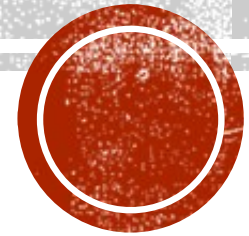


MAPREDUCE — CONT'D



MAPREDUCE — WHAT YOU NEED TO DO

- Programmers specify two functions:

map $(k, v) \rightarrow \langle k', v' \rangle^*$

reduce $(k', v') \rightarrow \langle k', v' \rangle^*$

- All values with the same key are reduced together
- The execution framework handles everything else!
- Not quite...usually, programmers also specify:

partition $(k', \text{number of partitions}) \rightarrow \text{partition for } k'$

 - Often a simple hash of the key, e.g., $\text{hash}(k') \bmod n$
 - Divides up key space for parallel reduce operations

combine $(k', v') \rightarrow \langle k', v' \rangle^*$

 - Mini-reducers that run in memory after the map phase
 - Used as an optimization to reduce network traffic



WORDCOUNT PSEUDOCODE

Map

```
map(String file, String doc)
{
    List<String> T = tokenize(doc);
    for each token in T
    {
        emit ((String)token, (Int) 1);
    }
}
```

Reduce

```
reduce(String token, List<Int> values)
{
    Integer sum = 0;
    for each value in values {
        sum = sum + value;
    }
    emit ((String)token, (Integer) sum);
}
```



WORDCOUNT – MAPREDUCE VERSION

Mapper

```
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();  
  
        public void map(Object key, Text value, Context context  
            ) throws IOException, InterruptedException {  
            StringTokenizer itr = new StringTokenizer(value.toString());  
            while (itr.hasMoreTokens()) {  
                word.set(itr.nextToken());  
                context.write(word, one);  
            }  
        }  
    }  
  
    public static class IntSumReducer  
        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);  
        }  
    }  
  
}
```

Reducer

Driver

```
public static void main(String[] args) throws Exception {  
    Configuration conf = new Configuration();  
    Job job = Job.getInstance(conf, "word count");  
    job.setJarByClass(WordCount.class);  
    job.setMapperClass(TokenizerMapper.class);  
    job.setCombinerClass(IntSumReducer.class);  
    job.setReducerClass(IntSumReducer.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);  
}
```



CODE EXPLAIN - MAPPER

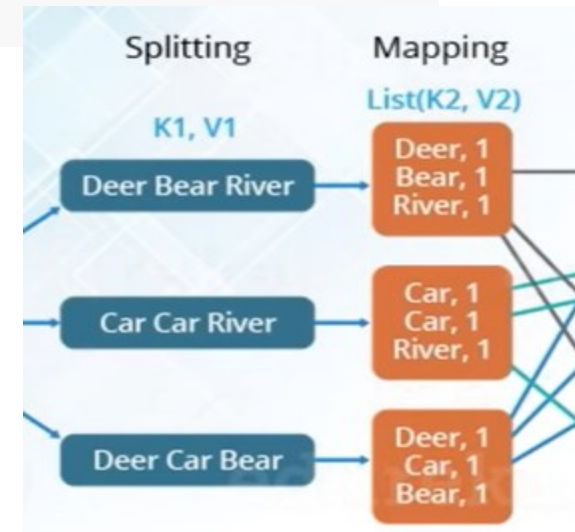
```
public static class TokenizerMapper
    extends Mapper<Object, Text, Text, IntWritable>{

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

    public void map(Object key, Text value, Context context
        ) throws IOException, InterruptedException {
        StringTokenizer itr = new StringTokenizer(value.toString());
        while (itr.hasMoreTokens()) {
            word.set(itr.nextToken());
            context.write(word, one);
        }
    }
}
```

Key-in, Value-in, Key-out, Value-out

- The Mapper implementation, via the map method, processes one line at a time, as provided by the specified *Text* Input Format
- It then splits the line into tokens separated by whitespaces, via the *StringTokenizer*, and
- emits a key-value pair of < <word>, 1>.



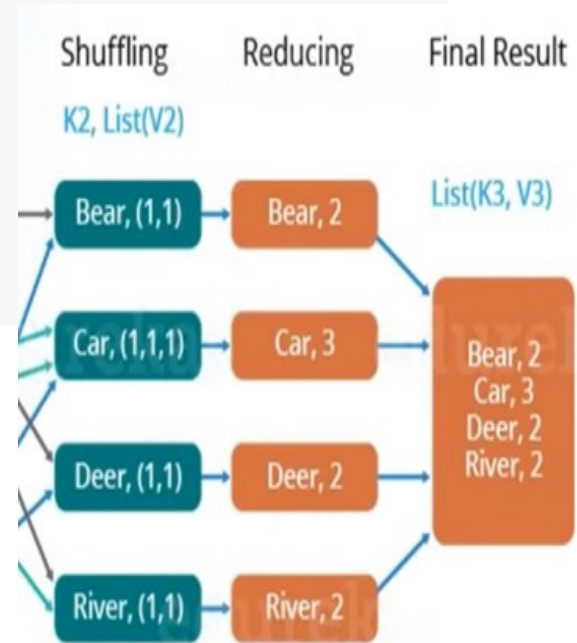
CODE EXPLAIN - REDUCER

```
public static class IntSumReducer
    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);
    }
}
```

- Hadoop takes care of the Shuffling process prior to the reducer gets executed.
- The Reducer implementation, via the reduce method just sums up the values, which are the occurrence counts for each key (i.e. words in this example).

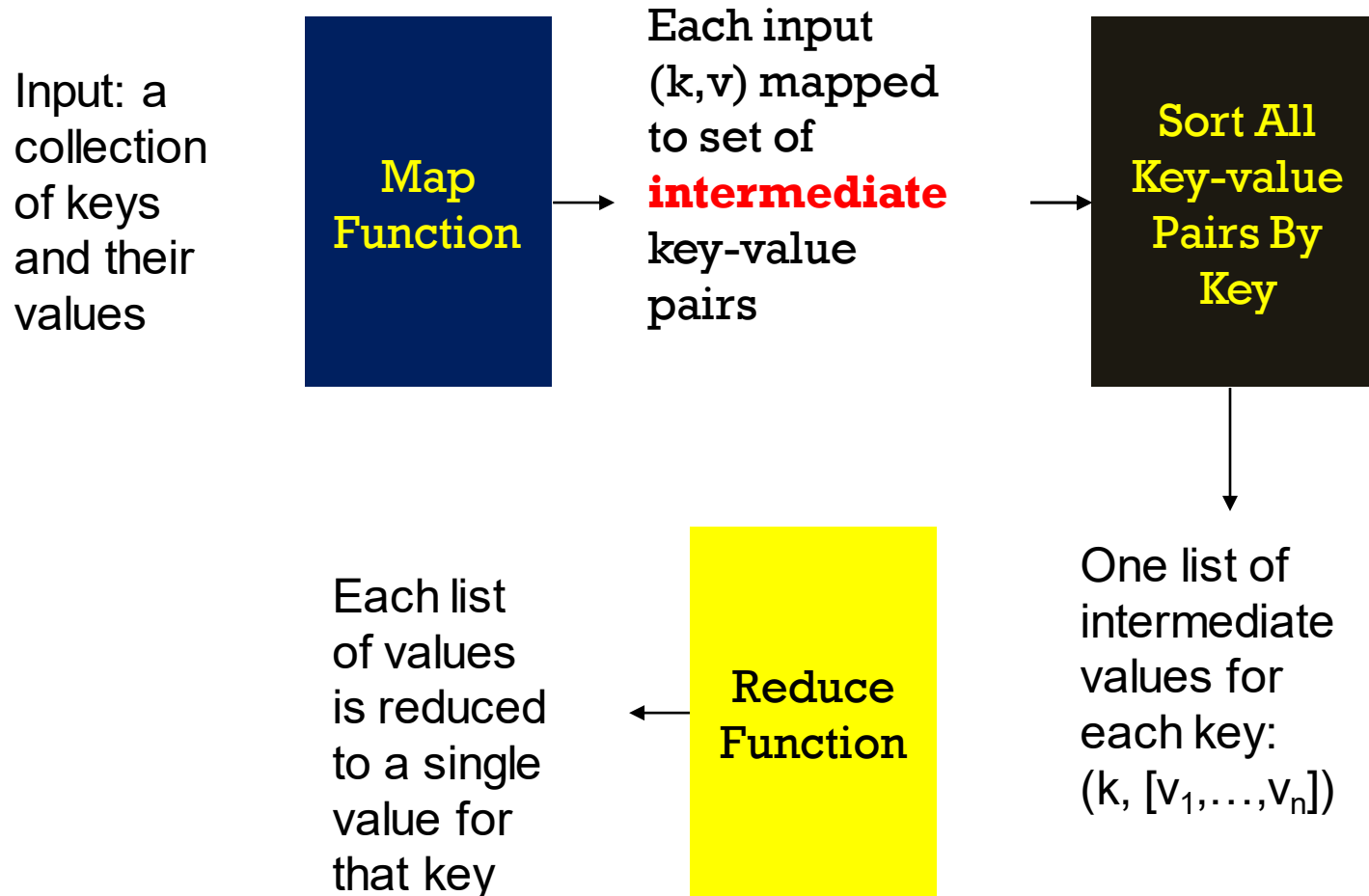


CODE EXPLAIN – MAIN (DRIVER)

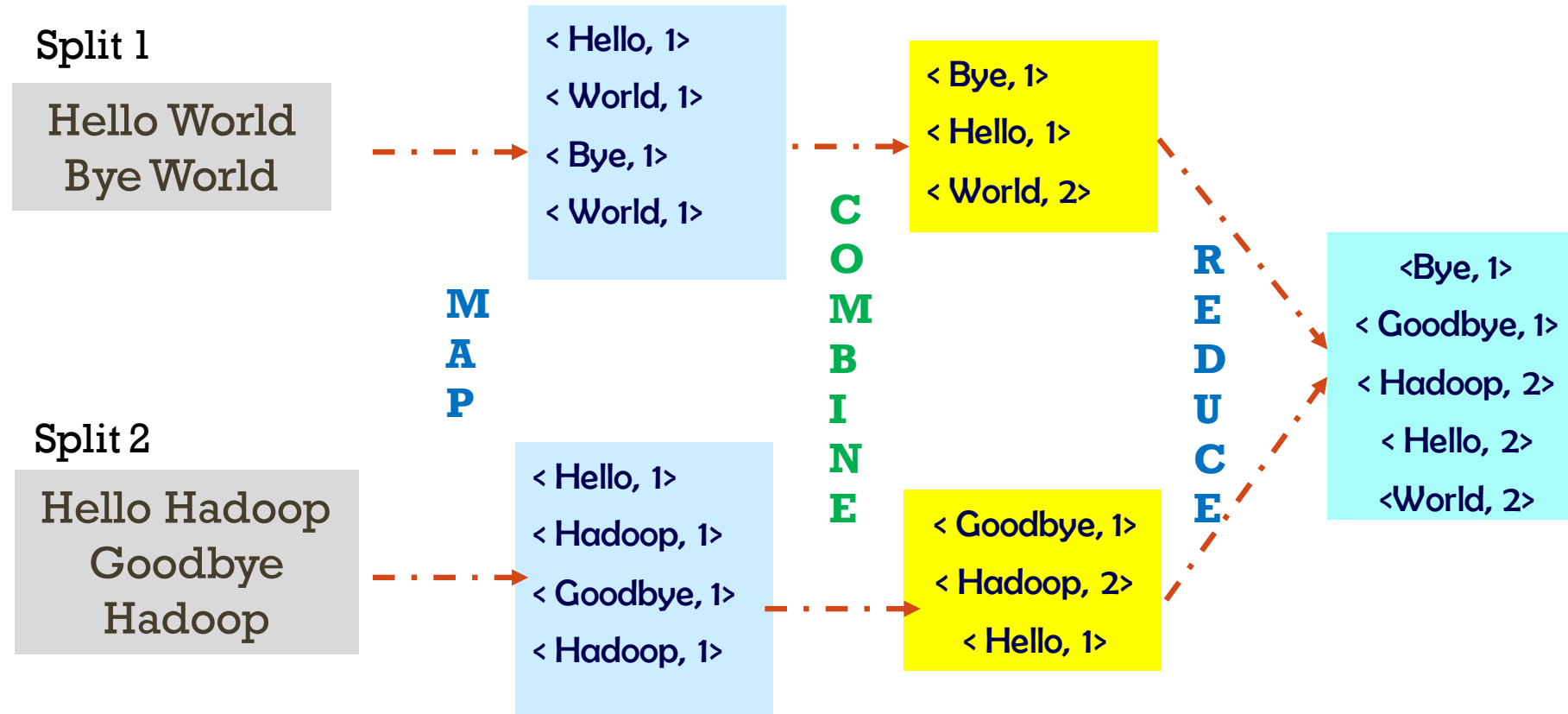
```
public static void main(String[] args) throws Exception {  
    Configuration conf = new Configuration();  
    Job job = Job.getInstance(conf, "word count");  
    job.setJarByClass(WordCount.class);  
    job.setMapperClass(TokenizerMapper.class);  
    job.setCombinerClass(IntSumReducer.class);  
    job.setReducerClass(IntSumReducer.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);  
}
```



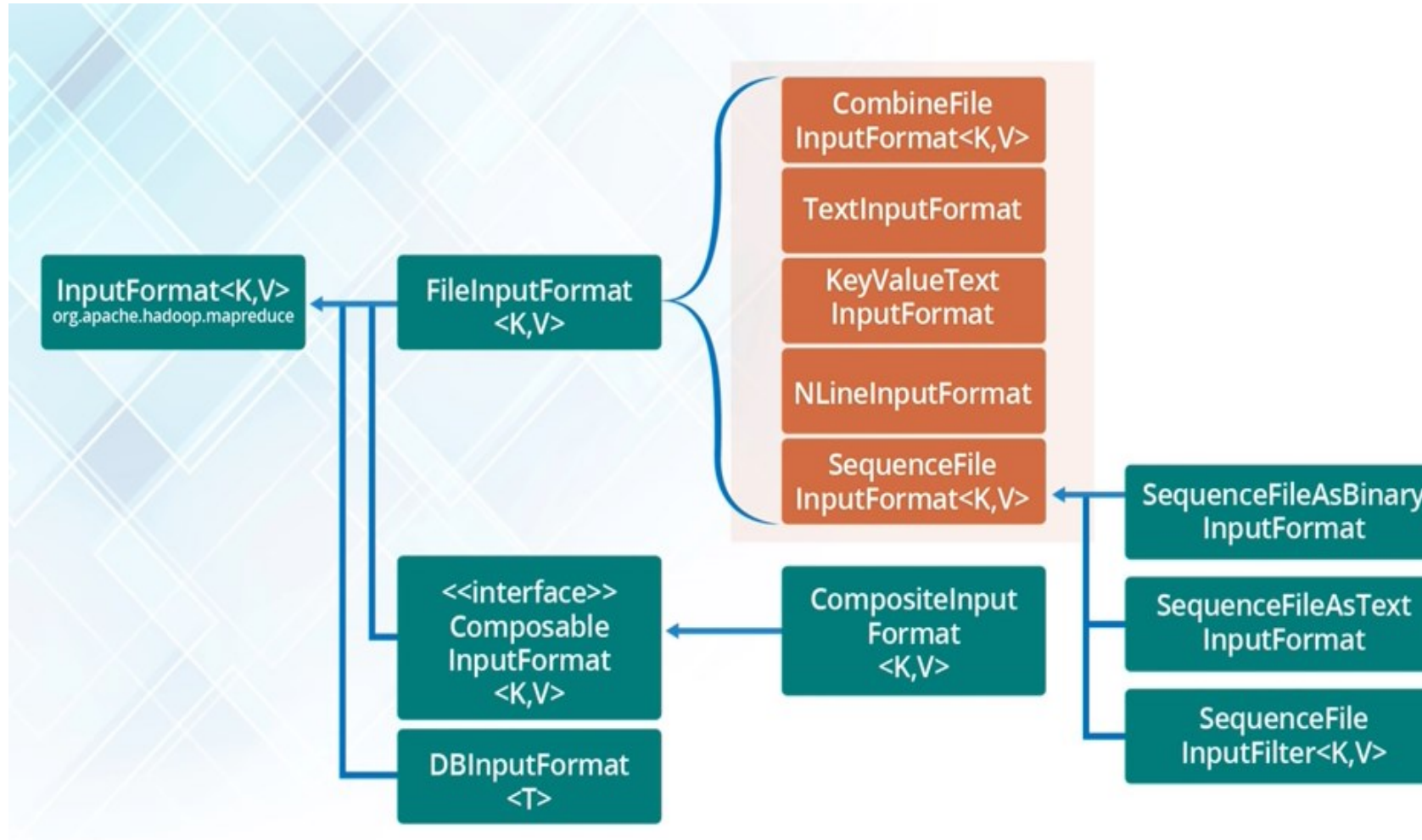
MAP-REDUCE FRAMEWORK FOR WORDCOUNT



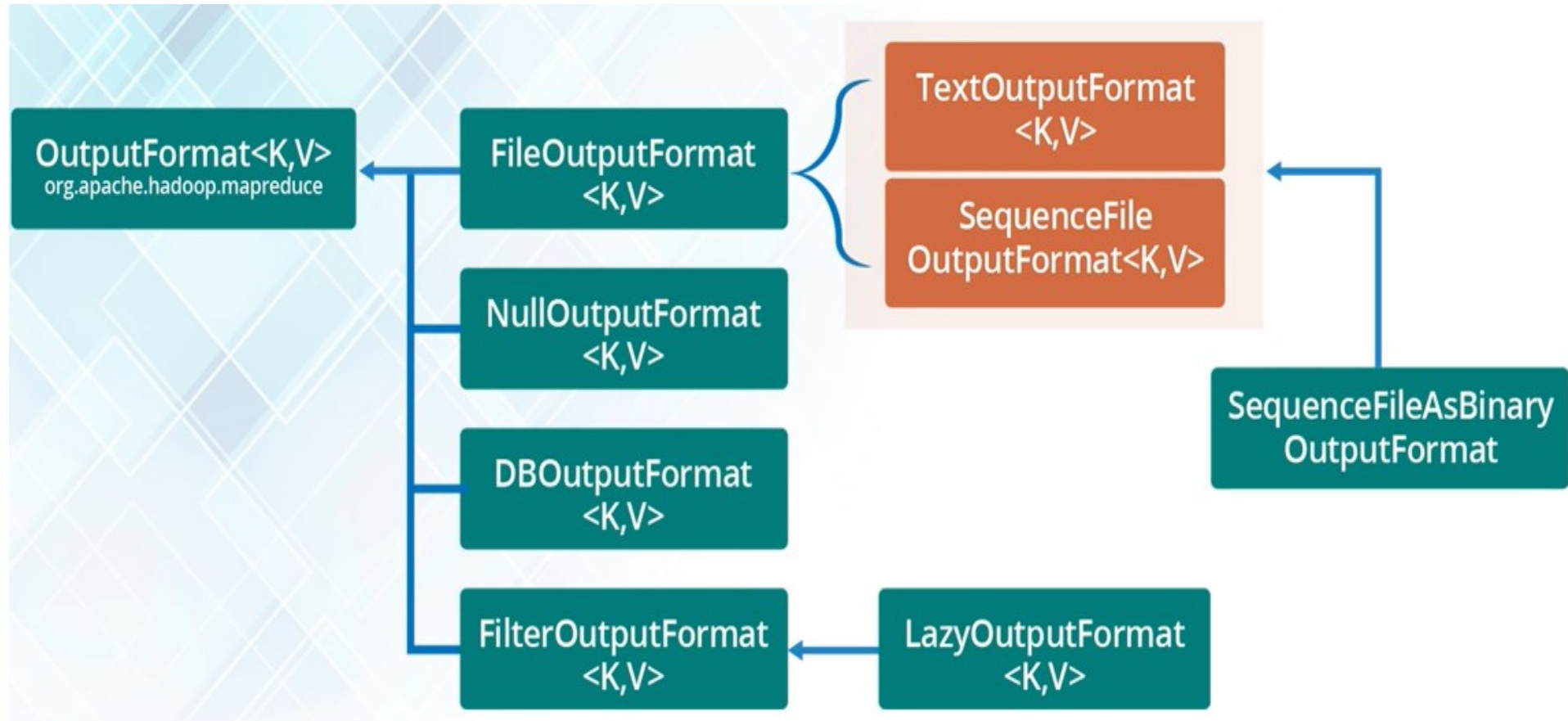
EXAMPLE – EXECUTION FLOW



MAPREDUCE INPUT FORMAT – CLASS HIERARCHY



MAPREDUCE OUTPUT FORMAT – CLASS HIERARCHY



HADOOP MAPREDUCE - JAVA MAIN PACKAGES

```
import java.io.IOException;  
import java.util.*;
```

```
import org.apache.hadoop.fs.Path;  
import org.apache.hadoop.conf.*;  
import org.apache.hadoop.io.*;
```

All these packages are present in
hadoop-common.jar

```
import org.apache.hadoop.mapreduce.*;  
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;  
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;  
import  
org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;  
import  
org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
```

All these
packages are
present in
**hadoop-mapreduce-
client-core.jar**

