



Guidelines and Grading Policy

You are required to submit a report describing your experiments: Make sure to include the experimental setup (Hardware and otherwise), graphs that visualize your program's performance (e.g. runtime as a function of the number of threads), and your own insightful analysis.

A program that does not compile will receive a **zero**. For every exercise, make sure you upload your solutions to your GitHub repository and maintain proper version control practices (commit your work regularly and meaningfully). You are also required to include the specifications of each method, in addition to the testing strategy and test-cases used to validate your program's correctness. The **report**, GitHub, specifications, and testing skills will account for 50% of the grade. The remaining 50% is for the program's correctness.

Counting Ones

In this assignment you will implement a program that counts and prints the total number of **ones** in an array of integers. The integers in the array are randomly generated with values between 0 and 5. Here's a sequential implementation of the program that you can use to evaluate your program's correctness:

```
1 int *array;
2 int length;
3 int count;
4
5 int count1s ()
6 {
7     int i;
8     count = 0;
9     for (i=0; i<length; i++)
10    {
11        if (array[i] == 1)
12        {
13            count++;
14        }
15    }
16    return count;
17 }
```

Your task is to experiment with different versions of a parallel implementation of the above program. You will follow a shared memory model using Posix threads. In your report, you should experiment with arrays of different sizes, namely 100, 10,000, 100,000, 1,000,000, 10,000,000, 100,000,000, and 1,000,000,000 integers. You should also experiment with different numbers of threads, namely 1, 2, 4, 8, 16, 32, and 64.

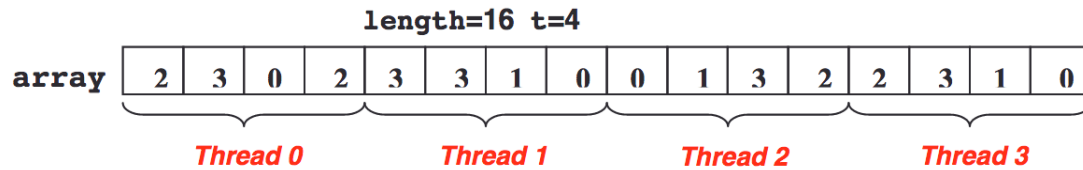


Figure 1: An example of an array being distributed across threads.

Race Condition (call your program `count_race.c`)

Start by implementing a parallel version that divides the array across the available threads. Your program should accept the number of threads as an argument. Figure 1 illustrates an example of how the array might be distributed across threads.

Note that a simple parallel implementation includes a single *count* variable that gets updated by the different threads, and is thus an incorrect implementation. Run this version of the program 100 times and report the number of times the program returns the right answer. Experiment with a different number of threads and report how the execution time changes as a result. For such experiments, you are expected to generate a plot to visualize your experiments in the submitted report.

Mutex (call your program `count_mutex.c`)

One of the ways to deal with race conditions is by using a locking mechanism that allows *mutual exclusion*, such as a *mutex*. Write a program that fixes the race condition by using a mutex to guard the critical section in the code. Perform the same experiments to ensure your program's correctness, and report the performance changes as a function of the number of threads (always include the sequential runtime as part of this comparison). In addition to the experiments and the runtime graphs, include your analysis and insight of the observed results.

Private Counts (call your program `count_private.c`)

Another way to deal with race conditions is by using private variables that are not shared across threads, and syncing the results after the threads are done with their local workload. Implement a third version of this program that introduces an array of private counts that will be accessed privately by the corresponding threads. This version should not use a mutex. Repeat the same experiments and include your analysis. Do you notice any performance improvement? Why or why not?

Padding Caches (call your program `count_cache.c`)

A cache is a special memory that is much faster to access than the computer's main memory or storage. The cache is made up of multiple levels, some of which (e.g. L1) reside on the CPU chip itself. The computer's operating system deploys many tactics to make use of this special memory to optimize programs' performance. As we should have established by now, writing parallel programs and successfully optimizing performance is not an easy task. One of the things a programmer should



be attentive to is the memory usage and bottlenecks. When a core updates a variable residing on its corresponding cache, the operating system proceeds to update the same variable on the rest of the caches corresponding to the different cores to ensure coherency. In a phenomenon called false sharing, variables that are only updated by one core may reside on the cache lines of other cores as well. This leads to extra work being done when one core updates the value of that variable. To avoid this extra work, we can make the variables occupy the entirety of the core's cache line. To do that, determine the size of the L1 cache on your computer, and create a struct that includes an integer which will be used as each thread's private count, in addition to a dummy *char* array with a size that fills the remaining bytes on the cache line. Update the private counts array to use instances of this struct. Repeat the previous experiments and include the results and analysis in your report.

In addition to your GitHub submission, add your programs to one folder: `username.a03.zip` (or `.rar`) and submit it to moodle.