

# Advanced Parallel Programming

## Exercise 3

Marvin Kaster



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

Summer term 2025  
05.06.2025

Please solve the following tasks by 05.06.2025. The results are not graded, but a solution is discussed on 05.06.2025.

---

### Task 1: Calculate the sum of a vector

---

Your task is to calculate the sum of a large vector. First, generate a random vector. Then, calculate the sum in serial. Finally, use `std::thread` to calculate the sum in parallel. One thread is responsible to calculate the sum of a subsection of the vector. As you cannot return the value directly as a return value, store it in a second vector where thread `i` will store its result at position `i`. Then, add all partial sums to calculate the final overall sum. Make sure that your serial and parallel implementation calculate the same value and measure how long each implementation needs. Create a plot for your time measurements for different number of threads and sizes of vectors.

We added an utility function for your convenience.

Hints:

- The constructor of `std::thread` will create copies of each passed parameter. Take a look at `std::ref` and `std::cref`<sup>1</sup> to avoid creating unnecessary copies of vectors.
- Take a look at `std::chrono`<sup>2</sup> for time measurements
- Take a look at `std::reduce`<sup>3</sup> for calculating a sum of a vector sequentially

```
1 #pragma once
2
3 #include <random>
4 #include <vector>
5
6 template<std::integral T>
7 static std::vector<T> create_random_vector(const std::size_t size=1e5, const int seed = 42) {
8     std::mt19937 mersenne_engine(seed);
9     std::uniform_int_distribution<T> dist {1, 10};
10    const auto gen = [&]() {
11        return dist(mersenne_engine);
12    };
13
14    std::vector<T> random_vector(size);
15    std::generate(random_vector.begin(), random_vector.end(), gen);
16
17    return random_vector;
18 }
19
20 template<typename T>
21 static void print_vector(const std::vector<T> &vector) {
22     for (const auto &i: vector) {
23         std::cout << i << " ";
24     }
25     std::cout << "\n";
26 }
```

<sup>1</sup><https://en.cppreference.com/w/cpp/utility/functional/ref>

<sup>2</sup>[https://en.cppreference.com/w/cpp/chrono/steady\\_clock/now](https://en.cppreference.com/w/cpp/chrono/steady_clock/now)

<sup>3</sup><https://en.cppreference.com/w/cpp/algorithm/reduce>

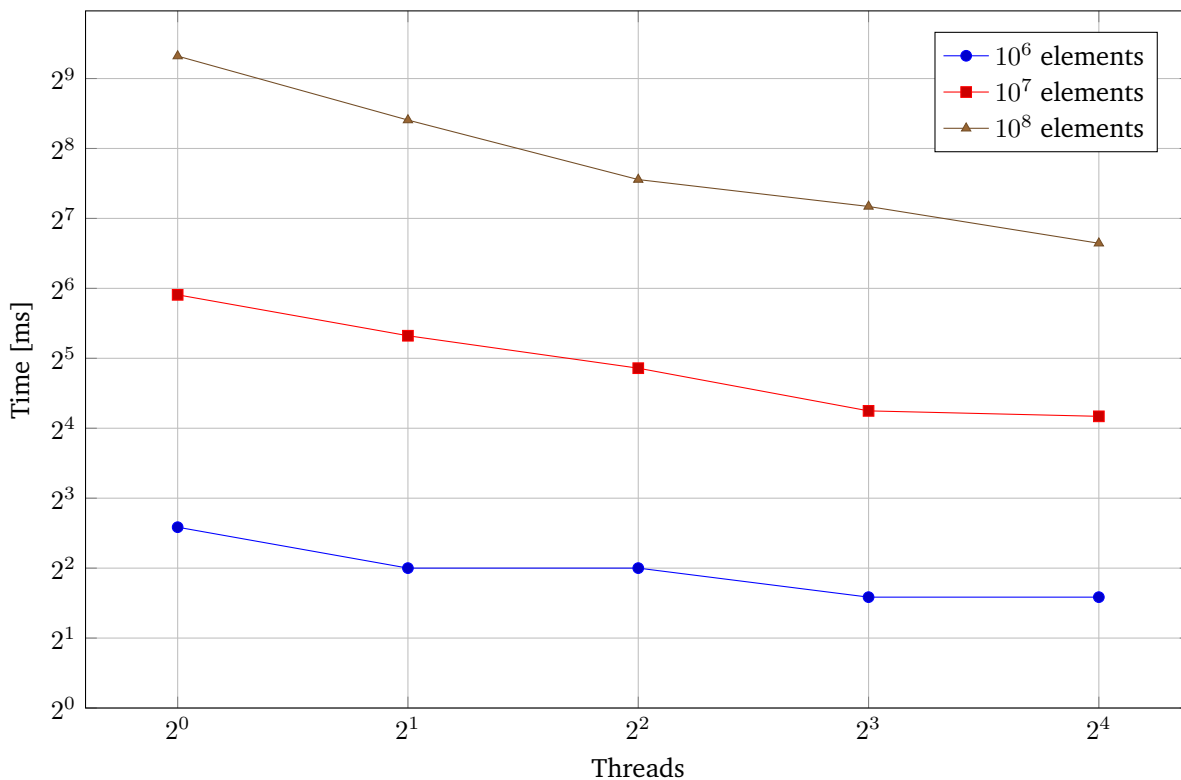


Figure 1: Execution time of the implementation (y-axis) vs number of threads used (x-axis) for different number of vector sizes.

26 }

#### Listing 1: Utility functions

Solution:

```

1  #include <algorithm>
2  #include <cassert>
3  #include <iostream>
4  #include <thread>
5  #include <vector>
6
7  #include "util.h"
8
9  using number_type = int;
10
11 void partial_sum(const std::vector<number_type> &input_vector, const std::size_t begin, const std::size_t
    end,
12                  number_type &result_output) {
13     assert(begin <= input_vector.size());
14     assert(end <= input_vector.size());
15     assert(begin <= end);
16     const auto sum = std::reduce(std::next(input_vector.begin(), begin), std::next(input_vector.begin(),
        end), 0,
17                                  std::plus<>());
18     result_output = sum;
19 }
20
21
22 int main() {
23     const std::size_t size = 1e7;

```

```

24  const std::size_t number_threads = 4U;
25  const auto vector = create_random_vector<number_type>(size);
26
27  //Serial implementation
28  const auto begin_serial = std::chrono::steady_clock::now();
29  const auto serial_sum = std::reduce(vector.begin(), vector.end(), 0, std::plus<>());
30  assert(size <= serial_sum && serial_sum <= size * 10U);
31  const auto end_serial = std::chrono::steady_clock::now();
32  const auto duration_serial = end_serial - begin_serial;
33  std::cout << "Serial result " << serial_sum << '\n';
34  std::cout << "Took " << std::chrono::duration_cast<std::chrono::milliseconds>(duration_serial).count()
    << "ms\n\n";
35
36
37  //Parallel implementation
38  const auto begin_parallel = std::chrono::steady_clock::now();
39  std::vector<number_type> results(number_threads);
40  std::vector<std::thread> threads;
41  threads.reserve(number_threads);
42  const auto number_numbers_per_thread = std::max(std::size_t(1U), size / number_threads);
43
44
45  for (auto thread_id = 0U; thread_id < number_threads; thread_id++) {
46      const std::size_t begin = std::min(thread_id * number_numbers_per_thread, size);
47      const std::size_t end = std::min((thread_id + 1) * number_numbers_per_thread, size);
48      auto thread = std::thread(partial_sum, std::cref(vector), begin, end, std::ref(results[thread_id])
49  );
50      threads.push_back(std::move(thread));
51  }
52
53  for (auto &thread: threads) {
54      thread.join();
55  }
56
57  const auto parallel_sum = std::reduce(results.begin(), results.end(), 0, std::plus<>());
58  const auto end_parallel = std::chrono::steady_clock::now();
59  const auto duration_parallel = end_parallel - begin_parallel;
60  std::cout << parallel_sum << '\n';
61  std::cout << "parallel result " << parallel_sum << '\n';
62  std::cout << "Took " << std::chrono::duration_cast<std::chrono::milliseconds>(duration_parallel).count()
    << "ms\n\n";
63  std::cout << "Used " << number_threads << " threads out of " << std::thread::hardware_concurrency() <<
    " cores\n\n";
64
65  return 0;
66 }

```

Listing 2: Vector addition

## Task 2: Implement a thread safe counter map

When we count the occurrences of unique objects in a list, we need some associative container that keeps track of the previous occurrences for each unique object. This could be implemented with a `std::unordered_map`. However, we would like to be able to go through large lists in parallel or through multiple lists in parallel. The implementation of `std::unordered_map` is not thread-safe, and therefore, we cannot share an `std::unordered_map` object with multiple threads. Your task is to implement a thread-safe counter. Internally, the data should be stored using a hash list to enable fast access to the elements.

The idea is to use the key's hash value to calculate the index in a fixed-size list. Note that this can result in the same index for different keys. Hence, we need a second list with a variable size for each index in the hash list, over which we have to iterate to find the entry with the matching key. Figure 2 shows a visualization of the data structure. Please implement the interface shown below.

```

1  template<typename KeyType, typename ValueType, typename HashType=std::hash<KeyType>>
2  class ThreadSafeCounterDict {
3  public:

```

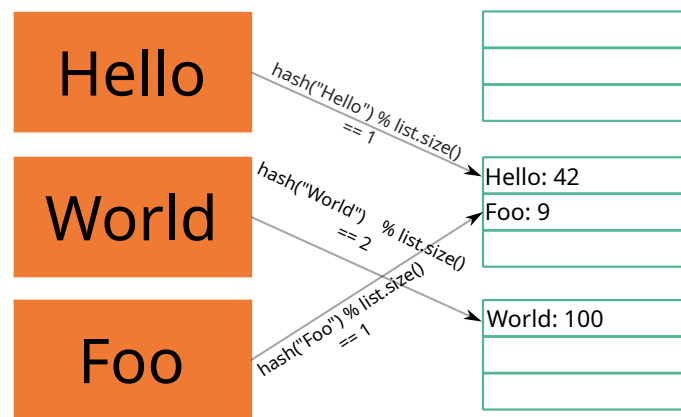


Figure 2: Visualization of a hash map. We calculate an index in our list by calculate the hash value of the key and apply a modulo operation. At this index we have a nested list to which we can add the value. If two keys point to the same index, we need to check all elements in that nested list.

```

4 CounterDict(const ValueType &default_value, const size_t list_size);
5
6 ValueType get_value_for(const KeyType &key) const;
7
8 ValueType increase_value_for_key(const KeyType &key, const ValueType &increment);
9
10 std::vector<std::pair<KeyType, ValueType>> convert_to_pairs();
11
12 };

```

### Solution:

Our counter dictionary interface has only three methods:

- Read the value of a key
- Increase the value of a key
- Return all key-value pairs

We follow simple thread-safety rules in our implementation. We avoid returning and storing references of objects to ensure modification from outside the class. As the bucket list has a fixed size and is not modified after the initialization in the constructor, we need no mutex for it. It is sufficient to use a single mutex per bucket. This mutex can use a `std::shared_lock` for read operations and a `std::unique_lock` for writing operations. The mutex must be locked at the beginning of the member function of the bucket before any reading or writing of the key-value pair list. It will be unlocked automatically after leaving the function.

If two threads want to add a value for the same key, a race condition could appear if it checks for the existence of the key first and then adds it without holding a lock the entire time. We avoid this problem by only providing a method to increase the value and provide a default value that is used when the key is not in the dictionary yet. In that way, we can hold a lock the entire time. Returning a default value in our read method solves the same problem in that we do not have to check if the key is in it.

Exceptions can only be raised in the `increase_value_for_key` function. The method `std::vector::emplace_back` could throw an exception (e.g., because of low memory), but it is exception-safe, leaving the vector in a valid state. Otherwise, exceptions can only be thrown in the operators from `ValueType` (e.g., the assignment operator, addition operator). It is up to the user to ensure thread safety for the provided type.

The `convert_to_pairs` method locks multiple mutexes and is the only point where a deadlock could occur. Hence, we must always lock the mutexes in the same order. We do it by locking the mutexes of the buckets with increasing indices.

```

1  #pragma once
2
3  #include <functional>
4  #include <memory>
5  #include <mutex>
6  #include <shared_mutex>
7  #include <tuple>
8  #include <vector>
9
10 template<typename KeyType, typename ValueType, typename HashType=std::hash<KeyType>>
11 class ThreadSafeCounterDict {
12     private:
13         /**
14          * Class that holds items for a single bucket
15          */
16         class Bucket {
17             //The actual data in the bucket
18             mutable std::vector<std::pair<KeyType, ValueType>> item_data;
19
20             /**
21              * Finds the entry for a certain key within this bucket
22              * @param key The key
23              * @return Iterator that point to the key-value pair or to the end of the vector when the key is not in
24              *         the bucket
25              */
26             std::vector<std::pair<KeyType, ValueType>>::iterator find_entry_for_key(const KeyType &key) const {
27                 return std::find_if(item_data.begin(), item_data.end(),
28                                     [&](const auto &item) { return item.first == key; });
29             }
30
31             public:
32             //The mutex for the bucket
33             mutable std::shared_mutex item_mutex;
34
35             /**
36              * Returns the current value for the key in this bucket
37              * @param key The key of the key-value pair
38              * @return The value belonging to the key
39              */
40             [[nodiscard]] ValueType get_value_for_key(const KeyType &key, const ValueType &initial_value) const {
41                 //We need a shared lock on this bucket while we read the value so that no thread can modify it
42                 //Other threads reading the same bucket is ok
43                 std::shared_lock<std::shared_mutex> lock(item_mutex);
44
45                 //Find the value belonging to the key
46                 const auto &it = find_entry_for_key(key);
47                 if (it != item_data.end()) {
48                     //Return a copy of the value
49                     return it->second;
50                 }
51                 //No value is assigned yet return the initial value
52                 return initial_value;
53             }
54
55             /**
56              * Increase the value belonging to the key in this bucket
57              * @param key The key of the key-value pair
58              * @param increment The number of which we increase the value
59              * @param initial_value The initial value belonging to the key when there is no value assigned to the
60              *         key yet
61              * @return The new value belonging to the key
62              */
63             ValueType
64             increase_value_for_key(const KeyType &key, const ValueType &increment, const ValueType &initial_value)
65             {
66                 //We need an unique lock on this bucket while we modify it
67                 //Other threads are not allowed to read the bucket while we modify it
68                 std::unique_lock<std::shared_mutex> lock(item_mutex);
69
70                 //Find the value belonging to the key

```

```

68     auto it = find_entry_for_key(key);
69     if (it != item_data.end()) {
70         //Read the old value
71         const auto old_value = it->second;
72         //Add the increment
73         const auto new_value = old_value + increment;
74         //Store the new value
75         it->second = new_value;
76         //Return a copy of the new value
77         return new_value;
78     }
79
80     //No value is assigned yet
81     //Add the increment to the initial value
82     const auto new_value = initial_value + increment;
83     item_data.emplace_back(key, new_value);
84     return new_value;
85 }
86
87 /**
88  * Returns a list of all stored key-value pairs in this bucket
89  * @return List of all key-value pairs
90  */
91 [[nodiscard]] std::vector<std::pair<KeyType, ValueType>> get_item_data() const {
92     //We need a shared lock while we copy the data to prevent other threads from modifying it meanwhile
93     std::shared_lock<std::shared_mutex> lock(item_mutex);
94
95     return item_data;
96 }
97 };
98
99 public:
100 /**
101  * Constructor
102  * @param init_value The initial object value for each key
103  * @param list_size The fixed size of the hash list
104  */
105 ThreadSafeCounterDict(const ValueType &initial_value, const size_t list_size) : initial_value(
106     initial_value),
107     buckets(list_size) {
108     for (auto i = 0; i < list_size; i++) {
109         buckets[i] = std::make_unique<Bucket>();
110     }
111 }
112
113 ThreadSafeCounterDict() = delete;
114
115 ThreadSafeCounterDict(ThreadSafeCounterDict &other) = delete;
116
117 ThreadSafeCounterDict &operator=(const ThreadSafeCounterDict &other) = delete;
118
119 /**
120  * Returns the current value for the key
121  * @param key The key of the key-value pair
122  * @return The value belonging to the key
123  */
124 [[nodiscard]] ValueType get_value_for_key(const KeyType &key) const {
125     //Find the bucket belonging to the key
126     const auto &bucket = get_bucket(key);
127
128     //Return the value of the bucket. The bucket handles the mutex itself
129     return bucket.get_value_for(key, initial_value);
130 }
131
132 /**
133  * Increase the value belonging to the key
134  * @param key The key of the key-value pair
135  * @param increment The number of which we increase the value
136  * @return The new value belonging to the key
137  */

```

```

137 ValueType increase_value_for_key(const KeyType &key, const ValueType &increment) {
138     //Find the bucket belonging to the key
139     auto &item_holder = get_bucket(key);
140
141     //Increase the value by the increment and return it. The bucket handles the mutex itself
142     return item_holder.increase_value_for_key(key, increment, initial_value);
143 }
144
145 /**
146  * Returns a list of all stored key-value pairs
147  * @return List of all key-value pairs
148  */
149 std::vector<std::pair<KeyType, ValueType>> convert_to_pairs() {
150     //We need to lock all buckets at the beginning as we want to iterate over all of them
151     std::vector<std::shared_lock<std::shared_mutex>> locks{};
152     for (const auto &bucket: buckets) {
153         locks.emplace_back(bucket->item_mutex);
154     }
155
156     //Copy the data from every bucket
157     std::vector<std::pair<KeyType, ValueType>> all_pairs{};
158     for (auto i = 0; i < buckets.size(); i++) {
159         const auto &item_data = buckets[i]->get_item_data();
160         //Copy each element of item data (key-value pairs) into the new vector
161         all_pairs.insert(all_pairs.end(), item_data.begin(), item_data.end());
162     }
163     return all_pairs;
164 }
165
166 private:
167
168     std::vector<std::unique_ptr<Bucket>> buckets{};
169     HashType hash_function{};
170     ValueType initial_value;
171
172     Bucket &get_bucket(const KeyType &key) const {
173         //Use hash function to determine the index of the bucket belonging to the key
174         const auto bucket_index = hash_function(key) % buckets.size();
175         return *buckets[bucket_index];
176     }
177 };

```

### Task 3: Tokenize text corpora in parallel

The web is full of unstructured human-readable text. Language models such as ChatGPT can only work on a list of words. While it seems straightforward to split a sentence represented as a string into a list of words (also strings) by just dividing them at the whitespace, the reality is much more complicated due to many exceptions. This process is called tokenization. We split an input sentence into a list of tokens (words). Your task is to tokenize multiple text corpora and count the occurrences of each token.

- Download the 20 Newsgroups data set<sup>4</sup>.
- The data from each newsgroup is in a separate directory. Each directory contains text files that contain an entry in the newsgroup.
- Process the data from each newsgroup in parallel in a separate `std::thread`.
- In each thread: Tokenize each file from the newsgroup and count the occurrences of each token with your implementation of `ThreadSafeCounterDict`.
- Run the program in parallel.
- Print the 50 most used tokens in the entire dataset to the console

<sup>4</sup><http://qwone.com/~jason/20Newsgroups/20news-18828.tar.gz>

- Implement and run the program in serial. You can use a `std::unordered_map` as a counter dictionary.
- Measure and compare the run time of the serial and parallel program
- Check the correctness of your program by comparing the results to the serial implementation
- Explain how you could further improve the performance of your program.

Take a look at `std::filesystem` for the browsing of directories, `std::ifstream` for the reading of a file, and `std::get_line` for the tokenization (separating by white spaces is sufficient). You can use the methods below to remove leading and trailing whitespaces from your token and convert them to lower case:

```

1 static inline void to_lower_case(std::string & s) {
2     std::transform(s.begin(), s.end(), s.begin(),
3     [](unsigned char c) { return std::tolower(c); });
4 }
5
6 static inline void ltrim(std::string &s) {
7     s.erase(s.begin(), std::find_if(s.begin(), s.end(), [](unsigned char ch) {
8         return !std::isspace(ch) && ch != '\n';
9     }));
10 }
11
12 static inline void rtrim(std::string &s) {
13     s.erase(std::find_if(s.rbegin(), s.rend(), [](unsigned char ch) {
14         return !std::isspace(ch) && ch != '\n';
15     }).base(), s.end());
16 }
17
18 static inline void trim(std::string &s) {
19     rtrim(s);
20     ltrim(s);
21 }
22

```

### Solution:

When a lot of threads write for the same token, they access the same bucket and, therefore, need to wait. Collecting the counts in a simple `std::unordered_map` in each thread and combining the results of all threads at the end could improve the performance further, but that depends on the word-count distribution. Usually, text corpora have words that repeat very often (e.g., 'the', 'and'). Hence, our implementation could benefit from this change.

See the implementation below:

```

1 #include <cassert>
2 #include <filesystem>
3 #include <fstream>
4 #include <iostream>
5 #include <shared_mutex>
6 #include <thread>
7 #include <tuple>
8 #include <unordered_map>
9 #include <vector>
10
11 #include "ThreadSafeCounterDict.h"
12 #include "Utils.h"
13
14 /**
15  * Function that is executed by each thread in the parallel implementation of the token counter. Counts the
16  * tokens for a single newsgroup
17  * @param counter_dict Reference to the thread safe counter dict that is shared by every thread
18  * @param directory The root directory of the newsgroup
19  */
20 void worker_func(ThreadSafeCounterDict<std::string, unsigned long> &counter_dict, const std::filesystem::
21     path &directory) {
22     //Iterate over each file from a single newsgroup
23     for (const auto &entry: std::filesystem::directory_iterator(directory)) {
24

```



```

22     if (entry.is_regular_file()) {
23         std::ifstream file(entry.path());
24         if (file.is_open()) {
25             //Iterate over each token in the file (separated by whitespaces)
26             std::string token;
27             while (std::getline(file, token, ' ')) {
28                 //Trim newline and whitespace characters from the beginning and end of the token
29                 trim(token);
30                 //Convert the token to lower case
31                 to_lower_case(token);
32                 if (token.empty()) {
33                     //Ignore empty tokens
34                     continue;
35                 }
36                 //Count the token by increasing its value by 1
37                 counter_dict.increase_value_for_key(token, 1);
38             }
39             file.close();
40         }
41     }
42 }
43 std::cout << "A thread finished\n";
44 }
45
46 /**
47  * Parallel implementation of the token counter
48  * @param root_directory The root directory of the dataset
49  * @return List of token-count pairs
50  */
51 std::vector<std::pair<std::string, unsigned long>> parallel(const std::filesystem::path &root_directory) {
52     ThreadSafeCounterDict<std::string, unsigned long> counter_dict{0, 661};
53
54     std::vector<std::thread> threads;
55
56     //Start a thread for every newsgroup
57     for (const auto &entry: std::filesystem::directory_iterator(root_directory)) {
58         if (entry.is_directory()) {
59             //Start thread with the worker function, the counter_dict reference and the newsgroup directory
60             std::thread thread(worker_func, std::ref(counter_dict), entry);
61             //Save thread object in list
62             threads.push_back(std::move(thread));
63         }
64     }
65
66     std::cout << "Started " << threads.size() << " threads\n";
67
68     //Wait for all threads to finish
69     for (auto &thread: threads) {
70         thread.join();
71     }
72
73     //Return vector of token-count pairs from the counter dictionary
74     return counter_dict.convert_to_pairs();
75 }
76
77 /**
78  * Serial implementation of the token counter
79  * @param root_directory The root directory of the dataset
80  * @return List of token-count pairs
81  */
82 std::vector<std::pair<std::string, unsigned long>> serial(const std::filesystem::path &root_directory) {
83     //Use an unordered_map for counting as it does not to be thread safe for the serial implementation
84     std::unordered_map<std::string, unsigned long> counter_dict{};
85
86     //Iterate over each file from a single newsgroup
87     for (const auto &directory: std::filesystem::directory_iterator(root_directory)) {
88         if (directory.is_directory()) {
89             for (const auto &file: std::filesystem::directory_iterator(directory)) {
90                 if (file.is_regular_file()) {
91                     //Read file

```

```

92     std::ifstream myfile(file.path());
93     if (myfile.is_open()) {
94         //Iterate over each token in the file (separated by whitespaces)
95         std::string token;
96         while (std::getline(myfile, token, ' ')) {
97             //Trim newline and whitespace characters from the beginning and end of the token
98             trim(token);
99             //Convert the token to lower case
100            to_lower_case(token);
101            if (token.empty()) {
102                //Ignore empty tokens
103                continue;
104            }
105            //Count the token by increasing its value by 1
106            counter_dict[token] = counter_dict.contains(token) ? counter_dict[token] + 1 : 1;
107        }
108        myfile.close();
109    }
110 }
111 }
112 }
113 }
114
115 //Return vector of token-count pairs from the counter dictionary
116 return {counter_dict.begin(), counter_dict.end()};
117 }
118 }
119
120 int main() {
121     std::filesystem::path root_directory = "20news-18828";
122
123     auto comp = [](const auto &p1, const auto &p2) {
124         if (p1.second == p2.second) return p1.first > p2.first;
125         return p1.second > p2.second;
126     };
127
128     //Run the serial program
129     auto begin = std::chrono::steady_clock::now();
130     auto all_pairs_serial = serial(root_directory);
131     auto end = std::chrono::steady_clock::now();
132     std::sort(all_pairs_serial.begin(), all_pairs_serial.end(),
133              comp);
134     std::cout << "Elapsed time serial = " << std::chrono::duration_cast<std::chrono::seconds>(end - begin).
135               << " [s]" << std::endl;
136
137     //Run the parallel program
138     begin = std::chrono::steady_clock::now();
139     auto all_pairs_parallel = parallel(root_directory);
140     end = std::chrono::steady_clock::now();
141     std::sort(all_pairs_parallel.begin(), all_pairs_parallel.end(),
142              comp);
143     std::cout << "Elapsed time parallel = " << std::chrono::duration_cast<std::chrono::seconds>(end - begin)
144               << " [s]" << std::endl;
145
146     assert(all_pairs_parallel == all_pairs_serial);
147
148     //Print the 50 most used tokens
149     for (auto i = 0; i < std::min(50ul, all_pairs_parallel.size()); i++) {
150         const auto &[token, count] = all_pairs_parallel[i];
151         const auto &[token_s, count_s] = all_pairs_serial[i];
152
153         std::cout << token << "(" << token_s << ")" << ": " << count << "(" << count_s << ")" << "\n";
154         assert(all_pairs_parallel[i] == all_pairs_serial[i]);
155     }
156
157     return 0;
158 }

```