# Word Co-ocurrence

-Team 1(DS)

# The Team →


Prashant Pandey


Hrithik Kumar


Yashwant reddy


Jammula Supriya

# Agenda →

Introduction

Problem Statement

Map Reduce

Code Implementation

# Map Reduce

- Map reduce is programming model for data processing

- Framework for Parallel computing

- Allows one to to Process Huge amount of data(terabytes and petabytes) on thousands of processors

# Problems → and Solutions

## Problem
Word Co-occurence

## what is word co-occurence?
Word co-occurrence refers to the frequency with which two or more words appear together in a given text or dataset.

## Using MapReduce
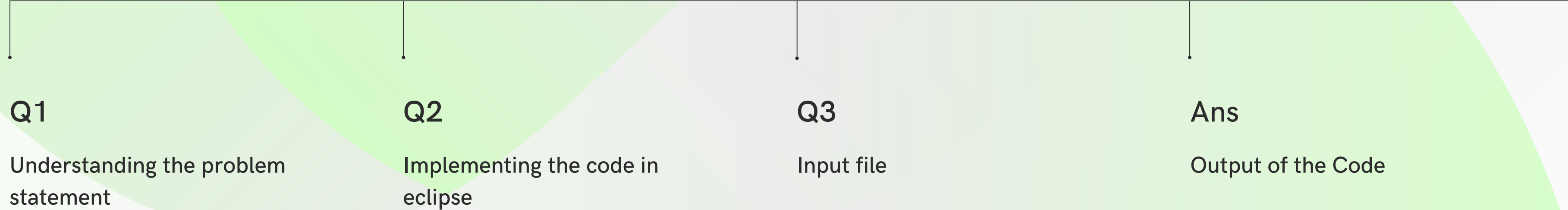Mapper phase ,Shuffle Phase, Reduce Phase

# Problems → and Solutions

The text from the input text file is tokenized into words to form a key value pair with all the words present in the input text file. The key is the word from the input file and value is '1'.
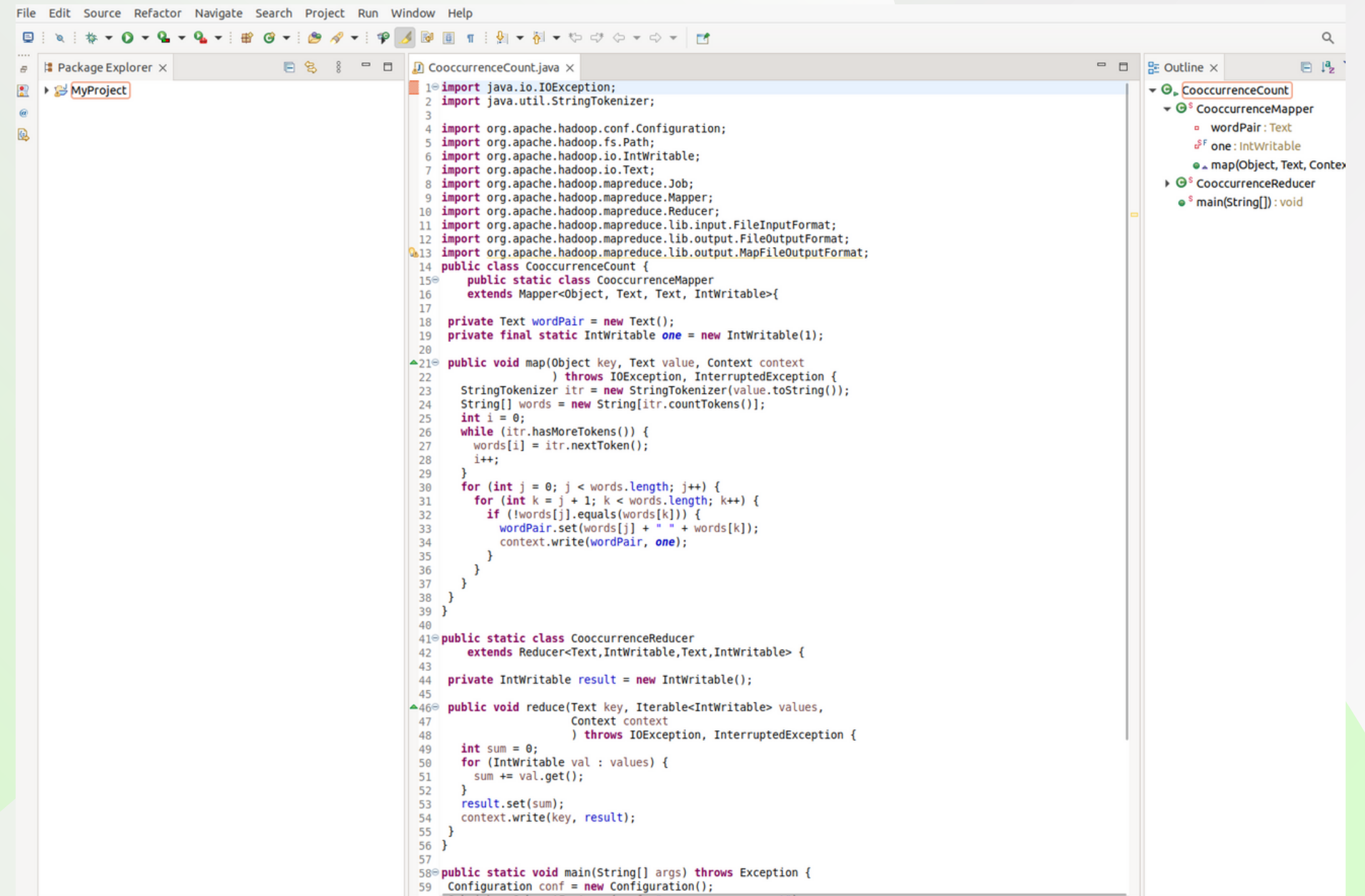
After the map phase execution is completed successfully, shuffle phase is executed automatically wherein the key-value pairs generated in the map phase are taken as input and then sorted in alphabetical order.es for a particular word.

In the reduce phase, all the keys are grouped together and the values for similar keys are added up to find the occurrences for a particular word

# Code → Implementation

**Q1**

Understanding the problem statement

**Q2**

Implementing the code in eclipse

**Q3**

Input file

**Ans**

Output of the Code

# Code Implementation

Package Explorer
CooccurrenceCount.java

Outline
- CooccurrenceCount
  - CooccurrenceMapper
    - wordPair : Text
    - one : IntWritable
    - map(Object, Text, Context
  - CooccurrenceReducer
  - main(String[]) : void

```java
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;
import org.apache.hadoop.mapreduce.lib.output.MapFileOutputFormat;
public class CooccurrenceCount {
    public static class CooccurrenceMapper
        extends Mapper<Object, Text, Text, IntWritable>{

    private Text wordPair = new Text();
    private final static IntWritable one = new IntWritable(1);

    public void map(Object key, Text value, Context context
                    ) throws IOException, InterruptedException {
        StringTokenizer itr = new StringTokenizer(value.toString());
        String[] words = new String[itr.countTokens()];
        int i = 0;
        while (itr.hasMoreTokens()) {
            words[i] = itr.nextToken();
            i++;
        }
        for (int j = 0; j < words.length; j++) {
            for (int k = j + 1; k < words.length; k++) {
                if (!words[j].equals(words[k])) {
                    wordPair.set(words[j] + " " + words[k]);
                    context.write(wordPair, one);
                }
            }
        }
    }
}

public static class CooccurrenceReducer
    extends Reducer<Text,IntWritable,Text,IntWritable> {

    private IntWritable result = new IntWritable();

    public void reduce(Text key, Iterable<IntWritable> values,
                       Context context
                       ) throws IOException, InterruptedException {
        int sum = 0;
        for (IntWritable val : values) {
            sum += val.get();
        }
        result.set(sum);
        context.write(key, result);
    }
}
public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
```

```java
public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "co-occurrence count");
    job.setJarByClass(CooccurrenceCount.class);
    job.setMapperClass(CooccurrenceMapper.class);
    job.setCombinerClass(CooccurrenceReducer.class);
    job.setReducerClass(CooccurrenceReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);
    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));
    System.exit(job.waitForCompletion(true) ? 0 : 1);
}
```

Writable          Smart Insert          1 : 28 : 27

# Output

```
247 pandey happy    6
248 pandey i        18
249 pandey is       3
250 pandey junk     6
251 pandey love     6
252 pandey name     3
253 pandey not      6
254 pandey person   6
255 pandey prashant         3
256 pandey to       6
257 person a        9
258 person am       9
259 person eat      6
260 person food     6
261 person good     3
262 person guy.     3
263 person guy.my   3
264 person happy    6
265 person i        15
266 person is       3
267 person junk     6
268 person love     6
269 person name     3
270 person not      3
271 person pandey   3
272 person prashant         3
273 person to       6
274 prashant a      12
275 prashant am     12
276 prashant eat    6
277 prashant food   6
278 prashant good   6
279 prashant guy.   3
280 prashant guy.my         3
281 prashant happy  6
282 prashant i      18
283 prashant is     3
284 prashant junk   6
285 prashant love   6
286 prashant name   3
287 prashant not    6
288 prashant pandey         6
289 prashant person         6
290 prashant to     6
291 to a    9
292 to am   9
293 to eat  6
294 to food         6
295 to good         3
296 to guy.         3
297 to guy.my       3
298 to happy        6
299 to i    12
300 to is   3
301 to junk         6
302 to love         3
303 to name         3
304 to not  3
305 to pandey       3
306 to person       3
307 to prashant     3
```

# Output File Stored

# Thank You