

Assignment 2: MapReduce

1 Set-Up

To obtain large text files to test the program with, I downloaded the first 10 long books I could think of from `archive.org` in `txt` file form. These were:

1. The Bible;
2. *War & Peace* by Leo Tolstoy;
3. Plutarch's *Lives*;
4. Herodotus' *Histories*;
5. *City of God* by Augustine of Hippo;
6. *Faust* by Goethe;
7. *Wealth of Nations* by Adam Smith;
8. *Capital* by Karl Marx;
9. The complete works of William Shakespeare;
10. *Structure & Interpretation of Computer Programs* by Harold Abelson & Gerald Jay Sussman.

2 Baseline Results

I modified the code to measure & output the time taken by each approach, in milliseconds. I also added timing for the different phases of the two MapReduce implementations, timing the map time, group time, and reduce time separately.

```
[andrew@arch] ~/currsem/CT4I4/assignments/assignment2/code P [master] +
% javac *.java && java MapReduceFiles ../data/augustine city of god.txt ../data/Faust-Goethe.txt ../data/herodotus histories.txt ../data/holy bible.txt ../data/marx capital.txt ../data/plutarch lives.txt ../da
ta/shakespeare complete works.txt ../data/structure and interpretation of computer programs.txt ../data/war and peace.txt ../data/wealth of nations.txt
Brute Force Results:
  Total Time: 1295

MapReduce Results:
  Map Time: 1387
  Group Time: 682
  Reduce Time: 414
  Total Time: 2483

Distributed MapReduce Results:
  Map Time: 726
  Group Time: 726
  Reduce Time: 31813
  Total Time: 31813

[andrew@arch] ~/currsem/CT4I4/assignments/assignment2/code P [master] +
% [✓ 38s]
```

Figure 1: Baseline results for my list of files (in milliseconds)

As can be seen from the above terminal screenshot, the brute force approach performed best with no modifications, followed by the non-distributed MapReduce, followed by the distributed MapReduce; this is to be expected, as the brute force approach is the simplest & requires the fewest iterations over the data and no complex data structures. The non-distributed MapReduce requires more intermediate data structure and more iterations over the data. Finally, the non-optimised version of the distributed MapReduce is the slowest because it spawns a thread for each word in the dataset, causing massive stress on the CPU and memory.

I also updated the code to use `ArrayLists` rather than `LinkedLists` to reduce memory overhead and have faster traversal.

```
[andrew@arch] ~/currsem/CT4I4/assignments/assignment2/code P [master] +
% javac *.java && java MapReduceFiles ../data/augustine city of god.txt ../data/Faust-Goethe.txt ../data/herodotus histories.txt ../data/holy bible.txt ../data/marx capital.txt ../data/plutarch lives.txt ../da
ta/shakespeare complete works.txt ../data/structure and interpretation of computer programs.txt ../data/war and peace.txt ../data/wealth of nations.txt
Brute Force Results:
  Total Time: 1387

MapReduce Results:
  Map Time: 623
  Group Time: 547
  Reduce Time: 282
  Total Time: 1372

Distributed MapReduce Results:
  Map Time: 374
  Group Time: 388
  Reduce Time: 31874
  Total Time: 31874

[andrew@arch] ~/currsem/CT4I4/assignments/assignment2/code P [master] +
% [✓ 37s]
```

Figure 2: Baseline results with `ArrayList` update (in milliseconds)

As can be seen from the above terminal screenshot, this has no affect on the brute force results (besides slight variance due to background processes running on my laptop) as this approach did not use LinkedLists anyway. The non-distributed MapReduce approach was significantly faster due to the faster iteration and lower memory overhead. The distributed MapReduce saw significant improvements in the map & group phases, but these were dwarfed by the still greatly inefficient reduce phase.

3 Testing the Updated Code

After implementing the requested changes in steps 2–6 of the assignment specification, I then implemented a grid-search function which tested a range of values for the number of lines of text per map thread and the number of words per reduce thread. The results of this grid-search were exported to a CSV file for analysis. I then wrote a Python script to visualise the parameter combinations using heatmaps.

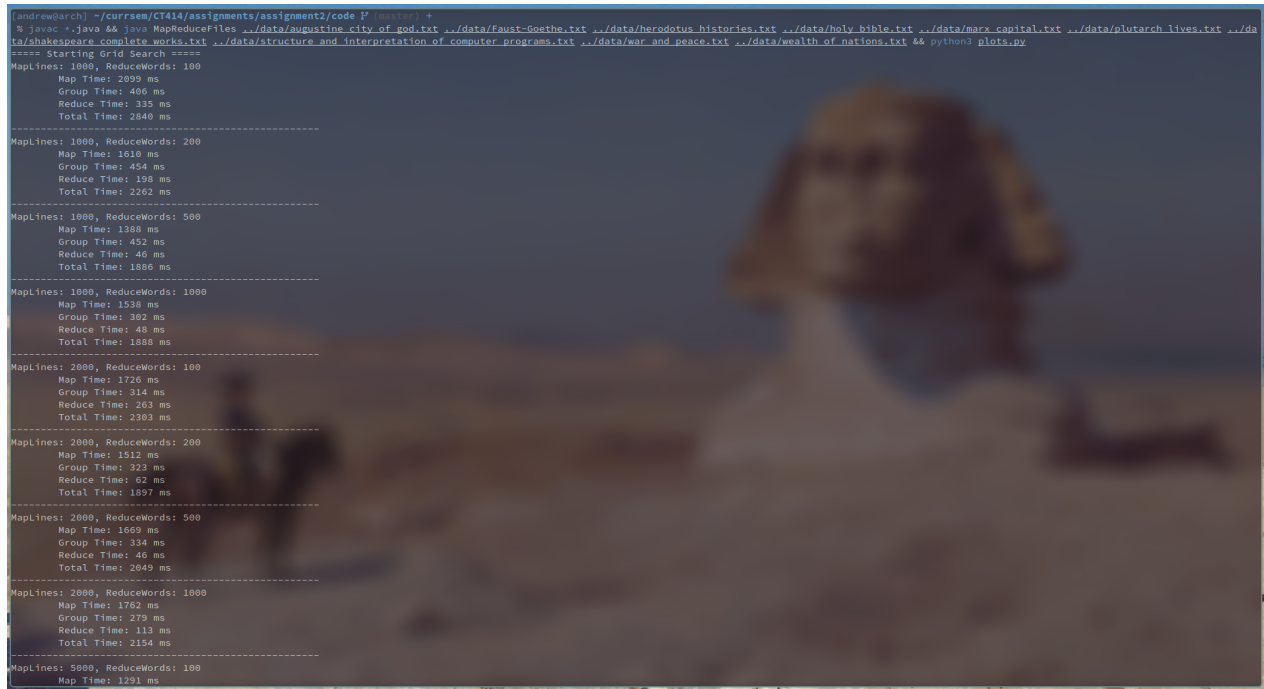


Figure 3: Running the grid-search and plotting the results

Map Lines	Reduce Words	Map Time (ms)	Group Time (ms)	Reduce Time (ms)	Total Time (ms)
1000	100	2099	406	335	2840
1000	200	1610	454	198	2262
1000	500	1388	452	46	1886
1000	1000	1538	302	48	1888
2000	100	1726	314	263	2303
2000	200	1512	323	62	1897
2000	500	1669	334	46	2049
2000	1000	1762	279	113	2154
5000	100	1291	331	92	1714
5000	200	1877	368	67	2312
5000	500	1640	396	41	2077
5000	1000	1439	365	193	1997
10000	100	1285	359	94	1738
10000	200	1598	359	98	2055
10000	500	1489	314	68	1871
10000	1000	1460	332	47	1839

Table 1: Results written to performance_results.csv

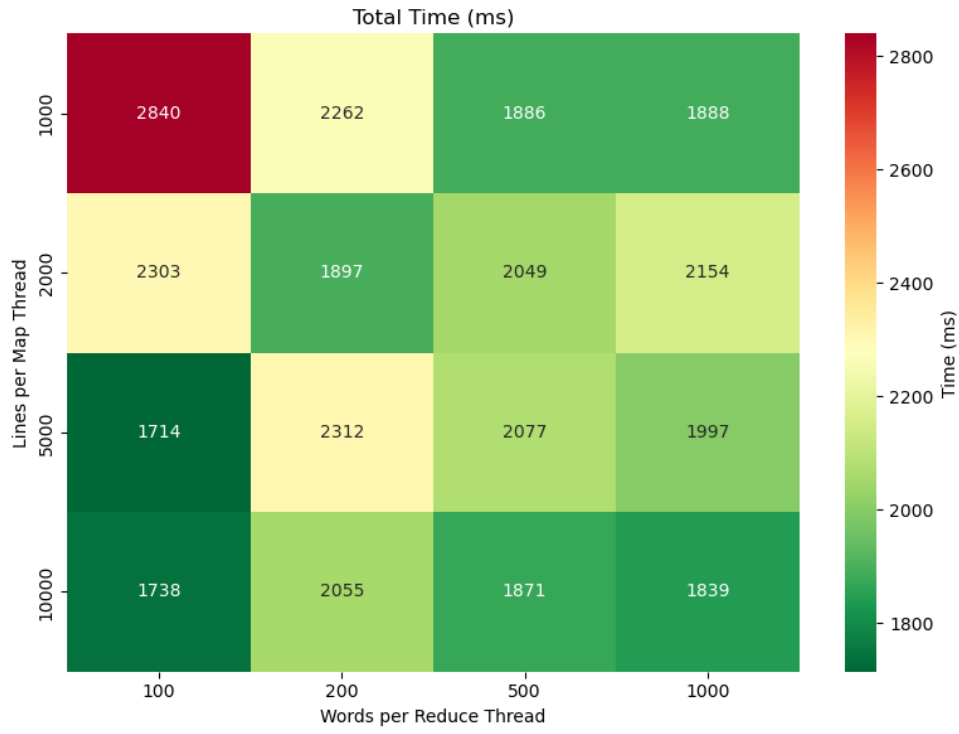


Figure 4: Heatmap of total time taken by each parameter combination

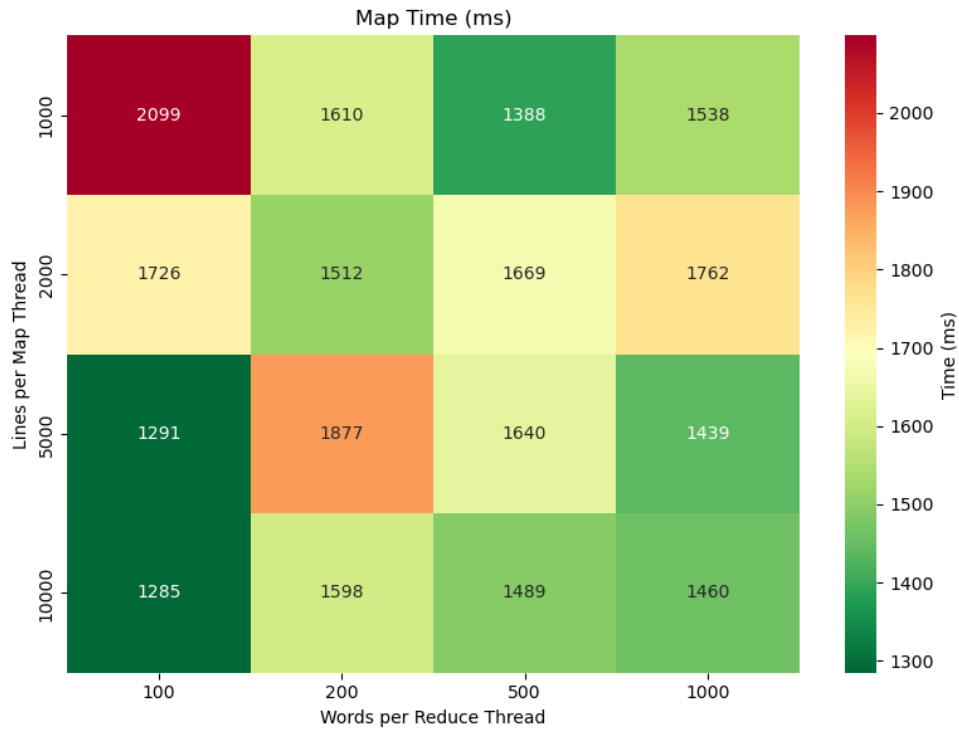


Figure 5: Heatmap of time taken during the map phase by each parameter combination

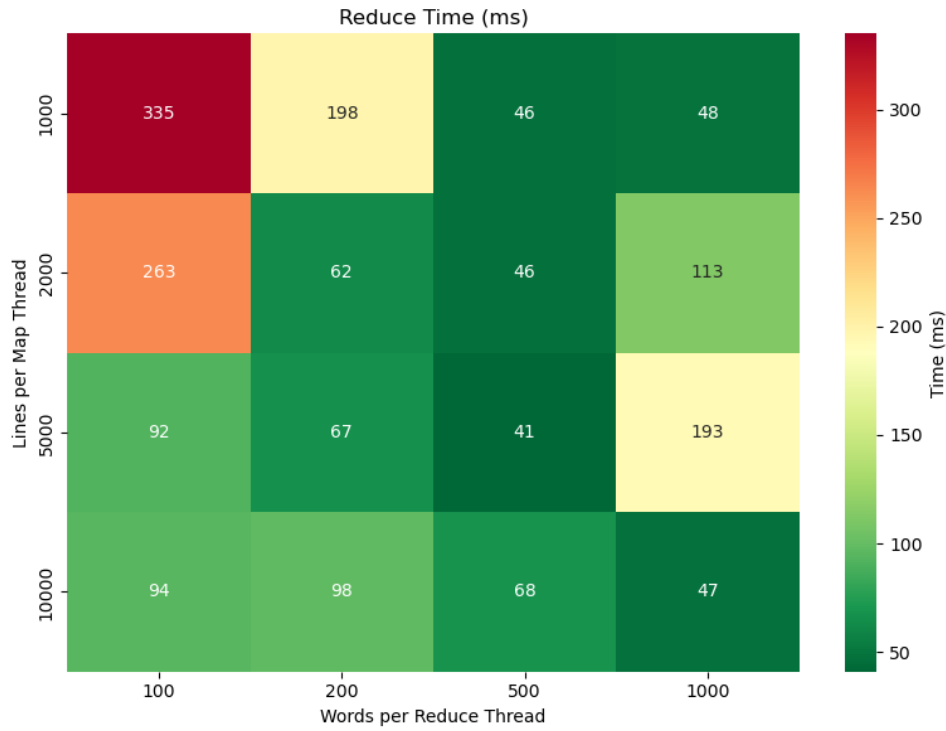


Figure 6: Heatmap of time taken during the reduce phase by each parameter combination

4 Appendix: Source Code

```

1 import java.util.*;
2 import java.io.*;
3
4 public class MapReduceFiles {
5
6     private static final String CSV_FILE = "performance_results.csv";
7
8     public static void main(String[] args) {
9         if (args.length < 1) {
10             System.err.println("Usage: java MapReduceFiles file1.txt file2.txt ... fileN.txt");
11             return;
12         }
13
14         Map<String, String> input = new HashMap<>();
15         try {
16             for (String filename : args) {
17                 input.put(filename, readFile(filename));
18             }
19         } catch (IOException ex) {
20             System.err.println("Error reading files: " + ex.getMessage());
21             ex.printStackTrace();
22             return;
23         }
24
25         int[] mapSizes = {1000, 2000, 5000, 10000};
26         int[] reduceSizes = {100, 200, 500, 1000};
27
28         System.out.println("==== Starting Grid Search =====");
29
30         try (PrintWriter writer = new PrintWriter(new FileWriter(CSV_FILE))) {
31             writer.println("MapLines,ReduceWords,MapTime,GroupTime,ReduceTime,TotalTime");
32         }

```

```

33     for (int mapSize : mapSizes) {
34         for (int reduceSize : reduceSizes) {
35             runDistributedMapReduce(input, mapSize, reduceSize, writer);
36         }
37     }
38
39 } catch (IOException e) {
40     System.err.println("Error writing to CSV file: " + e.getMessage());
41 }
42
43 System.out.println("==== Grid Search Complete =====");
44 System.out.println("Results saved to: " + CSV_FILE);
45 }
46
47 public static void runDistributedMapReduce(Map<String, String> input, int linesPerMapThread, int
↵ wordsPerReduceThread, PrintWriter csvWriter) {
48     final Map<String, Map<String, Integer>> output = new HashMap<>();
49
50     // MAP Phase
51     long mapStartTime = System.currentTimeMillis();
52     List<MappedItem> mappedItems = Collections.synchronizedList(new ArrayList<>());
53
54     final MapCallback<String, MappedItem> mapCallback = new MapCallback<>() {
55         public synchronized void mapDone(String file, List<MappedItem> results) {
56             mappedItems.addAll(results);
57         }
58     };
59
60     List<Thread> mapCluster = new ArrayList<>();
61     for (Map.Entry<String, String> entry : input.entrySet()) {
62         final String file = entry.getKey();
63         final String[] lines = entry.getValue().split("\\r?\\n");
64
65         for (int i = 0; i < lines.length; i += linesPerMapThread) {
66             int end = Math.min(i + linesPerMapThread, lines.length);
67             final List<String> chunk = new ArrayList<>();
68             for (int j = i; j < end; j++) {
69                 chunk.addAll(splitLongLine(lines[j]));
70             }
71
72             Thread t = new Thread(() -> map(file, chunk, mapCallback));
73             mapCluster.add(t);
74             t.start();
75         }
76     }
77
78     for (Thread t : mapCluster) {
79         try {
80             t.join();
81         } catch (InterruptedException e) {
82             throw new RuntimeException(e);
83         }
84     }
85
86     long mapTotalTime = System.currentTimeMillis() - mapStartTime;
87
88     // GROUP Phase
89     long groupStartTime = System.currentTimeMillis();
90     Map<String, List<String>> groupedItems = new HashMap<>();
91     for (MappedItem item : mappedItems) {
92         groupedItems.computeIfAbsent(item.getWord(), k -> new ArrayList<>()).add(item.getFile());

```

```

93     }
94     long groupTotalTime = System.currentTimeMillis() - groupStartTime;
95
96     // REDUCE Phase
97     long reduceStartTime = System.currentTimeMillis();
98     final ReduceCallback<String, String, Integer> reduceCallback = (word, result) -> {
99         synchronized (output) {
100             output.put(word, result);
101         }
102     };
103
104     List<Thread> reduceCluster = new ArrayList<>();
105     List<Map<String, List<String>>> reduceChunks = new ArrayList<>();
106     Map<String, List<String>> currentChunk = new HashMap<>();
107     int count = 0;
108
109     for (Map.Entry<String, List<String>> entry : groupedItems.entrySet()) {
110         currentChunk.put(entry.getKey(), entry.getValue());
111         count++;
112         if (count >= wordsPerReduceThread) {
113             reduceChunks.add(currentChunk);
114             currentChunk = new HashMap<>();
115             count = 0;
116         }
117     }
118     if (!currentChunk.isEmpty()) reduceChunks.add(currentChunk);
119
120     for (final Map<String, List<String>> chunk : reduceChunks) {
121         Thread t = new Thread(() -> {
122             for (Map.Entry<String, List<String>> entry : chunk.entrySet()) {
123                 reduce(entry.getKey(), entry.getValue(), reduceCallback);
124             }
125         });
126         reduceCluster.add(t);
127         t.start();
128     }
129
130     for (Thread t : reduceCluster) {
131         try {
132             t.join();
133         } catch (InterruptedException e) {
134             throw new RuntimeException(e);
135         }
136     }
137
138     long reduceTotalTime = System.currentTimeMillis() - reduceStartTime;
139     long totalTime = mapTotalTime + groupTotalTime + reduceTotalTime;
140
141     // Print & Log
142     System.out.println("MapLines: " + linesPerMapThread + ", ReduceWords: " + wordsPerReduceThread);
143     System.out.println("\tMap Time: " + mapTotalTime + " ms");
144     System.out.println("\tGroup Time: " + groupTotalTime + " ms");
145     System.out.println("\tReduce Time: " + reduceTotalTime + " ms");
146     System.out.println("\tTotal Time: " + totalTime + " ms");
147     System.out.println("-----");
148
149     csvWriter.printf("%d,%d,%d,%d,%d,%d\n",
150         linesPerMapThread, wordsPerReduceThread,
151         mapTotalTime, groupTotalTime, reduceTotalTime, totalTime);
152     csvWriter.flush();
153 }

```

```

154
155 public static void map(String file, List<String> lines, MapCallback<String, MappedItem> callback) {
156     List<MappedItem> results = new ArrayList<>();
157     for (String line : lines) {
158         String[] words = line.trim().split("\\s+");
159         for (String word : words) {
160             word = word.replaceAll("[^a-zA-Z]", "").toLowerCase();
161             if (!word.isEmpty()) {
162                 results.add(new MappedItem(word, file));
163             }
164         }
165     }
166     callback.mapDone(file, results);
167 }
168
169 public static void reduce(String word, List<String> list, ReduceCallback<String, String, Integer>
↵ callback) {
170     Map<String, Integer> reducedList = new HashMap<>();
171     for (String file : list) {
172         reducedList.put(file, reducedList.getOrDefault(file, 0) + 1);
173     }
174     callback.reduceDone(word, reducedList);
175 }
176
177 public interface MapCallback<E, V> {
178     void mapDone(E key, List<V> values);
179 }
180
181 public interface ReduceCallback<E, K, V> {
182     void reduceDone(E e, Map<K, V> results);
183 }
184
185 private static class MappedItem {
186     private final String word;
187     private final String file;
188
189     public MappedItem(String word, String file) {
190         this.word = word;
191         this.file = file;
192     }
193
194     public String getWord() {
195         return word;
196     }
197
198     public String getFile() {
199         return file;
200     }
201
202     @Override
203     public String toString() {
204         return "[" + word + ", " + file + "]";
205     }
206 }
207
208 private static String readFile(String pathname) throws IOException {
209     File file = new File(pathname);
210     StringBuilder fileContents = new StringBuilder((int) file.length());
211     Scanner scanner = new Scanner(new BufferedReader(new FileReader(file)));
212     String lineSeparator = System.getProperty("line.separator");
213

```

```

214     try {
215         while (scanner.hasNextLine()) {
216             fileContents.append(scanner.nextLine()).append(lineSeparator);
217         }
218         return fileContents.toString();
219     } finally {
220         scanner.close();
221     }
222 }
223
224 private static List<String> splitLongLine(String line) {
225     List<String> result = new ArrayList<>();
226     while (line.length() > 80) {
227         int splitAt = line.lastIndexOf(' ', 80);
228         if (splitAt <= 0) splitAt = 80;
229         result.add(line.substring(0, splitAt));
230         line = line.substring(splitAt).trim();
231     }
232     if (!line.isEmpty()) result.add(line);
233     return result;
234 }
235 }

```

Listing 1: MapReduceFiles.java

```

1  import pandas as pd
2  import seaborn as sns
3  import matplotlib.pyplot as plt
4
5  df = pd.read_csv('performance_results.csv')
6
7  def save_heatmap(metric, title, filename):
8      pivot = df.pivot(index='MapLines', columns='ReduceWords', values=metric)
9      plt.figure(figsize=(8, 6))
10     sns.heatmap(
11         pivot,
12         annot=True,
13         fmt="d",
14         cmap="RdYlGn_r",
15         cbar_kws={'label': 'Time (ms)'}
16     )
17     plt.title(title)
18     plt.ylabel("Lines per Map Thread")
19     plt.xlabel("Words per Reduce Thread")
20     plt.tight_layout()
21     plt.savefig(filename)
22     plt.close()
23     print(f"Saved: {filename}")
24
25 save_heatmap('TotalTime', 'Total Time (ms)', '../latex/images/total_time_heatmap.png')
26 save_heatmap('MapTime', 'Map Time (ms)', '../latex/images/map_time_heatmap.png')
27 save_heatmap('ReduceTime', 'Reduce Time (ms)', '../latex/images/reduce_time_heatmap.png')

```

Listing 2: plots.py