

Assignment 1

Bjørn Christian Weinbach Thameez Bodhanya

September 9, 2021

Concurrency and Non-Determinism

The code provided on studium is shown below

```
1 #include <iostream>
2 #include <mutex>
3 #include <thread>
4
5
6 std::mutex mutex;
7
8 void loop(int n)
9 {
10     mutex.lock();
11     std::cout << "Task " << n << " is running." << std::endl;
12     mutex.unlock();
13
14     mutex.lock();
15     std::cout << "Task " << n << " is terminating." << std::endl;
16     mutex.unlock();
17 }
18
19 int main()
20 {
21     std::thread t1(loop, 1);
22     std::thread t2(loop, 2);
23     std::thread t3(loop, 3);
24     std::thread t4(loop, 4);
25     std::thread t5(loop, 5);
26     std::thread t6(loop, 6);
27     std::thread t7(loop, 7);
28     std::thread t8(loop, 8);
29
30     t1.join();
31     t2.join();
32     t3.join();
33     t4.join();
34     t5.join();
35     t6.join();
36     t7.join();
37     t8.join();
38 }
```

```

39     return 0;
40 }

```

Listing 1: non-determinism.cpp

And when we compile this program as told in the exercise, we get an output similar to this:

```

1 Task 1 is running.
2 Task 3 is running.
3 Task 3 is terminating.
4 Task 1 is terminating.
5 Task 4 is running.
6 Task 5 is running.
7 Task 6 is running.
8 Task 7 is running.
9 Task 7 is terminating.
10 Task 5 is terminating.
11 Task 4 is terminating.
12 Task 6 is terminating.
13 Task 2 is running.
14 Task 2 is terminating.
15 Task 8 is running.
16 Task 8 is terminating.

```

But exactly what tasks starts running and terminating in what order is observed to be stochastic. This is known as **non-determinism** that is known to arise when we have concurrent programs. The non-determinism of concurrent programs is important to take into account as this could lead to the final output also becoming non-deterministic. Therefore we need mutual exclusion flags to ensure that the concurrent program outputs the correct solution for the problem that it is solving. In the program above this mutual exclusion flag is raised when outputting that a task is running and when it is terminating. Mutual exclusion flags solves some problems but will also lead to other problems which we will see later. Due to the problems discussed above, the exact output is undefined and we could see other orders printed out.

Shared-Memory Concurrency

We have, as stated in the assignment description, three threads that operate on a shared integer variable named **x**

The first thread **t1** operates the following code

```

1 void inc()
2 {
3     bool r = true;
4     while (r)
5     {
6         mutex.lock();
7         ++x;
8         r = run;
9         mutex.unlock();
10    }

```

```
11 }
```

Listing 2: Thread 1 source code

This thread runs an 'infinite' loop on the boolean variable *r* which is true until the boolean value *r* is equal to *run*, which we will see later in the main function. Thread 1 raises its mutual exclusion flag before it increments *x* and assigns *run* to *r* before it lowers its flag.

The second thread **t2** operates the following code

```
1 void dec()
2 {
3     bool r = true;
4     while (r)
5     {
6         mutex.lock();
7         --x;
8         r = run;
9         mutex.unlock();
10    }
11 }
```

Listing 3: Thread 2 source code

This thread does the same thing as *t1* but opposite. It loops infinitely and raises its flag when decrementing *x* and assigning *run* to *r*.

The last thread **t3** operates the following code

```
1 void print()
2 {
3     bool r = true;
4     while (r)
5     {
6         mutex.lock();
7         std::cout << x << std::endl;
8         r = run;
9         mutex.unlock();
10    }
11 }
```

Listing 4: Thread 3 source code

This thread also loops infinitely on *r* until that gets assigned false by *run*. It raises its mutual exclusion flag before outputting the value of *x* and when assigning *run* to *r*.

Lastly, let's see the main method, how it uses these threads and what we would expect as output.

```
1 int main()
2 {
3     std::thread t1(inc);
4     std::thread t2(dec);
5     std::thread t3(print);
6
7     std::this_thread::sleep_for(std::chrono::seconds(1));
8
9     mutex.lock();
```

```

10     run = false;
11     mutex.unlock();
12
13     t1.join();
14     t2.join();
15     t3.join();
16
17     return 0;
18 }

```

Listing 5: Main method

We see that the threads are started and the program sleeps for 1 second. In that time, all three threads are operating on the same shared variable. Given how the mutual exclusion flags are handled. When one thread loops another thread may be given access to the variable. Since the order of what threads will get access and for what time during the second of sleep is non-deterministic, the variable x should also have a non-deterministic value when outputted.

After five trial runs, the last outputted value of x was: 57031, -44762, -45637, 220150 and 11428. So we see that the final value of x can be considered stochastic and non-deterministic.

Race Conditions and Data Races

A race condition occurs when a program depends on the timing of one or more processes to function correctly. A data race is a special case of race condition where a shared variable is being operated on by several threads. The first program *non-determinism.cpp* does have a race condition which gives different outputs based on the order of the threads. No data race exists though as there is no shared variable that the threads operate on.

The second program, *shared-variable.cpp* does produce different outputs dependent on the sequence of execution on the threads and we therefore have a race condition. The operations on the shared variable are atomic and therefore the program does not have a data race on the shared variable. This has been verified by compiling the program with 'fsanitize=thread' and running it.

In summary: Both programs are examples of race conditions, none of them are examples of data races due to atomic operations on shared data.

Multicore Architectures

Explore the CPU

```

1 -bash-4.1$ lscpu
2 Architecture:          x86_64
3 CPU op-mode(s):        32-bit, 64-bit
4 Byte Order:             Little Endian
5 CPU(s):                 32

```

```

6 On-line CPU(s) list: 0-31
7 Thread(s) per core: 2
8 Core(s) per socket: 8
9 Socket(s): 2
10 NUMA node(s): 4
11 Vendor ID: AuthenticAMD
12 CPU family: 21
13 Model: 1
14 Model name: AMD Opteron(TM) Processor 6274
15 Stepping: 2
16 CPU MHz: 1400.000
17 BogoMIPS: 4399.38
18 Virtualization: AMD-V
19 L1d cache: 16K
20 L1i cache: 64K
21 L2 cache: 2048K
22 L3 cache: 6144K
23 NUMA node0 CPU(s): 0-7
24 NUMA node1 CPU(s): 8-15
25 NUMA node2 CPU(s): 16-23
26 NUMA node3 CPU(s): 24-31

```

Listing 6: Output from lscpu command

Based on the output above it seems that the CPU has dualthreaded cores and 8 cores per cpu. With our two physical cpus this would make 16 cores and 32 virtual CPUs. But this is an AMD processor and the architecture is a bit different than it first seems. According to cpu-world.com the CPU actually has 16 physical cores. It has L1 data caches for each core, and shared instruction caches for every two cores. The L2 cache is shared between two cores as well. L3 Cache is shared among 8 cores. See figure 1 for a visual illustration of a single AMD Opetron 6274.

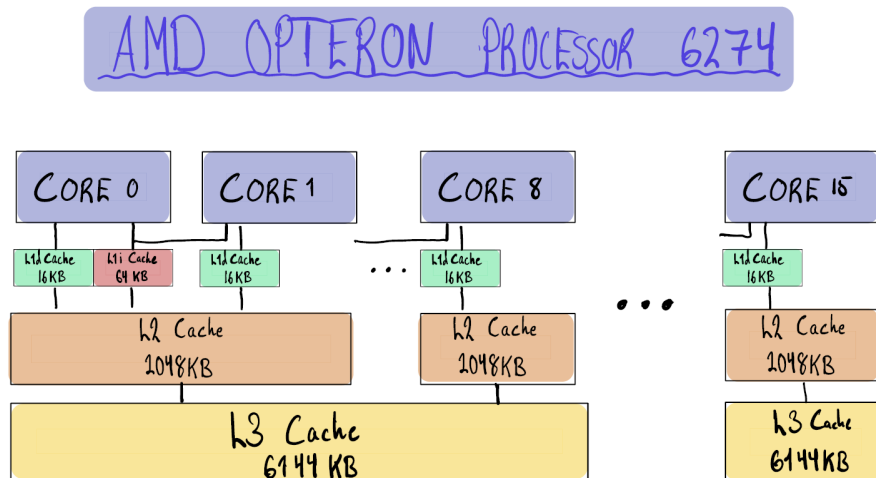


Figure 1: Cache Hierarchy of AMD Opteron Processor 6274 on gullviva.it.uu.se

We see that this processor has a lot more physical cores. In summary we have:

- There are in total 2 physical CPUs. Each with 16 physical cores.
- Each physical core has an L1d of 16KB
- There is an L1i Cache of 64KB for every two cores.
- There is a L2 Cache of 2048KB for every two cores.
- Eight physical cores share one L3 cache of size 6144KB.

Compared with the Intel cpu the cache levels are the same. The memory split is what is the main different between the two cpus. There are more physical cores in the AMD cpu but the instruction memory is shared between the cores. The CPU we are using also has a lot more cores than the intel CPU.

Performance Measurements

Script used to run program is shown in listing 7.

```
1  #!/bin/bash
2
3  set +e
4  rm results.txt
5
6  set -eo pipefail
7
8  for t in {1,2,4,8,16,32}; do
9      for n in {1,2,4,8,16,32}; do
10         echo "($t,$n)" >> results.txt
11         ./performance $t $n >> results.txt
12     done
13 done
```

Listing 7: Bash script to try all values of T and N

Runtimes

```
1  (1,1)
2  Finished in 1.21908 seconds (wall clock).
3  (1,2)
4  Finished in 2.43608 seconds (wall clock).
5  (1,4)
6  Finished in 4.90978 seconds (wall clock).
7  (1,8)
8  Finished in 9.84641 seconds (wall clock).
9  (1,16)
10 Finished in 19.692 seconds (wall clock).
11 (1,32)
12 Finished in 39.3932 seconds (wall clock).
13 (2,1)
```

```
14 Finished in 0.616555 seconds (wall clock).
15 (2,2)
16 Finished in 1.2252 seconds (wall clock).
17 (2,4)
18 Finished in 2.45917 seconds (wall clock).
19 (2,8)
20 Finished in 4.94952 seconds (wall clock).
21 (2,16)
22 Finished in 9.88132 seconds (wall clock).
23 (2,32)
24 Finished in 19.7601 seconds (wall clock).
25 (4,1)
26 Finished in 0.309365 seconds (wall clock).
27 (4,2)
28 Finished in 0.616561 seconds (wall clock).
29 (4,4)
30 Finished in 1.22603 seconds (wall clock).
31 (4,8)
32 Finished in 2.46182 seconds (wall clock).
33 (4,16)
34 Finished in 4.97059 seconds (wall clock).
35 (4,32)
36 Finished in 9.9639 seconds (wall clock).
37 (8,1)
38 Finished in 0.161366 seconds (wall clock).
39 (8,2)
40 Finished in 0.316883 seconds (wall clock).
41 (8,4)
42 Finished in 0.63344 seconds (wall clock).
43 (8,8)
44 Finished in 1.24348 seconds (wall clock).
45 (8,16)
46 Finished in 2.49317 seconds (wall clock).
47 (8,32)
48 Finished in 5.01298 seconds (wall clock).
49 (16,1)
50 Finished in 0.083777 seconds (wall clock).
51 (16,2)
52 Finished in 0.161929 seconds (wall clock).
53 (16,4)
54 Finished in 0.319613 seconds (wall clock).
55 (16,8)
56 Finished in 0.635709 seconds (wall clock).
57 (16,16)
58 Finished in 1.25654 seconds (wall clock).
59 (16,32)
60 Finished in 2.52042 seconds (wall clock).
61 (32,1)
62 Finished in 0.040681 seconds (wall clock).
63 (32,2)
64 Finished in 0.075294 seconds (wall clock).
65 (32,4)
66 Finished in 0.151482 seconds (wall clock).
67 (32,8)
68 Finished in 0.298942 seconds (wall clock).
69 (32,16)
70 Finished in 0.592317 seconds (wall clock).
```

```

71 (32,32)
72 Finished in 1.18443 seconds (wall clock).

```

Listing 8: Reported runtimes for all values of T and N

Plotted graph

The plotted times on threads, size of array and its affect on time is shown In figure 2

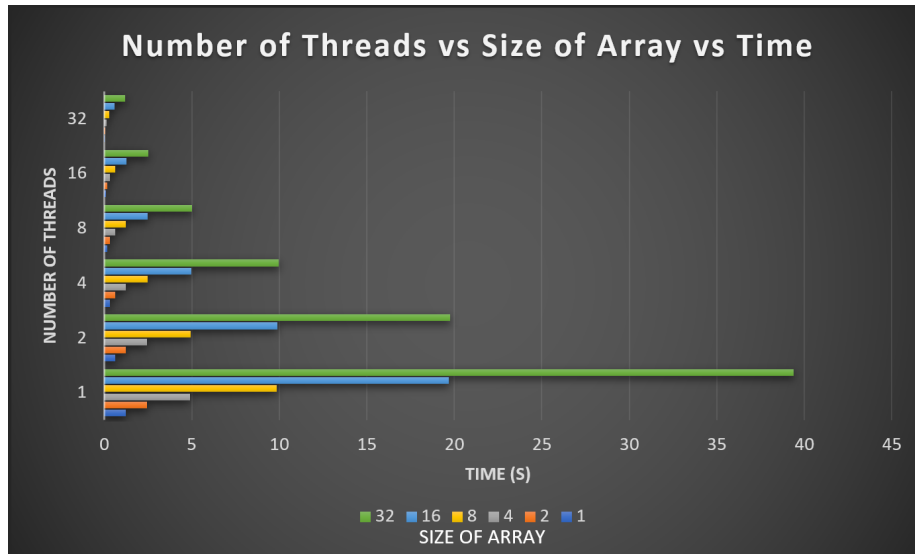


Figure 2: Threads and array size vs time.

Discussion and explanation

A quick look at the graph shows us that as expected as we increase the size of the array (while keeping thread size constant), the time taken to run the program increases almost linearly (it takes twice the amount of time to run the program for array size n^2). This pattern continues no matter how many threads we add to the program (and I presume represents the work that cannot be paralleled). Once we start varying the amount of threads, things get interesting, as we see an almost inverse relationship to above. For every thread we add, we see the time taken halved (it takes 0.5 the amount of time of array size n with threads $2t$). And this is a great demonstration of the speed parallelizing an application can yield, as the amount of time required to complete array size 32 with 32 threads is lower than the amount of time taken to run the application for array size one on a single thread.

Dining Philosophers

Dining philosophers is a known example problem in computer science formulated by Edsger Dijkstra in 1965 to illustrate a problem that can arise from concurrent programs. In the assignment there has been provided a program which simulates the dining philosophers problem. After compilation and a few trial runs we get a output similar to this which halts:

```
1 Philosopher 0 is thinking.
2 Philosopher 0 picked up her left fork.
3 Philosopher 0 picked up her right fork.
4 Philosopher 0 is eating.
5 Philosopher 0 is putting down her right fork.
6 Philosopher 0 is putting down her left fork.
7 Philosopher 0 is thinking.
8 Philosopher 0 picked up her left fork.
9 Philosopher 0 picked up her right fork.
10 Philosopher 0 is eating.
11 Philosopher 0 is putting down her right fork.
12 Philosopher 0 is putting down her left fork.
13 Philosopher 0 is thinking.
14 Philosopher 0 picked up her left fork.
15 Philosopher 0 picked up her right fork.
16 Philosopher 0 is eating.
17 Philosopher 1 is thinking.
18 Philosopher 0 is putting down her right fork.
19 Philosopher 0 is putting down her left fork.
20 Philosopher 0 is thinking.
21 Philosopher 0 picked up her left fork.
22 Philosopher 1 picked up her left fork.
```

Listing 9: Output from the dining philosophers program.

The program above had two philosophers and two forks. What has happened is that philosopher 0 picks up her left fork and philosopher 1 picks up her left fork. This means that both philosophers are at step 2 and will wait indefinitely for the other to put down their fork. This means that we have a deadlock, a state in which each member waits for the other member, including itself, to take action.

To solve the dining philosophers problem, we have implemented the Arbitrator solution. This means that we introduced a mutex that can be thought of as a waiter that does not allow a philosopher to pick up both forks or none. This comes at the cost of lost parallelism as well as adding a new mutex.

The change that we did is shown in the listing below:

```
1 arbitrator.lock();
2 left->lock();
3 std::cout << "Philosopher " << n << " picked up her left fork."
  << std::endl;
4
5 right->lock();
6 std::cout << "Philosopher " << n << " picked up her right fork."
  << std::endl;
```

```
7 arbitrator.unlock();
```

Listing 10: Arbitrator/Waiter code modification