

Multithreaded Word Count Using Mapper–Reducer Style

Hoang Minh Quan
Student ID: 23BI14371

December 5, 2025

Abstract

This report presents a small C++ application that counts word frequencies in a text file using a mapper–reducer style design. The program reads an input file, splits the text into chunks that are processed in parallel by mapper threads, and then combines the intermediate results with reducer threads before writing the global word counts to an output file.

Contents

1	Introduction	2
2	Problem Description	2
3	System Overview	2
3.1	High-level Design	2
3.2	Illustration	3
4	Implementation Details	3
4.1	Word Normalisation and Splitting	3
4.2	Mapper Threads	3
4.3	Reducer Threads	3
4.4	Source Code Listing	3
5	Example Run	7
6	Conclusion	7

1 Introduction

Counting how often each word appears in a text is a simple but representative problem for data-processing frameworks. In this assignment, we implement a standalone C++ program that performs word counting in parallel on a single machine. The program is inspired by the MapReduce model: we use several *mapper* threads to create local word-count maps and several *reducer* threads to merge these partial results into a final map.

The implementation uses the C++ standard library, threads, and `std::unordered_map` to store word frequencies. Input and output are handled via standard text files.

2 Problem Description

The goal is to build a command-line tool with the following behaviour:

- read a text file given by the user;
- normalise words (convert to lowercase and remove non-alphabetic characters);
- count how many times each distinct word occurs;
- allow the user to configure how many mapper and reducer threads are used;
- store the word counts in a separate output file.

The provided example input file contains the sentence

```
why do programmers prefer dark mode because light attracts bugs
```

and the expected output lists each word once together with its frequency, for example:

```
attracts 1
because 1
bugs 1
dark 1
do 1
light 1
mode 1
prefer 1
programmers 1
why 1
```

3 System Overview

3.1 High-level Design

The program follows four main stages:

1. **Load file:** the entire input file is read into a single string.
2. **Map phase:** the text is partitioned into chunks, each processed by a mapper thread producing a local (`word`, `count`) map.
3. **Shuffle phase:** intermediate pairs are grouped so that all counts of the same word are sent to the same reducer.
4. **Reduce phase:** reducer threads accumulate final word counts and the result is written to the output file.

3.2 Illustration

Figure 1 illustrates the architecture of the application. A single input text is transformed into several mapper chunks, aggregated via reducer threads, and finally written as a sorted list of words and frequencies.

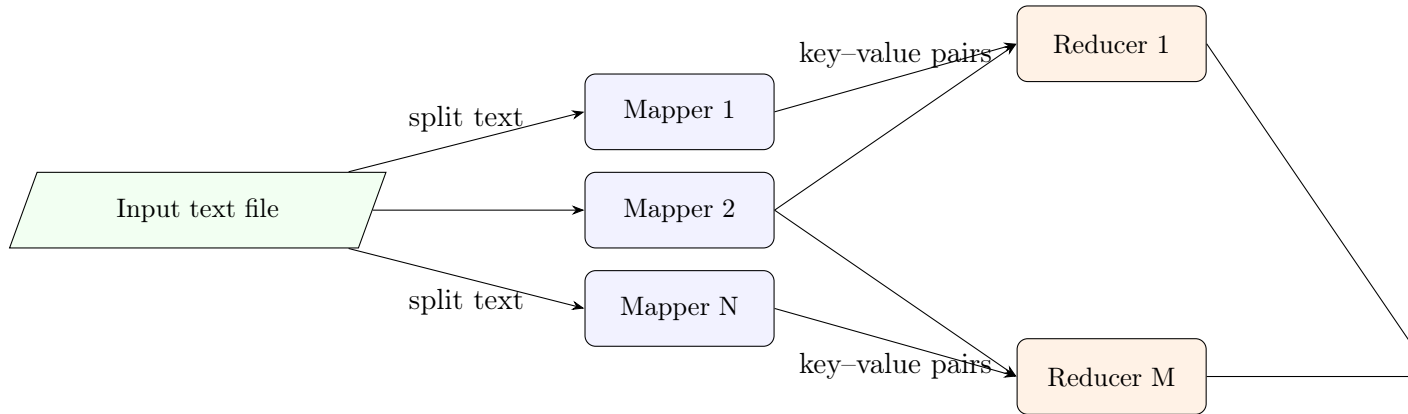


Figure 1: Mapper-reducer style architecture for the word-count tool.

4 Implementation Details

4.1 Word Normalisation and Splitting

The function `normalize` converts characters to lowercase and removes anything that is not an alphabetic letter. A second helper function, `split_words`, walks through the input text, extracts sequences of non-whitespace characters, normalises them, and returns a vector of clean words.

4.2 Mapper Threads

Each mapper thread receives a substring (“chunk”) of the original text, calls `split_words` on that chunk, then builds a local `unordered_map<std::string,int>` that counts occurrences of each word in the chunk.

To avoid cutting words in half, the code adjusts the chunk boundary to the next whitespace character before launching the mapper.

4.3 Reducer Threads

After all mapper threads have finished, the program distributes the intermediate `(word, count)` pairs to reducers. A simple hash partitioning scheme is used:

$$\text{reducer_id} = \text{hash}(\text{word}) \bmod \text{num_reducers}. \quad (1)$$

Each reducer accumulates the counts it receives into its own `unordered_map`. Finally, these per-reducer maps are merged into a single global map.

4.4 Source Code Listing

Listing 1 shows the full implementation of the program.

```
1 #include <algorithm>
2 #include <cctype>
3 #include <fstream>
4 #include <iostream>
```

```

5  #include <mutex>
6  #include <string>
7  #include <thread>
8  #include <unordered_map>
9  #include <utility>
10 #include <vector>
11
12 using MapperOutput = std::unordered_map<std::string, int>;
13 using ReducerOutput = std::unordered_map<std::string, int>;
14
15 std::mutex io_mutex;
16
17 std::string normalize(const std::string &w) {
18     std::string res;
19     for (char c : w) {
20         if (std::isalpha(static_cast<unsigned char>(c))) {
21             res.push_back(static_cast<char>(std::tolower(c)));
22         }
23     }
24     return res;
25 }
26
27 std::vector<std::string> split_words(const std::string &text) {
28     std::vector<std::string> words;
29     std::string current;
30     for (char c : text) {
31         if (std::isspace(static_cast<unsigned char>(c))) {
32             if (!current.empty()) {
33                 std::string norm = normalize(current);
34                 if (!norm.empty()) words.push_back(norm);
35                 current.clear();
36             }
37             } else {
38                 current.push_back(c);
39             }
40         }
41         if (!current.empty()) {
42             std::string norm = normalize(current);
43             if (!norm.empty()) words.push_back(norm);
44         }
45         return words;
46     }
47
48 void mapper_worker(const std::string &chunk, MapperOutput &out_map) {
49     auto words = split_words(chunk);
50     for (const auto &w : words) {
51         ++out_map[w];
52     }
53 }
54
55 void reducer_worker(const std::vector<std::pair<std::string, int>> &
56     pairs,
57     ReducerOutput &out_map) {
58     for (const auto &p : pairs) {
59         out_map[p.first] += p.second;
60     }
61 }

```

```

62 int main(int argc, char **argv) {
63     if (argc < 3) {
64         std::cerr << "Usage: " << argv[0]
65             << " <input_file> <output_file> [num_mappers] [
                num_reducers]\n";
66         return 1;
67     }
68
69     std::string input_file = argv[1];
70     std::string output_file = argv[2];
71     int num_mappers = (argc >= 4) ? std::stoi(argv[3]) : 4;
72     int num_reducers = (argc >= 5) ? std::stoi(argv[4]) : 2;
73
74     if (num_mappers <= 0 || num_reducers <= 0) {
75         std::cerr << "Number of mappers and reducers must be > 0\n";
76         return 1;
77     }
78
79     std::ifstream in(input_file, std::ios::binary);
80     if (!in.is_open()) {
81         std::cerr << "Cannot open input file: " << input_file << "\n";
82         return 1;
83     }
84     std::string text((std::istreambuf_iterator<char>(in)),
85                     std::istreambuf_iterator<char>());
86     in.close();
87
88     if (text.empty()) {
89         std::cerr << "Input file is empty.\n";
90         return 1;
91     }
92
93     std::vector<std::thread> mapper_threads;
94     std::vector<MapperOutput> mapper_outputs(num_mappers);
95
96     size_t chunk_size = text.size() / num_mappers;
97     size_t start = 0;
98
99     for (int i = 0; i < num_mappers; ++i) {
100         size_t end = (i == num_mappers - 1) ? text.size() : start +
            chunk_size;
101
102         if (end < text.size()) {
103             while (end < text.size() &&
104                 !std::isspace(static_cast<unsigned char>(text[end])))
105                 ++end;
106             }
107         }
108
109         std::string chunk = text.substr(start, end - start);
110         mapper_threads.emplace_back(mapper_worker, chunk,
111                                     std::ref(mapper_outputs[i]));
112         start = end;
113     }
114
115     for (auto &t : mapper_threads) t.join();
116

```

```

117     std::vector<std::vector<std::pair<std::string, int>>> reducer_inputs
118         (
119             num_reducers);
120
121     for (const auto &m_out : mapper_outputs) {
122         for (const auto &kv : m_out) {
123             const std::string &word = kv.first;
124             int count = kv.second;
125             std::size_t h = std::hash<std::string>{}(word);
126             int rid = static_cast<int>(h % num_reducers);
127             reducer_inputs[rid].emplace_back(word, count);
128         }
129     }
130
131     std::vector<std::thread> reducer_threads;
132     std::vector<ReducerOutput> reducer_outputs(num_reducers);
133
134     for (int r = 0; r < num_reducers; ++r) {
135         reducer_threads.emplace_back(reducer_worker,
136                                     std::cref(reducer_inputs[r]),
137                                     std::ref(reducer_outputs[r]));
138     }
139
140     for (auto &t : reducer_threads) t.join();
141
142     std::unordered_map<std::string, int> final_counts;
143     for (const auto &r_out : reducer_outputs) {
144         for (const auto &kv : r_out) {
145             final_counts[kv.first] += kv.second;
146         }
147     }
148
149     std::vector<std::pair<std::string, int>> sorted(final_counts.begin(),
150                                                    final_counts.end());
151     std::sort(sorted.begin(), sorted.end(),
152               [](const auto &a, const auto &b) { return a.first < b.first; });
153
154     std::ofstream out(output_file);
155     if (!out.is_open()) {
156         std::cerr << "Cannot open output file: " << output_file << "\n";
157         return 1;
158     }
159
160     for (const auto &kv : sorted) {
161         out << kv.first << " " << kv.second << "\n";
162     }
163     out.close();
164
165     std::cout << "Word count written to: " << output_file << "\n";
166     return 0;

```

Listing 1: C++ implementation of the multithreaded word-count tool

5 Example Run

With the example sentence given earlier, the program can be compiled and executed as:

`mpic++` is not required here; we simply use `g++` or `clang++`:

```
g++ -std=c++17 -pthread wordcount.cpp -o wordcount
./wordcount input.txt output.txt 4 2
```

The resulting `output.txt` contains one line per word and its frequency, as shown previously. For this short input, each word appears exactly once.

6 Conclusion

This small project shows how the mapper-reducer idea can be implemented using only the C++ standard library. Even though the example input is short, the design naturally scales to larger text files: more mapper and reducer threads can be enabled to use the available CPU cores. Possible extensions include:

- adding command-line options to filter stop words;
- supporting case-sensitive or case-insensitive counting modes;
- exporting the results in formats such as CSV or JSON for further analysis.