

# Parallel Programming: OpenMP Lab 1

## 1 Hello, World!(Using OpenMP)

```
1 #include <iostream>
2 #include <ctime>
3 #include <omp.h>
4 using namespace std;
5
6 void main() {
7
8     clock_t t1, t2;
9     double wt1, wt2;
10
11
12     int threadsNo ,id;
13
14     t1 = clock();
15     wt1 = omp_get_wtime();
16
17     #pragma omp parallel
18     {
19         threadsNo = omp_get_num_threads(); //returns the number of threads that has been spawned
20         int id = omp_get_thread_num(); //returns the thread ID
21
22         cout<<"Hello, World! from thread "<<id<<" (out of total "<<threadsNo<<" threads)"<<endl;
23     }
24     t2 = clock();
25     wt2 = omp_get_wtime();
26
27
28     cout << "CPU Time (in seconds) = " << (double)(t2 - t1) / CLOCKS_PER_SEC <<endl;
29     cout << "wTime (in seconds) = " << wt2 - wt1 << endl;
30
31     system("pause");
32 }
```

### Output

```
Hello, World! from thread 0 (out of total 4 threads)
Hello, World! from thread 3 (out of total 4 threads)
Hello, World! from thread 2 (out of total 4 threads)
Hello, World! from thread 1 (out of total 4 threads)
CPU Time (in seconds) = 0.029
wTime (in seconds) = 0.0294827
```

- **First**, we talk about the function `omp_get_wtime()` which returns the time in seconds (as `double`), starting from some universal time (some number of years ago). You will typically see very large values ( if you displayed the value of `wt1` and `wt2`).

**wTime** is the *elapsed real time*, *wall-clock time*, or *wall time* which measures the actual time taken from the start of a computer program to the end.

- **Second**, we talk about the function `clock()` (requires to include `ctime` library) which returns CPU clock time (measured in clock ticks) since the program is started (as `clock_t`, so we need to cast it to `double`).

**CPU time** measures only the time during which the processor is actively working on a certain task, i.e. it does not include waiting for input/output (I/O) operations, entering low-power (idle) mode, or multitasking delays. *In contrast to the wTime.*

`CLOCKS_PER_SEC`: is the number of CPU clock ticks per second.

- **Third**, We discuss the problems in the code above:

1. **Problem 1**, Every thread is getting the value of the variable `ThreadsNo`, while it is enough to get its value by only one thread and share this value with the rest of threads.

**Solution 1**, Using `if` statement

```
1 void main() {
2
3     int threadsNo ,id;
4
5     #pragma omp parallel
6     {
7         id = omp_get_thread_num(); //returns the thread ID
8
9         if(id == 0)
10             threadsNo = omp_get_num_threads(); //returns the number of threads that has been spawned
11
12         cout<<"Hello, World! from thread "<<id<<" (out of total "<<threadsNo<<" threads)"<<endl;
13     }
14 }
```

### Output

```
Hello, World! from thread Hello, World! from thread 2 (out of total -858993460 threads)
1 (out of total -858993460 threads)
Hello, World! from thread 3 (out of total -858993460 threads)
Hello, World! from thread 0 (out of total 4 threads)
CPU Time (in seconds) = 0.032
wTime (in seconds) = 0.03279
```

Now this solution caused the following **problem** :

We said that only thread 0 will get the value of *threadNo* and share it with the rest of the threads but what happened (in this run) is the rest of the threads were scheduled before thread 0 so they didn't have the value of *threadNo* so they displayed garbage. And this problem has 3 solutions :

**Solution 1.1**, Using **barrier** construct with the **if** statement

```
1 void main() {
2
3     int threadsNo ,id;
4
5     #pragma omp parallel
6     {
7         int id = omp_get_thread_num(); //returns the thread ID
8
9         if(id == 0)
10             threadsNo = omp_get_num_threads(); //returns the number of threads that has been spawned
11
12         #pragma omp barrier //wait until all the threads hit this barrier and then continue
13         cout<<"Hello, World! from thread "<<id<<" (out of total "<<threadsNo<<" threads)"<<endl;
14     }
15 }
```

### Output

```
Hello, World! from thread 3 (out of total 4 threads)
Hello, World! from thread 2 (out of total 4 threads)
Hello, World! from thread 1 (out of total 4 threads)
Hello, World! from thread 0 (out of total 4 threads)
CPU Time (in seconds) = 0.042
wTime (in seconds) = 0.0431469
```

**Solution 1.2**, Making two separate parallel regions,

- 1) one for getting the number of threads using **if** statement, and
- 2) the second one for displaying the **cout** message

```
1 void main() {
2
3     int threadsNo, id;
4     #pragma omp parallel
5     {
6         id = omp_get_thread_num();
7
8         if(id == 0)
9             threadsNo = omp_get_num_threads();
10
11     } /* end of first parallel region */
12
13     #pragma omp parallel
14     {
15         cout<<"Hello, World! from thread "<<id<<" (out of total "<<threadsNo<<" threads)"<<endl;
16     } /* end of second parallel region */
17 }
```

### Output

```
Hello, World! from thread: 3 from 4
Hello, World! from thread: 3 from 4
Hello, World! from thread: 3 from 4
Hello, World! from thread: 3 from 4
```

Race Condition problem has happened. See **Problem 2** for the solution in this case! <sup>7</sup>

**Solution 1.3**, Using **single** construct instead of **barrier** and **if** together

```
1 void main() {
2
3     int threadsNo ,id;
4
5     #pragma omp parallel
6     {
7         int id = omp_get_thread_num(); //returns the thread ID
8
9         //Any thread that is scheduled first get the number of threads and share it with the other
          threads
10        #pragma omp single
11        threadsNo = omp_get_num_threads(); //returns the number of threads that has been spawned
12
13        cout<<"Hello , World! from thread "<<id<<" (out of total "<<threadsNo<<" threads)"<<endl;
14    }
15 }
```

#### Output

```
Hello, World! from thread 2 (out of total 4 threads)
Hello, World! from thread 1 (out of total 4 threads)
Hello, World! from thread 0 (out of total 4 threads)
Hello, World! from thread 3 (out of total 4 threads)
CPU Time (in seconds) = 0.033
wTime    (in seconds) = 0.0332637
```

2. A variable defined outside (inside) the parallel region is by default **shared** (**private**).

**Data race (Race Condition):** is the condition when a (**shared**) variable is, *at the same time*, accessed by different threads, resulting in an undefined behavior.

**Problem 2**, The variable *id* is **shared** among all the threads which could lead to a race condition.

```
1 #include <iostream>
2 #include <ctime>
3 #include <omp.h>
4 using namespace std;
5
6 void main() {
7
8     clock_t t1 , t2;
9     double wt1, wt2;
10
11
12     int threadsNo ,id;
13
14     t1 = clock();
15     wt1 = omp_get_wtime();
16
17     #pragma omp parallel
18     {
19         id = omp_get_thread_num(); //returns the thread ID
20
21         for(int i = 0; i < 100000; i++); //long wait so that other threads get scheduled (to make the race
          condition happen "almost" every time)
22
23         #pragma omp single
24         threadsNo = omp_get_num_threads(); //returns the number of threads that has been spawned
25
26         cout<<"Hello , World! from thread "<<id<<" (out of total "<<threadsNo<<" threads)"<<endl;
27     }
28     t2 = clock();
29     wt2 = omp_get_wtime();
30
31
32     cout << "CPU Time (in seconds) = " << (double)(t2 - t1) / CLOCKS_PER_SEC <<endl;
33     cout << "wTime    (in seconds) = " << wt2 - wt1 << endl;
34
35
36     system("pause");
37 }
```

#### Output

```
Hello, World! from thread 3 (out of total 4 threads)
Hello, World! from thread 3 (out of total 4 threads)
Hello, World! from thread 3 (out of total 4 threads)
Hello, World! from thread 3 (out of total 4 threads)
CPU Time (in seconds) = 0.031
wTime    (in seconds) = 0.0311536
```

**Solution 2**, Making the variable *id* private by either defining it inside the parallel region OR using the **private** (for *id*) and **shared** (if we want) clauses

```
1 void main() {
2
3     int threadsNo ,id;
4
5     #pragma omp parallel private(id)
6     {
7         id = omp_get_thread_num();
8
9         for(int i = 0; i < 100000; i++); //long wait so that other threads get scheduled (to make the race
          condition happen "almost every time")
10
11        #pragma omp single
12        threadsNo = omp_get_num_threads();
13
14        cout<<"Hello, World! from thread "<<id<<" (out of total "<<threadsNo<<" threads)"<<endl;
15    }
16 }
```

#### Output

```
Hello, World! from thread 3 (out of total 4 threads)
Hello, World! from thread 0 (out of total 4 threads)
Hello, World! from thread 1 (out of total 4 threads)
Hello, World! from thread 2 (out of total 4 threads)
CPU Time (in seconds) = 0.032
wTime    (in seconds) = 0.0334435
```

γ In case making two Parallel Regions, Using the **threadprivate** clause globally (for *id*) after defining the *id* variable.

**threadprivate** : Specifies that variables are replicated, with each thread having its own copy.  
Each copy of a **threadprivate** variable is initialized once prior to the first reference to that copy.

```
1 int id;
2 #pragma omp threadprivate(id) //initialize id variable for every thread that will be spawned
3 /*
4     to solve repeating getting id command in both parallel regions
5
6     it must be used globally, i.e. not inside the main function or any other function
7 */
8 void main() {
9
10    int threadsNo;
11    #pragma omp parallel
12    {
13        id = omp_get_thread_num(); //first reference of the id variable copy for each thread
14
15        if(id == 0)
16            threadsNo = omp_get_num_threads();
17
18    }/* end of first parallel region */
19
20    #pragma omp parallel
21    {
22        cout<<"Hello, World! from thread: "<<id<<" from "<<threadsNo<<endl;
23    }/* end of second parallel region */
24 }
```

#### Output

```
Hello, World! from thread: 0 from 4
Hello, World! from thread: 2 from 4
Hello, World! from thread: 3 from 4
Hello, World! from thread: 1 from 4
```

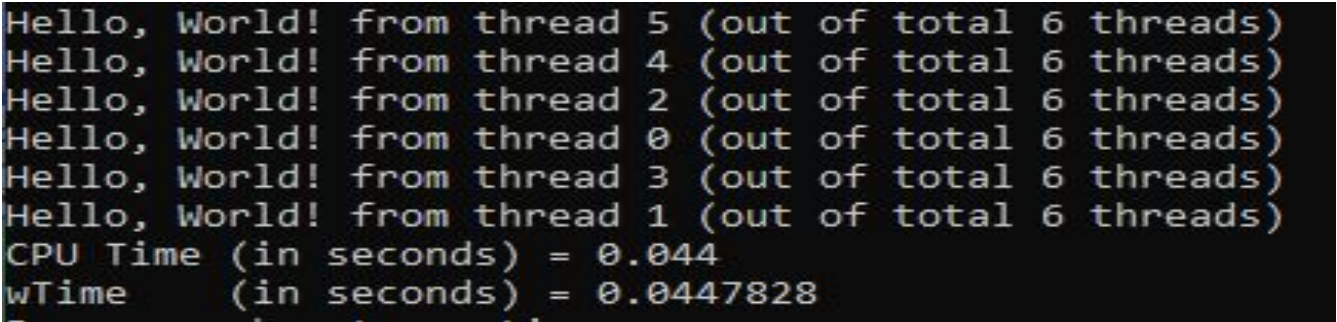
But making more than one parallel region is not a desirable thing because many forking and joining are occurred.



So the best code is:

```
1 #include <iostream>
2 #include <ctime>
3 #include <omp.h>
4 using namespace std;
5
6 void main() {
7
8     clock_t t1, t2;
9     double wt1, wt2;
10
11     int threadsNo ,id;
12
13     t1 = clock();
14     wt1 = omp_get_wtime();
15     #pragma omp parallel num_threads(6) default(shared) private(id)
16     {
17         id = omp_get_thread_num();
18
19         for(int i = 0; i < 100000; i++); //long wait so that other threads get scheduled (to make the race condition
            happen "almost" every time)
20
21         #pragma omp single
22         threadsNo = omp_get_num_threads();
23
24         cout<<"Hello, World! from thread "<<id<<" (out of total "<<threadsNo<<" threads)"<<endl;
25     }
26     t2 = clock();
27     wt2 = omp_get_wtime();
28
29
30     cout << "CPU Time (in seconds) = " << (double)(t2 - t1) / CLOCKS_PER_SEC <<endl;
31     cout <<"wTime      (in seconds) = " << wt2 - wt1 << endl;
32
33     system("pause");
34 }
```

Output



Note:

- `num_threads` clause is used to determine the number of threads to be spawned in the parallel region.
- `shared` clause is used to explicitly say that a specific variable is shared among all the threads.
- `default` clause is used to determine explicitly whether the default (i.e. if it is not explicitly determined) data-sharing attributes of the variables is private or shared.
- You can write as many clauses as you need and separate them by *space* or *comma*.  
For example: in the above program you can write:

```
#pragma omp parallel num_threads(6) default(shared) private(id)
#pragma omp parallel num_threads(6), default(shared), private(id)
```

## 2 Sum of an array

### 2.1 Sequential

```
1 #include <iostream>
2 #include <iomanip>
3 #include <ctime>
4 #include <omp.h>
5 using namespace std;
6
7 const int n = 100000000;
8
9 void main() {
10     clock_t t1, t2;
11     double wt1, wt2;
12     int *A, sum = 0;
13     A = new int[n];
14
15     for(int i = 0; i < n; i++){
16         A[i] = 1;
17     }
18
19     t1 = clock();
20     wt1 = omp_get_wtime();
21     for(int i = 0; i < n; i++){
22         sum += A[i];
23     }
24     t2 = clock();
25     wt2 = omp_get_wtime();
26
27     cout << setprecision(10);
28     cout << "sum = " << sum << endl;
29     cout << "CPU Time (in seconds) = " << (double)(t2 - t1) / CLOCKS_PER_SEC << endl;
30     cout << "wTime (in seconds) = " << wt2 - wt1 << endl;
31
32     system("pause");
33 }
```

Output

```
sum = 100000000
CPU Time (in seconds) = 0.335
wTime (in seconds) = 0.3352936
```

### 2.2 Parallel

#### 2.2.1 First attempt

```
1 void main() {
2
3     int *A, sum = 0, mySum;
4     A = new int[n];
5     //initializing the array with ones
6     for(int i = 0; i < n; i++){
7         A[i] = 1;
8     }
9
10    int threadsNo = 4;
11
12    #pragma omp parallel num_threads(threadsNo)
13    {
14        #pragma omp single
15        cout << "Using " << threadsNo << " threads:" << endl;
16
17        for(int i = 0; i < n; i++){
18            sum += A[i];
19        }
20    }
21 }
```

Output

```
Using 4 threads:
sum = 108732882
CPU Time (in seconds) = 1.52
wTime (in seconds) = 1.5197758
```

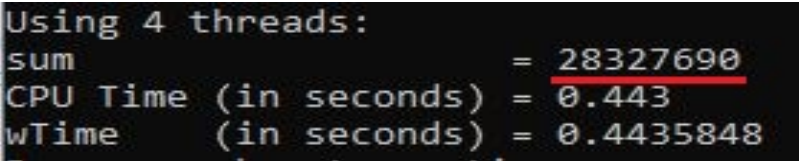
The sum is wrong because every thread that has been spawned, according to the use of `#pragma omp parallel` construct, will execute the for loop from  $i = 0$  to  $i = n - 1$  (i.e. every thread did all the work that was supposed to be divided among the threads). Note that, the time is more than the sequential version. So, we need a way to divide the work among the threads.

2.2.2 Second attempt

We divide the array entries among the threads using each thread *id* and make each thread add *only* its own entries to the shared variable *sum*.

```
1 void main() {
2
3     int *A, sum = 0;
4     A = new int[n];
5     //initializing the array with ones
6     for(int i = 0; i < n; i++){
7         A[i] = 1;
8     }
9
10    int threadsNo = 4;
11    double len;
12
13    #pragma omp parallel num_threads(threadsNo)
14    {
15        #pragma omp single
16        {
17            cout<<"Using "<<threadsNo<<" threads:"<<endl;
18            len = (double)n / threadsNo;
19        }
20        int id = omp_get_thread_num(); //returns the thread ID
21        int from = len * id;
22        int to = from + len;
23
24        for(int i = from; i < to; i++){
25            sum += A[i];
26        }
27    }
28 }
```

Output



Again the sum is wrong. Why this time?  
This is because the *race condition* on the shared variable *sum*.  
The variable *sum* must be kept shared, so how we can solve the race condition problem on a shared variable?!

2.2.3 Third attempt

We can use one of the following two construct

- 1. **atomic**: Ensures a specific storage location is accessed atomically.  
the word *atom* has a Greek origin that means *indivisible*. So, an atomic operation is the operation that is indivisible by any other operation (i.e. no other operation can happen within an atomic operation). In concurrent programming, it is said that an atomic process runs completely independently of any other processes.

**atomic** is used to solve the race condition only in the case it happens in one of the following statements:

- **Update**: like  
 $x ++, x --, ++x, --x, x \text{ bin\_op} = \text{expr}, x = x \text{ bin\_op} \text{ expr}$
- **read**:  $v = x$
- **write**:  $x = \text{expr}$   
where:  
**bin\_op** is one of the following binary operators (+, \*, -, /, &, ^, |, <<, >>), and  
**expr** is an expression of scalar type that does not reference the variable x.

- 2. **critical**: Restricts execution of the associated structured block to a single thread at a time.

**critical** is used with all kinds of statemets.

In our program here, both *atomic* and *critical* can be used since the *race condition* problem is in the statement

$sum += A[i]$

```
1 void main() {
2
3     int *A, sum = 0, mySum;
4     A = new int[n];
5     //initializing the array with ones
6     for(int i = 0; i < n; i++){
7         A[i] = 1;
8     }
9
10    int threadsNo = 4;
11    double len;
12
13    #pragma omp parallel num_threads(threadsNo)
14    {
```

```

15 #pragma omp single
16 {
17     cout<<"Using "<<threadsNo<<" threads:"<<endl;
18     len = (double)n / threadsNo;
19 }
20 int id = omp_get_thread_num(); //returns the thread ID
21 int from = len * id;
22 int to = from + len;
23
24 for(int i = from; i < to; i++){
25     #pragma omp critical
26     //OR
27     //#pragma omp atomic
28     sum += A[i];
29 }
30 }
31 }

```

#### Output

```

Using 4 threads:
sum = 100000000
CPU Time (in seconds) = 25.065
wTime (in seconds) = 25.0642369

```

Now the sum is correct. But why does it take such a long time?!

This is because when using **critical** or **atomic** constructs on the variable *sum*, it locks, updates, and unlocks the variable *sum*. This process happens as many as the entries of the array (100000000 times!) which is very time consuming.

#### 2.2.4 Fourth attempt

We make a private variable *mySum* for each thread to calculate the summation of the entries of this thread instead of summing directly into the variable *sum*. Then, we add the distinct *mySum* variables to the shared variable *sum* with using **critical** or **atomic** constructs to avoid the race condition problem.

Now the **lock-update-unlock** process on the variable *sum* only happens as many as the number of threads (Here, 4 times!).

```

1 void main() {
2
3     int *A, sum = 0, mySum;
4     A = new int[n];
5     //initializing the array with ones
6     for(int i = 0; i < n; i++){
7         A[i] = 1;
8     }
9
10    int threadsNo = 4;
11    double len;
12
13    #pragma omp parallel num_threads(threadsNo) private(mySum)
14    {
15        #pragma omp single
16        {
17            cout<<"Using "<<threadsNo<<" threads:"<<endl;
18            len = (double)n / threadsNo;
19        }
20        int id = omp_get_thread_num(); //returns the thread ID
21        int from = len * id;
22        int to = from + len;
23
24        mySum = 0;
25        for(int i = from; i < to; i++){
26            mySum += A[i];
27        }
28
29        #pragma omp critical
30        //OR
31        //#pragma omp atomic
32        sum += mySum;
33    }
34 }

```

#### Output

```

Using 4 threads:
sum = 100000000
CPU Time (in seconds) = 0.127
wTime (in seconds) = 0.1267415

```

Great!  
But OpenMP could do better...

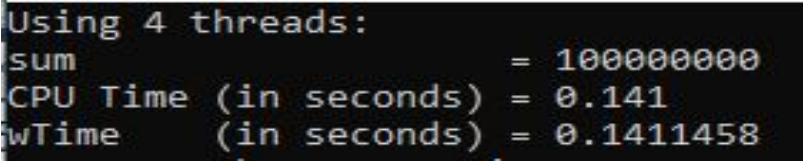


2.2.5 Fifth attempt

OpenMP provides its own for loop `#pragma omp for` that is written above the for loop in your program and OpenMP does the job of dividing the data among the threads. So, there is no need for the variables *from* and *to* and the for loop will go again from  $i = 0$  to  $i = n$

```
1 void main() {
2
3     int *A, sum = 0, mySum;
4     A = new int[n];
5     //initializing the array with ones
6     for(int i = 0; i < n; i++){
7         A[i] = 1;
8     }
9
10    int threadsNo = 4;
11
12    #pragma omp parallel num_threads(threadsNo) private(mySum)
13    {
14        #pragma omp single
15        cout<<"Using " << threadsNo << " threads: " << endl;
16
17        mySum = 0;
18        #pragma omp for
19        for(int i = 0; i < n; i++){
20            mySum += A[i];
21        }
22        #pragma omp critical
23        //OR
24        // #pragma omp atomic
25        sum += mySum;
26    }
27 }
```

Output



2.2.6 Sixth attempt

We introduce `reduction` clause. Reduction means reducing a long expression in the right hand side of an equation into only one variable in the left hand side of it. such as:

$$sum = A[0] + A[1] + A[2] + \dots + A[n - 1] + A[n]$$

The expression  $A[0] + A[1] + A[2] + \dots + A[n - 1] + A[n]$  is reduced into the variable *sum*.

- The `reduction` clause is written as follows:

$$reduction(op : var)$$

where:

**op** is the operation of reduction and it is only one of the these operations (+, -, \*, max, min).

**var** is the variable in which the reduction will happen.

So in our example we write:

$$reduction(+ : sum)$$

- The `reduction` clause is written next to:

`#pragma omp parallel` , then it will return its result (*var*) at the end of the parallel region, or `#pragma omp for` , then it will return its result (*var*) at the end of loop region.

- The `reduction` clause works (we are talking about summation) as follows:  
It does all the work that we have done until now:
  - it defines a private partial sum *mySum* for every thread,
  - adds the entries of each thread to its partial sum variable,
  - adds the partial sums to the shared variable *sum*, and
  - prevents any race condition from happening on the shared variable.

```
1 void main() {
2
3     int *A, sum = 0;
4     A = new int[n];
5     //initializing the array with ones
6     for(int i = 0; i < n; i++){
7         A[i] = 1;
8     }
9
10    int threadsNo = 4;
11
12    #pragma omp parallel num_threads(threadsNo) reduction(+:sum)
13    {
14        #pragma omp single
15        cout<<"Using " << threadsNo << " threads: " << endl;
16
17        #pragma omp for
18        for(int i = 0; i < n; i++){
19            sum += A[i];
20        }
21    }
22 }
```

Output

```
Using 4 threads:
sum                = 100000000
CPU Time (in seconds) = 0.151
wTime             (in seconds) = 0.1520898
```

OR

```
1 void main() {
2
3     int *A, sum = 0;
4     A = new int[n];
5     //initializing the array with ones
6     for(int i = 0; i < n; i++){
7         A[i] = 1;
8     }
9
10    int threadsNo = 4;
11
12    #pragma omp parallel num_threads(threadsNo)
13    {
14        #pragma omp single
15        cout<<"Using " << threadsNo << " threads: " << endl;
16
17        #pragma omp for reduction(+:sum)
18        for(int i = 0; i < n; i++){
19            sum += A[i];
20        }
21    }
22 }
```

Output

```
Using 4 threads:
sum                = 100000000
CPU Time (in seconds) = 0.15
wTime             (in seconds) = 0.1505957
```

It seems that our program only has a for loop and nothing important more is going on in the parallel region. So, we can write `#pragma omp parallel for` and then write the for loop below it directly as follows:

```
1 void main() {
2
3     int *A, sum = 0;
4     A = new int[n];
5     //initializing the array with ones
6     for(int i = 0; i < n; i++){
7         A[i] = 1;
8     }
9
10    int threadsNo = 4;
11
12    #pragma omp parallel for num_threads(threadsNo) reduction(+:sum)
13    for(int i = 0; i < n; i++){
14        sum += A[i];
15    }
16 }
```

Output

```
sum                = 100000000
CPU Time (in seconds) = 0.139
wTime             (in seconds) = 0.1385489
```