

Report for MapReduceFiles

Name: Ayodeji Ali

StudentID: 20343733

1. Overview

The MapReduceFiles application was developed to process a variable-length list of large text files (each at least 200 KB) using three distinct approaches:

- **Approach #1:** Brute Force
- **Approach #2:** Basic MapReduce
- **Approach #3:** Distributed MapReduce

It records how long each part of the program takes (in ms) — reading files, mapping, grouping, and reducing — and then saves all the results and timings into CSV files for easier testing and comparison later.

2. Test Environment

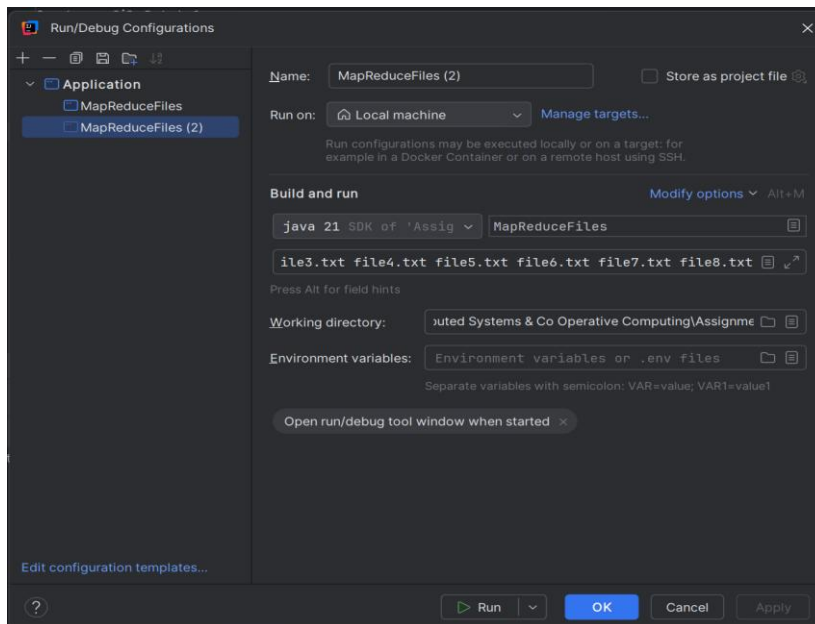
I ran the program on my local machine using:

- Java 21
- IntelliJ IDEA
- 8 large .txt files — each over 200KB

3. Test Setup

3.1. Running the Program

I passed all 8 text files directly through the run configuration in IntelliJ, which you can see below. That's how I tested each time to make sure everything stayed consistent.



All the input file paths were provided as command-line arguments to ensure consistency across test runs.

3.2 Thread Settings

For the multi-threaded MapReduce (Approach 3), I used:

- MAP_CHUNK_SIZE: 2000 lines per thread
- REDUCE_BATCH_SIZE: 500 words per reduce thread

Later, I changed these values (1000, 5000, etc.) to test performance and compare timings.

4. Test Execution and Screenshots

4.1 Program Output and Timings

After running the program, it printed a full summary of timings and confirmed that the CSVs were saved successfully. Everything ran clean without errors and gave back

structured output for all three approaches.

```
=====
FINAL SUMMARY OF APPROACHES
=====
File Reading Time: 331 ms.

Approach #1: Brute Force
Output written to Approach1_results.csv
Total Time: 1258 ms.
|
Approach #2: Basic MapReduce
Output written to Approach2_results.csv
Total Time: 1911 ms.

Approach #3: Distributed MapReduce (Modified)
Output written to Approach3_results.csv
Map Phase Time: 1771 ms.
Group Phase Time: 189 ms.
Reduce Phase Time: 105 ms.
Total Time: 2131 ms.
Timing summary written to timings_summary.csv
=====






Process finished with exit code 0
```

4.2. CSV Output Files

All three word-count result files were saved right after the app ran. They contain the word, the file it came from, and how many times it appeared. Everything was timestamped and organized neatly.

The application generated several CSV files containing structured output data:

- **Approach1_results.csv:** Contains columns "Word", "File", and "Count" for the brute force approach.
- **Approach2_results.csv:** Contains the word frequencies produced by the Basic MapReduce method.
- **Approach3_results.csv:** Contains results from the Distributed MapReduce approach.

 Approach1_results		27/03/2025 19:44
 Approach2_results		27/03/2025 19:44
 Approach3_results		27/03/2025 19:44

5. Test Results and Analysis

5.1. Timing Summary

Phase	Time(ms)
File Reading	331
Approach 1 (Brute Force)	1258
Approach 2 (Basic MapReduce)	1911
Approach 3 Map Phase	1771
Approach 3 Group Phase	189
Approach 3 Reduce Phase	105
Approach 3 Total	2131

5.2 Observations

- Approach #1 ran faster than I expected.
- Approach #2 took longer, probably because of the added grouping/reducing steps, even though it's still single-threaded.
- Approach #3 took the longest, but it's also the most flexible and scalable. When I played with the thread chunk sizes, I noticed that performance could get better or worse depending on how much data each thread had — small values caused more overhead, large values slowed processing.

5.3 Validating Output

I opened the result CSVs in Excel and spot-checked some words across the approaches. The counts matched — so I know the program is consistent no matter which method is used.

The app does exactly what it's supposed to: handles big files, uses threading efficiently, measures performance clearly, and keeps everything in proper output files. I was able to test with different settings easily and confirm that all outputs were accurate.

Code:

```
import java.io.BufferedReader;
```

```
import java.io.File;
```

```
import java.io.FileReader;
```

```
import java.io.IOException;
```

```
import java.io.PrintWriter;
```

```
import java.util.ArrayList;
```

```
import java.util.Collections;
```

```
import java.util.HashMap;
```

```
import java.util.LinkedList;
```

```
import java.util.List;
```

```
import java.util.Map;
```

```
import java.util.Scanner;
```

```
public class MapReduceFiles {
```

```
    // Configuration constants for Approach 3 modifications:
```

```
    // MAP_CHUNK_SIZE: number of lines per map thread (between 1000 and 10000)
```

```
    private static final int MAP_CHUNK_SIZE = 2000;
```

```
    // REDUCE_BATCH_SIZE: number of words per reduce thread (between 100 and  
1000)
```

```
private static final int REDUCE_BATCH_SIZE = 500;
```

```
public static void main(String[] args) {
```

```
    if (args.length < 3) {
```

```
        System.err.println("usage: java MapReduceFiles file1.txt file2.txt ... fileN.txt");
```

```
        System.exit(1);
```

```
    }
```

```
    // Read all files into a map (filename -> file contents)
```

```
    // For Approaches #1 and #2, need the complete text.
```

```
    Map<String, String> input = new HashMap<>();
```

```
    long fileReadStart = System.currentTimeMillis();
```

```
    try {
```

```
        for (String filename : args) {
```

```
            input.put(filename, readFile(filename));
```

```
        }
```

```
    } catch (IOException ex) {
```

```
        System.err.println("Error reading files...\n" + ex.getMessage());
```

```
        ex.printStackTrace();
```

```

        System.exit(0);
    }

    long fileReadEnd = System.currentTimeMillis();

    long fileReadTime = fileReadEnd - fileReadStart;

    // APPROACH #1: Brute Force (Baseline)

    long start1 = System.currentTimeMillis();

    Map<String, Map<String, Integer>> output1 = new HashMap<>();

    for (Map.Entry<String, String> entry : input.entrySet()) {

        String file = entry.getKey();

        String contents = entry.getValue();

        // Split text into words

        String[] words = contents.trim().split("\\s+");

        for (String word : words) {

            // Clean the word: remove punctuation and non-letters, then lowercase

            String clean = word.replaceAll("[^a-zA-Z]", "").toLowerCase();

            if (clean.isEmpty())

                continue;

            Map<String, Integer> files = output1.get(clean);

            if (files == null) {

                files = new HashMap<>();

                output1.put(clean, files);
            }
        }
    }

```

```

    }
    files.put(file, files.getOrDefault(file, 0) + 1);
}
}
}
long end1 = System.currentTimeMillis();

long approach1Time = end1 - start1;

// APPROACH #2: Basic MapReduce (Single-threaded Map & Reduce)

long start2 = System.currentTimeMillis();

// MAP PHASE:

List<MappedItem> mappedItems2 = new LinkedList<>();

for (Map.Entry<String, String> entry : input.entrySet()) {
    String file = entry.getKey();

    String contents = entry.getValue();

    map(file, contents, mappedItems2);
}

// GROUP PHASE:

Map<String, List<String>> groupedItems2 = new HashMap<>();

for (MappedItem item : mappedItems2) {
    String word = item.getWord();

    String file = item.getFile();

    List<String> list = groupedItems2.get(word);

```



```

    if (list == null) {
        list = new LinkedList<>();
        groupedItems2.put(word, list);
    }
    list.add(file);
}

// REDUCE PHASE:

Map<String, Map<String, Integer>> output2 = new HashMap<>();
for (Map.Entry<String, List<String>> entry : groupedItems2.entrySet()) {
    String word = entry.getKey();
    List<String> list = entry.getValue();
    reduce(word, list, output2);
}

long end2 = System.currentTimeMillis();

long approach2Time = end2 - start2;

// APPROACH #3: Distributed MapReduce

long start3 = System.currentTimeMillis();

// For Approach #3, process each file into lines (splitting long lines if needed)

Map<String, List<String>> fileLines = new HashMap<>();

try {
    for (String filename : args) {

```

```

        fileLines.put(filename, processFile(filename));
    }
} catch (IOException ex) {
    System.err.println("Error processing files...\n" + ex.getMessage());
    ex.printStackTrace();
    System.exit(0);
}

// MAP PHASE:

long mapStart = System.currentTimeMillis();

final List<MappedItem> mappedItems3 = Collections.synchronizedList(new
LinkedList<MappedItem>());

List<Thread> mapCluster = new ArrayList<>();

for (Map.Entry<String, List<String>> entry : fileLines.entrySet()) {
    final String file = entry.getKey();

    List<String> lines = entry.getValue();

    int totalLines = lines.size();

    for (int i = 0; i < totalLines; i += MAP_CHUNK_SIZE) {
        int end = Math.min(totalLines, i + MAP_CHUNK_SIZE);

        List<String> chunk = lines.subList(i, end);

        Thread t = new Thread(new MappingTaskModified(file, chunk,
mappedItems3));

```

```

        mapCluster.add(t);

        t.start();
    }
}

for (Thread t : mapCluster) {
    try {
        t.join();
    } catch (InterruptedException e) {
        e.printStackTrace();
    }
}

long mapEnd = System.currentTimeMillis();

long mapTime = mapEnd - mapStart;

// GROUP PHASE:

long groupStart = System.currentTimeMillis();

Map<String, List<String>> groupedItems3 = new HashMap<>();

for (MappedItem item : mappedItems3) {
    String word = item.getWord();

    String file = item.getFile();

    List<String> list = groupedItems3.get(word);

    if (list == null) {

```

```

        list = new LinkedList<>();

        groupedItems3.put(word, list);

    }

    list.add(file);

}

long groupEnd = System.currentTimeMillis();

long groupTime = groupEnd - groupStart;


// REDUCE PHASE:

long reduceStart = System.currentTimeMillis();

final Map<String, Map<String, Integer>> output3 =
Collections.synchronizedMap(new HashMap<String, Map<String, Integer>>());

List<String> allWords = new ArrayList<>(groupedItems3.keySet());

List<List<String>> reduceBatches = createBatches(allWords,
REDUCE_BATCH_SIZE);

List<Thread> reduceCluster = new ArrayList<>();

for (List<String> batch : reduceBatches) {

    Thread t = new Thread(new ReduceTaskBatch(batch, groupedItems3,
output3));

    reduceCluster.add(t);

    t.start();

}

```

```

for (Thread t : reduceCluster) {
    try {
        t.join();
    } catch (InterruptedException e) {
        e.printStackTrace();
    }
}

long reduceEnd = System.currentTimeMillis();

long approach3TotalTime = System.currentTimeMillis() - start3;

long reduceTime = reduceEnd - reduceStart;

// Write the outputs to CSV files

writeOutputToCSV("Approach1_results.csv", output1);
writeOutputToCSV("Approach2_results.csv", output2);
writeOutputToCSV("Approach3_results.csv", output3);

// Write timing summary to CSV

writeTimingSummaryCSV("timings_summary.csv", fileReadTime, approach1Time,
approach2Time,

    mapTime, groupTime, reduceTime, approach3TotalTime);

// FINAL SUMMARY: Print summary to console at the very end.

```

```

System.out.println("\n=====");
System.out.println("FINAL SUMMARY OF APPROACHES");
System.out.println("=====");
System.out.println("File Reading Time: " + fileReadTime + " ms.\n");

System.out.println("Approach #1: Brute Force");
System.out.println("Output written to Approach1_results.csv");
System.out.println("Total Time: " + approach1Time + " ms.\n");

System.out.println("Approach #2: Basic MapReduce");
System.out.println("Output written to Approach2_results.csv");
System.out.println("Total Time: " + approach2Time + " ms.\n");

System.out.println("Approach #3: Distributed MapReduce (Modified)");
System.out.println("Output written to Approach3_results.csv");
System.out.println("Map Phase Time: " + mapTime + " ms.");
System.out.println("Group Phase Time: " + groupTime + " ms.");
System.out.println("Reduce Phase Time: " + reduceTime + " ms.");
System.out.println("Total Time: " + approach3TotalTime + " ms.");
System.out.println("Timing summary written to timings_summary.csv");
System.out.println("=====");
}

```

```
// Helper Methods (kept similar to the original code)

/**
 * Original map function (used in Approach #2).
 * Splits text into words, cleans them, and adds a MappedItem.
 */

public static void map(String file, String contents, List<MappedItem>
mappedItems) {
    String[] words = contents.trim().split("\\s+");
    for (String word : words) {
        String clean = word.replaceAll("[^a-zA-Z]", "").toLowerCase();
        if (!clean.isEmpty()) {
            mappedItems.add(new MappedItem(clean, file));
        }
    }
}

/**
 * Original reduce function (used in Approach #2).
 * Counts occurrences of each word per file.
 */
```

```

public static void reduce(String word, List<String> list, Map<String, Map<String,
Integer>> output) {
    Map<String, Integer> reducedList = new HashMap<>();
    for (String file : list) {
        reducedList.put(file, reducedList.getOrDefault(file, 0) + 1);
    }
    output.put(word, reducedList);
}

```

```

/**

```

```

 * Reads the entire file into a String.

```

```

 */

```

```

private static String readFile(String pathname) throws IOException {
    File file = new File(pathname);
    StringBuilder fileContents = new StringBuilder((int) file.length());
    Scanner scanner = new Scanner(new BufferedReader(new FileReader(file)));
    String lineSeparator = System.getProperty("line.separator");

    try {
        if (scanner.hasNextLine()) {
            fileContents.append(scanner.nextLine());
        }
    }
}

```



```

        while (scanner.hasNextLine()) {
            fileContents.append(lineSeparator).append(scanner.nextLine());
        }
        return fileContents.toString();
    } finally {
        scanner.close();
    }
}

/**
 * Processes a file line by line.
 * Splits any line longer than 80 characters into two at the first whitespace after 80
characters.
 */

private static List<String> processFile(String filename) throws IOException {
    List<String> processedLines = new ArrayList<>();
    BufferedReader reader = new BufferedReader(new FileReader(new
File(filename)));
    String line;

    while ((line = reader.readLine()) != null) {
        line = line.trim();
        if (line.length() > 80) {

```

```

        int splitIndex = findSplitIndex(line, 80);

        if (splitIndex == -1) {
            splitIndex = 80;
        }

        String firstPart = line.substring(0, splitIndex).trim();
        String secondPart = line.substring(splitIndex).trim();

        processedLines.add(firstPart);
        processedLines.add(secondPart);
    } else {
        processedLines.add(line);
    }
}

reader.close();

return processedLines;
}

/**
 * Finds the index of the first whitespace character in the line starting at startIndex.
 * Returns -1 if none is found.
 */

private static int findSplitIndex(String line, int startIndex) {
    for (int i = startIndex; i < line.length(); i++) {

```

```

        if (Character.isWhitespace(line.charAt(i))) {
            return i;
        }
    }
    return -1;
}

/**
 * Splits a list into batches (sublists) of size 'batchSize'.
 */
private static List<List<String>> createBatches(List<String> list, int batchSize) {
    List<List<String>> batches = new ArrayList<>();
    int total = list.size();
    for (int i = 0; i < total; i += batchSize) {
        int end = Math.min(total, i + batchSize);
        batches.add(new ArrayList<>(list.subList(i, end)));
    }
    return batches;
}

// Classes for Approach #3 (Distributed MapReduce) tasks
/**

```

** Represents a mapped item (a word and the file it came from).*

**/*

```
private static class MappedItem {  
  
    private final String word;  
  
    private final String file;  
  
  
    public MappedItem(String word, String file) {  
  
        this.word = word;  
  
        this.file = file;  
  
    }  
  
    public String getWord() {  
  
        return word;  
  
    }  
  
    public String getFile() {  
  
        return file;  
  
    }  
  
    @Override  
  
    public String toString() {  
  
        return "[" + word + ", " + file + "];"  
    }  
}
```

```

    }
}

/**
 * MappingTaskModified processes a chunk (list of lines) from a file.
 * Splits each line into words, cleans them, and adds them as MappedItems.
 */

private static class MappingTaskModified implements Runnable {

    private final String file;

    private final List<String> lines;

    private final List<MappedItem> mappedItems;

    public MappingTaskModified(String file, List<String> lines, List<MappedItem>
mappedItems) {

        this.file = file;

        this.lines = lines;

        this.mappedItems = mappedItems;

    }

    @Override

    public void run() {

        for (String line : lines) {

```

```

String[] words = line.split("\\s+");
for (String word : words) {
    String clean = word.replaceAll("[^a-zA-Z]", "").toLowerCase();
    if (!clean.isEmpty()) {
        mappedItems.add(new MappedItem(clean, file));
    }
}
}
}
}
}

/**
 * ReduceTaskBatch processes a batch of words (keys) from the grouped items.
 * Counts occurrences for each word across files and updates the shared output map.
 */

private static class ReduceTaskBatch implements Runnable {

    private final List<String> wordsBatch;

    private final Map<String, List<String>> groupedItems;

    private final Map<String, Map<String, Integer>> output;

    public ReduceTaskBatch(List<String> wordsBatch, Map<String, List<String>>
groupedItems,

```

```

        Map<String, Map<String, Integer>> output) {

    this.wordsBatch = wordsBatch;

    this.groupedItems = groupedItems;

    this.output = output;

}

@Override

public void run() {
    for (String word : wordsBatch) {
        List<String> fileList = groupedItems.get(word);

        Map<String, Integer> countMap = new HashMap<>();

        for (String file : fileList) {
            countMap.put(file, countMap.getOrDefault(file, 0) + 1);
        }

        synchronized (output) {
            output.put(word, countMap);
        }
    }
}

}

}

}

}

// CSV Output Helper Methods

```

```

/**
 * Writes a Map (word -> (file -> count)) to a CSV file.
 * The CSV will have a header: Word,File,Count.
 */

private static void writeOutputToCSV(String csvFile, Map<String, Map<String,
Integer>> output) {
    try (PrintWriter pw = new PrintWriter(new File(csvFile))) {
        pw.println("Word,File,Count");
        for (Map.Entry<String, Map<String, Integer>> entry : output.entrySet()) {
            String word = entry.getKey();
            Map<String, Integer> fileCounts = entry.getValue();
            for (Map.Entry<String, Integer> fc : fileCounts.entrySet()) {
                pw.println(word + "," + fc.getKey() + "," + fc.getValue());
            }
        }
    } catch (IOException e) {
        e.printStackTrace();
    }
}

/**

```


** Writes a timing summary to a CSV file with headers "Phase,Time(ms)".*

**/*

```
private static void writeTimingSummaryCSV(String csvFile, long fileReadTime, long
approach1Time,
                                     long approach2Time, long mapTime, long groupTime,
                                     long reduceTime, long approach3TotalTime) {
try (PrintWriter pw = new PrintWriter(new File(csvFile))) {
    pw.println("Phase,Time(ms)");
    pw.println("File Reading," + fileReadTime);
    pw.println("Approach 1 (Brute Force)," + approach1Time);
    pw.println("Approach 2 (Basic MapReduce)," + approach2Time);
    pw.println("Approach 3 Map Phase," + mapTime);
    pw.println("Approach 3 Group Phase," + groupTime);
    pw.println("Approach 3 Reduce Phase," + reduceTime);
    pw.println("Approach 3 Total," + approach3TotalTime);
} catch (IOException e) {
    e.printStackTrace();
}
}
```