# PART – 4

Write one MapReduce program using each of the classes that extend FileInputFormat<k,v>
(CombineFileInputFormat, FixedLengthInputFormat, KeyValueTextInputFormat,
NLineInputFormat, SequenceFileInputFormat, TextInputFormat)



## CombineFileInputFormat:

This code example uses the CombineFileInputFormat class, which allows combining multiple small input files into larger splits to improve efficiency.

The WordCountMapper class reads each line from the input file and tokenizes it into words.

The WordCountReducer class receives the words as keys and counts their occurrences.

The program counts the occurrence of each word in the input files and writes the result to the output file.

The significance of using CombineFileInputFormat is to optimize the processing of small files by combining them into larger input splits, reducing overhead.

**Code: (Full code is submitted in the part-4 file with reducer & main method)**

public class CombineFileInputFormatExample {

```java
public static class WordCountMapper extends Mapper<LongWritable, Text, Text,
IntWritable> {

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

    private Text word = new Text();

public void map(LongWritable 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);

        }

    }

  }
```

## FixedLengthInputFormat:

This code example uses the FixedLengthInputFormat class, which reads fixed-length records from the input file.

The WordCountMapper class reads each character from the input file and emits it as a separate key-value pair.

The WordCountReducer class receives the characters as keys and counts their occurrences.

The program counts the occurrence of each character in the input file and writes the result to the output file.

The significance of using FixedLengthInputFormat is to process input files with fixed-length records, such as files with fixed-width columns.

**Code: (Full code is submitted in the part-4 file with reducer & main method)**

```java
public class FixedLengthInputFormatExample {

public static class WordCountMapper extends Mapper<LongWritable, Text, Text,
IntWritable> {

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

    private Text word = new Text();
```

```
 public void map(LongWritable key, Text value, Context context) throws IOException,
InterruptedException {

        String line = value.toString();

        word.set(line);

        context.write(word, one);

    }

  }
```

## KeyValueTextInputFormat:

This code example uses the KeyValueTextInputFormat class, which reads input files in key-value pair format, where each line contains a key and a value separated by a delimiter (default is tab).

The WordCountMapper class reads each value from the input file and tokenizes it into words.

The WordCountReducer class receives the words as keys and counts their occurrences.

The program counts the occurrence of each word in the input file and writes the result to the output file.

The significance of using KeyValueTextInputFormat is to process input files with key-value pairs, such as log files or configuration files.

**Code: (Full code is submitted in the part-4 file with reducer & main method)**

```
public class KeyValueTextInputFormatExample {

  public static class WordCountMapper extends Mapper<Text, Text, Text, IntWritable> {

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

    private Text word = new Text();

public void map(Text key, Text value, Context context) throws IOException,
InterruptedException {

        String line = value.toString();

        String[] words = line.split("\\s+"); \\ Regex one or more whitespace characters (spaces,
tabs, or line breaks).

        for (String word : words) {

            this.word.set(word);

            context.write(this.word, one);
```

```
        }

    }

}
```

## NLineInputFormat:

This code example uses the NLineInputFormat class, which reads N lines of input as a split.

The WordCountMapper class reads each line from the input file and tokenizes it into words.

The WordCountReducer class receives the words as keys and counts their occurrences.

The program counts the occurrence of each word in the input file and writes the result to the output file.

The significance of using NLineInputFormat is to control the number of lines per input split, which can be useful when the input data has a specific structure or when you want to process a specific number of lines at once.

**Code: (Full code is submitted in the part-4 file with reducer & main method)**

```
public class NLineInputFormatExample {

    public static class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {

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

        private Text word = new Text();

        public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {

            String line = value.toString();

            String[] words = line.split("\\s+"); \\ Regex one or more whitespace characters (spaces, tabs, or line breaks).

            for (String word : words) {

                this.word.set(word);

                context.write(this.word, one);

            }

        }

    }
```

## SequenceFileInputFormat:

This code example uses the SequenceFileInputFormat class, which reads input files stored in the SequenceFile format, which is a binary file format that stores key-value pairs.

The WordCountMapper class reads each key from the input file and tokenizes it into words.

The WordCountReducer class receives the words as keys and counts their occurrences.

The program counts the occurrence of each word in the input file and writes the result to the output file.

The significance of using SequenceFileInputFormat is to process input files stored in the SequenceFile format, which is commonly used in Hadoop for efficient data storage and retrieval.

**Code: (Full code is submitted in the part-4 file with reducer & main method)**

```
public class SequenceFileInputFormatExample {

    public static class WordCountMapper extends Mapper<Text, IntWritable, Text,
IntWritable> {

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

        public void map(Text key, IntWritable value, Context context) throws IOException,
InterruptedException {

            context.write(key, one);

        }

    }
```

## TextInputFormat:

This code example uses the TextInputFormat class, which reads input files as plain text files, where each line is considered a record.

The WordCountMapper class reads each line from the input file and tokenizes it into words.

The WordCountReducer class receives the words as keys and counts their occurrences.

The program counts the occurrence of each word in the input file and writes the result to the output file.

The significance of using TextInputFormat is to process plain text files where each line represents a record, which is a common format for textual data processing.

**Code: (Full code is submitted in the part-4 file with reducer & main method)**

```
public class TextInputFormatExample {
```

```java
    public static class WordCountMapper extends Mapper<LongWritable, Text, Text,
IntWritable> {

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

        private Text word = new Text();

        public void map(LongWritable 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);

            }

        }

    }
```