

# Advanced Programming

## Mini-Introduction to Parallel Programming in C++

Parallel programming allows a program to execute multiple tasks simultaneously, improving performance, especially on multi-core processors. C++ provides several tools and libraries for parallel programming, enabling developers to write efficient, concurrent code.

This document serves as supplementary material to introduce basic parallel programming concepts in C++, focusing on threads (`std::thread` class), and OpenMP

## 1 Threads and the `std::thread` Class

A thread is a sequence of instructions that can be executed independently of other code.

### 1.1 Creating and Managing Threads

Threads in C++ are managed using the ‘`std::thread`’ class. Let’s see an example:

```
1  #include <thread>
2  #include <iostream>
3
4  // define a function to be executed by a thread
5  void complete_msg() {
6      std::cout << "task" << std::endl;
7  }
8
9  int main()
10 {
11     std::thread t(complete_msg); // Create a new thread that runs our function
12     std::cout << "Hello from: ";
13     t.join(); // Wait for the thread to finish
14     return 0;
15 }
```

The output of this code is `Hello from: task`. In this case a thread `t` is created and runs the function `complete_msg`, however, to get the output from the thread we need to use the function `join()`. This function waits for the thread to finish before continuing with the execution of the rest of the code. This is important when the output of functions running on a thread are required further in the code. Alternatively, we can let the thread run independently of the rest of the code by detaching threads. This can be useful for executing some tasks (e.g writing files) that do not affect execution.

```
1  #include <iostream>
2  #include <fstream>
3  #include <thread>
4
5  void write_hello() {
6      // Create and open a text file
7      std::ofstream msg_file("msg.txt");
8
9      // Write to the file
10     msg_file << "Hello from file";
11 }
```

```

12     // Close the file
13     msg_file.close();
14 }
15
16 int main() {
17     std::thread t(write_hello);
18     t.detach(); // Detach the thread
19     // Main thread continues without waiting for t
20     std::cout << "Hello from terminal ";
21     return 0;
22 }

```

**Important:** Detaching a thread means the main program will not wait for the thread to complete, which can lead to undefined behaviour if the thread accesses resources that are destroyed before it finishes.

## 1.2 Multiple threads

So far we have seen an example using threads for a single task. Let's see some examples creating more than one threads.

```

1  #include <thread>
2  #include <iostream>
3
4  // define functions to be executed by threads
5  void complete_msg1() {
6      std::cout << "task 1" << std::endl;
7  }
8
9  void complete_msg2() {
10     std::cout << "task 2" << std::endl;
11 }
12
13 int main()
14 {
15     // Create new threads that run our functions
16     std::thread t1(complete_msg1);
17     std::thread t2(complete_msg2);
18     std::cout << "Hello from: ";
19     t1.join(); // Wait for the thread to finish
20     std::cout << "Hello from: ";
21     t2.join(); // Wait for the thread to finish
22
23     return 0;
24 }

```

This approach is simple, but it's a little unpractical for largest number of tasks. Additionally, the output can be different as expected by just reading the code (Try yourself).

Alternatively, we can create a container of threads:

```

1  #include <thread>
2  #include <iostream>
3  #include <vector>
4
5  // define a function to be executed by threads
6  // include an argument to identify the task in the
7  // output
8  void complete_msg(int i) {
9      std::cout << "Hello from task: " << i << std::endl;
10 }

```

```

11
12 // define a function to create n threads
13 // this function returns an std::vector
14 // with each element of type std::thread
15 std::vector<std::thread> create_threads(int n)
16 {
17     // declare a vector of threads
18     std::vector<std::thread> t_vec(n);
19     for (int i = 0; i < n; i++) {
20         // define each element of the vector as a thread that runs our function
21         // note that in the second parameter we are passing the argument
22         // required by the function complete_msg
23         t_vec[i] = std::thread(complete_msg, i);
24     }
25     return t_vec ;
26 }
27
28
29 int main()
30 {
31     // define the number of threads
32     int num_threads = 12;
33     // create and run the threads
34     std::vector<std::thread> t_vec = create_threads(num_threads);
35     // remember to use join to see the output
36     // we need to use join on each thread
37     for (auto& t_i : t_vec) {
38         t_i.join();
39     }
40     return 0;
41 }

```

Try yourself and check the output, particularly the order of the output in your screen.

Using threads through the `std::thread` class provides an option for executing tasks concurrently. This approach requires appropriate memory management, control flow and exception handling, making actual implementations complex. Additionally, creation of threads makes use of some resources, hence, it's not advisable to use threads for simple tasks.

## 2 OpenMP

OpenMP allows for multi-threaded shared memory parallelisation through compiler's directives.

**Important:** Include the `-fopenmp` flag for compilation, i.e `g++ -std=c++17 -fopenmp`

Let's see an example:

```

1  #include <iostream>
2  #include <omp.h> // include the required library
3
4  // define a function to be executed by threads
5  // include an argument to identify the task in the
6  // output
7  void complete_msg(int i)
8  {
9      std::cout << "Hello from task: " << i << std::endl;
10 }
11
12 int main()
13 {
14     // define a parallel region

```

```

15     #pragma omp parallel
16     {
17         // get thread id
18         int thread_id = omp_get_thread_num();
19         // print message on screen
20         complete_msg(thread_id);
21     }
22     return 0;
23 }

```

The output of this code should be equivalent to the output of the previous example using the `std::threads` class. However, this option is shorter and simpler. The key is the use of the directive `#pragma omp parallel`. This directive creates a parallel section of code executed by all threads. The general structure for defining a parallel region is:

```

1  // serial code
2
3  // definition of parallel section
4  #pragma omp parallel // some clauses can be added here
5  {
6      // section of code to execute by threads
7  }
8
9  // more serial code

```

Contrary to `std::threads`, OpenMP automates threads creation and tasks parallelisation. By default, OpenMP creates as many threads as cores available in your computer. However, you can set the number of threads inside the code using `omp_set_num_threads(int)` or using the environment variable `OMP_NUM_THREADS`.

```

1  #include <iostream>
2  #include <omp.h>
3  #include <unistd.h>
4
5  void complete_msg(int i)
6  {
7      std::cout << "Hello from task: " << i << std::endl;
8  }
9
10 int main()
11 {
12     // define a parallel region with 5 threads
13     #pragma omp parallel num_threads(5)
14     {
15         // get thread id
16         int thread_id = omp_get_thread_num();
17         // print message on screen
18         complete_msg(thread_id);
19     }
20     return 0;
21 }

```

We have added the clause `num_threads(5)` to the `#pragma omp parallel` directive. This can be useful for some sections of code that should run on a specific number of threads, however, it can be restrictive. Alternatively, we can use the environment variable `OMP_NUM_THREADS`, e.g for a single run of the code:

```
OMP_NUM_THREADS=4 ./main
```

Since `OMP_NUM_THREADS` is an environment variable, you can set it for a terminal session (e.g using VSCode terminal) or permanently for a computer by adding it in the `.bashrc` file.

Additionally, OpenMP offers a set of directives suitable for particular applications, e.g `for` loops. For this example we will create a header file (`parallel_test.h`):

```

1  #include <vector>
2  #include <cmath>
3  #include <chrono>
4  #include <ctime>
5  #include <omp.h>
6
7  // define a class for measuring execution time of serial and parallel calculations
8  class Timer
9  {
10 public:
11     void start()
12     {
13         startTime = std::chrono::system_clock::now();
14     }
15
16     void stop()
17     {
18         endTime = std::chrono::system_clock::now();
19     }
20
21     double elapsedMilliseconds()
22     {
23
24         return std::chrono::duration_cast<std::chrono::milliseconds>(endTime -
25             startTime).count();
26     }
27
28     double elapsedSeconds()
29     {
30         return elapsedMilliseconds() / 1000.0;
31     }
32
33 private:
34     std::chrono::time_point<std::chrono::system_clock> startTime;
35     std::chrono::time_point<std::chrono::system_clock> endTime;
36 };
37
38 // define a class for calculating the average of values in a vector
39 class SumVector
40 {
41 public:
42     std::vector<double> v;
43     SumVector(std::vector<double> input_vector)
44     {
45         v = input_vector;
46     }
47     // serial
48     double serial()
49     {
50         double ave = 0.;
51         for (int i = 0; i < v.size(); ++i)
52         {
53             ave += v[i] / v.size();
54         }
55         return ave;
56     }
57     // parallel
58     double parallel()
59     {

```

```

60     double ave = 0.;
61     #pragma omp parallel for num_threads(4)
62     for (int i = 0; i < v.size(); ++i)
63     {
64         ave += v[i] / v.size();
65     }
66     return ave;
67 }
68 };

```

Note the directive is now *#pragma omp parallel for*. Additionally we have set the loop to run on 4 threads. Now, let's move to our cpp file:

```

1  #include <iostream>
2  #include <vector>
3  #include <numeric>
4  #include <algorithm>
5  #include "parallel_test.h"
6
7  int main()
8  {
9      const int n = 1000000;
10     std::vector<double> v(n);
11     int i = 0;
12     std::generate(v.begin(), v.end(), [&i]
13                 { return ++i; });
14
15     // create objects
16     SumVector sum_serial(v);
17     SumVector sum_parallel(v);
18     Timer timer_serial;
19     Timer timer_parallel;
20
21     // Calculations
22     // serial
23     timer_serial.start();
24     double result_serial = sum_serial.serial();
25     timer_serial.stop();
26     // parallel
27     timer_parallel.start();
28     double result_parallel = sum_parallel.parallel();
29     timer_parallel.stop();
30
31     // Outputs
32     // serial
33     std::cout << "Result serial: " << result_serial << std::endl;
34     std::cout << "Time: " << timer_serial.elapsedMilliseconds() << "ms" << std::endl;
35     std::cout << "Time: " << timer_serial.elapsedSeconds() << "s" << std::endl;
36     // parallel
37     std::cout << "Result parallel: " << result_parallel << std::endl;
38     std::cout << "Time: " << timer_parallel.elapsedMilliseconds() << "ms" << std::endl;
39     std::cout << "Time: " << timer_parallel.elapsedSeconds() << "s" << std::endl;
40
41     return 0;
42 }

```

If you compile and run this code you will notice a particular output. The result from the parallel calculation is wrong (smaller than the expected value), and there's no significant reduction in time (it is probably slower than serial). This is because each of the 4 threads defined are carrying out the task independently, without considering

the updates to `ave` from other threads and finishing at different times, leading to overwriting the contribution of another thread to the final value. Additionally, since threads are not allowed to write a value to a variable at the same time, each thread has to wait until another thread finishes before writing the value to a variable, resulting in an increased time. This behaviour is known as *data race*. There are different ways for handling a data race, depending on the type of operations involved in the task. In this particular case we can use the `reduction` clause for the `#pragma omp parallel for` directive, hence, we just need to change a line in our header file:

```
#pragma omp parallel for num_threads(4) reduction(+: ave)
```

After this, the parallel calculation should output the correct result with a significant speed-up. You can find more information about OpenMP directives and clauses in the [OpenMP cheat sheet](#).

### 3 Conclusion and Further Reading

Parallel programming in C++ can significantly enhance the performance of your applications, particularly on multi-core processors. However, it requires careful management of resources to avoid common pitfalls such as race conditions and deadlocks.

For further reading, consider the following resources:

- Anthony Williams, *C++ Concurrency in Action*.
- Herb Sutter, *Effective Concurrency*.
- Bjarne Stroustrup, *The C++ Programming Language* (C++11 and beyond).