputPath(job, new Path(otherArgs[0]));
        FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
}
```

TopN Mapper Code:

```
package samples.topn;
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.io.IntWritable;
```

TopN Reducer Code:

```
package samples.topn;
import java.io.IOException;
import java.util.HashMap;
import java.util.Map;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;
import utils.MiscUtils;

public class TopNReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    private Map<Text, IntWritable> countMap = new HashMap<>();

    public void reduce(Text key, Iterable<IntWritable> values, Context context) throws
IOException, InterruptedException {
        int sum = 0;
        for (IntWritable val : values)
            sum += val.get();
        this.countMap.put(new Text(key), new IntWritable(sum));
    }

    protected void cleanup(Context context) throws IOException, InterruptedException {
        Map<Text, IntWritable> sortedMap = MiscUtils.sortByValues(this.countMap);
        int counter = 0;
    }
}
```

```

    for (Text key : sortedMap.keySet()) {
        if (counter++ == 20)
            break;
        context.write(key, sortedMap.get(key));
    }
}
}

```

Output screenshots:

```

hadoop@bmscece-HP-Elite-Tower-600-G9-Desktop-PC:~$ jps
12082 NameNode
19238 Jps
13014 NodeManager
15100 org.eclipse.equinox.launcher_1.6.1000.v20250227-1734.jar
12845 ResourceManager
12557 SecondaryNameNode
12255 DataNode
hadoop@bmscece-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -copyFromLocal /home/hadoop/Desktop/TopN.txt /bda_
hadoop/tinput.txt
copyFromLocal: '/bda_hadoop/tinput.txt': No such file or directory: 'hdfs://localhost:9000/bda_hadoop/tinput.txt'
hadoop@bmscece-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -copyFromLocal /home/hadoop/Desktop/TopN.txt /bda_
hadoop/tinput.txt
hadoop@bmscece-HP-Elite-Tower-600-G9-Desktop-PC:~$ hadoop jar /home/hadoop/Desktop/TopN.jar TopN.TNDriver /bda_
_hadoop/tinput.txt /bda_hadoop/toutput
2025-05-26 14:59:03,334 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-05-26 14:59:03,372 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-05-26 14:59:03,372 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2025-05-26 14:59:03,426 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. I
mplement the Tool interface and execute your application with ToolRunner to remedy this.
2025-05-26 14:59:03,472 INFO input.FileInputFormat: Total input files to process : 1
2025-05-26 14:59:03,497 INFO mapreduce.JobSubmitter: number of splits:1
2025-05-26 14:59:03,554 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local1824101299_0001
2025-05-26 14:59:03,554 INFO mapreduce.JobSubmitter: Executing with tokens: []
2025-05-26 14:59:03,609 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2025-05-26 14:59:03,610 INFO mapreduce.Job: Running job: job_local1824101299_0001
2025-05-26 14:59:03,610 INFO mapred.LocalJobRunner: OutputCommitter set in conf: null
2025-05-26 14:59:03,614 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting
to FileOutputCommitterFactory
2025-05-26 14:59:03,614 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-26 14:59:03,614 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders un
der output directory:false, ignore cleanup failures: false
2025-05-26 14:59:03,614 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.F
ileOutputCommitter
2025-05-26 14:59:03,654 INFO mapred.LocalJobRunner: Waiting for map tasks
2025-05-26 14:59:03,655 INFO mapred.LocalJobRunner: Starting task: attempt_local1824101299_0001_m_000000_0
2025-05-26 14:59:03,664 INFO output.PathOutputCommitterFactory: No output committer factory defined, defaulting
to FileOutputCommitterFactory
2025-05-26 14:59:03,664 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2025-05-26 14:59:03,664 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders un
der output directory:false, ignore cleanup failures: false
2025-05-26 14:59:03,670 INFO mapred.Task: Using ResourceCalculatorProcessTree : [ ]
2025-05-26 14:59:03,672 INFO mapred.MapTask: Processing split: hdfs://localhost:9000/bda_hadoop/tinput.txt:0+95
2025-05-26 14:59:03,701 INFO mapred.MapTask: (EQUATOR) 0 kvl 26214396(104857584)
2025-05-26 14:59:03,701 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
2025-05-26 14:59:03,701 INFO mapred.MapTask: soft limit at 83886080
2025-05-26 14:59:03,701 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
2025-05-26 14:59:03,701 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
2025-05-26 14:59:03,702 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapO
utputBuffer
2025-05-26 14:59:03,738 INFO mapred.LocalJobRunner:
2025-05-26 14:59:03,739 INFO mapred.MapTask: Starting flush of map output

```

```

File System Counters
  FILE: Number of bytes read=10682
  FILE: Number of bytes written=1291808
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=190
  HDFS: Number of bytes written=40
  HDFS: Number of read operations=15
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=4
  HDFS: Number of bytes read erasure-coded=0
Map-Reduce Framework
  Map input records=3
  Map output records=15
  Map output bytes=154
  Map output materialized bytes=190
  Input split bytes=108
  Combine input records=0
  Combine output records=0
  Reduce input groups=5
  Reduce shuffle bytes=190
  Reduce input records=15
  Reduce output records=5
  Spilled Records=30
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=0
  Total committed heap usage (bytes)=1052770304
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=95
File Output Format Counters
  Bytes Written=40
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -ls /bda_hadoop/toutput
Found 2 items
-rw-r--r--  1 hadoop supergroup          0 2025-05-26 14:59 /bda_hadoop/toutput/_SUCCESS
-rw-r--r--  1 hadoop supergroup        40 2025-05-26 14:59 /bda_hadoop/toutput/part-r-000000
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$ hdfs dfs -cat /bda_hadoop/toutput/part-r-000000
banana 5
apple  4
fruit  3
mango  2
kiwi   1
hadoop@bmscecse-HP-Elite-Tower-600-G9-Desktop-PC:~$

```