



OpenMP Race Conditions

Dr. Christopher Marcotte — Durham University

Race conditions

One of the archetypal problems in simple parallel programming is the effective parallelization of sums, with independent summands. These occur in the dot (inner) product, and thus appear ubiquitously, e.g. in matrix-multiplication.

Today we will take a look at a simple example,

$$S = \sum_{n=1}^N n = \frac{N(N+1)}{2},$$

for which we have an analytical solution.

A simple serial program to implement this sum might look like this:

serial.c

```
#include <stdio.h>

int main(){
    const int N = 40000;
    int sum = 0;
    for (int n=1; n<N+1; n++){
        sum+=n;
    }
    printf("Result is %d. It should be %d.\n", sum, N*(N+1)/2);
    return 0;
}
```

If we parallelise the `for` loop in the usual way using OpenMP, we will run into a data race. This arises because multiple threads are reading from and writing to the variable `int sum` simultaneously.

Challenge

Parallelize the above program **incorrectly** with a `#pragma omp parallel for` and run it with multiple threads to convince yourself that this straight-forward approach is unreliable.

There are several ways to approach a solution. We will show a manual way to do so, and then describe a few others for you to attempt (you may wish to come back to this exercise after a later lecture).

Hint

Funnily enough, the result of the race condition depends a bit on your processor architecture, clock speed, and how many OpenMP threads you have used relative to the number of physical cores. Obviously, producing indeterminate garbage is not what we want from computers*.

If you get a correct answer for an incorrect implementation, try increasing the number of threads with `OMP_NUM_THREADS`. E.g. on my Apple Silicon based Mac, I could reliably get correct results from incorrect

*Despite all the time, money, and effort it took to generate software that can produce effectively indeterminate garbage.



implementations with small numbers of threads. I had to go up to 1024 threads to get reliably wrong answers with some of these wrong implementations.

Manual

The first approach is entirely manual. We allocate an array `int partialsums[4]` in which to accumulate the partial sums. We compute each partial sum *taking care with the OpenMP thread number (index)*. Then we accumulate into `int sum` with a serial `for` loop over the partial sum array.

manual.c

c

```
#include <stdio.h>
#include <omp.h>
int main(){
    const int N = 40000;
    int partialsums[4] = {0,0,0,0};
    #pragma omp parallel for
    for (int n=1; n<N+1; n++){
        partialsums[omp_get_thread_num()]+=n;
    }
    int sum = 0;
    for (int n=0; n<4; n++){
        sum += partialsums[n];
    }
    printf("Result is %d. It should be %d.\n",
        sum, N*(N+1)/2);
    return 0;
}
```

Exercise

The parallel fraction of the serial program is $1 - f = 0$, identically. Estimate the parallel fraction of the manual implementation program as a function of N (assuming all additions cost the same and allocations are *free*).

Use `#pragma omp [...]` `default(none)` with the `private(...)` and `shared(...)` clauses to make the data availability of each variable fully explicit.

Will this code run without OpenMP? How many threads should it be run with? Will it run with fewer than four threads? Will it give a correct answer?

We allocate a four element array in which to accumulate the partial sums – what happens if we use more than four OpenMP threads? Will it give a correct answer?

Adapt this code to work with *any* number of OpenMP threads. Estimate the new parallel fraction. Compare it to the parallel fraction you estimated earlier.

Solution

The manual implementation has 4 serial adds, and $\frac{N}{4}$ adds done in parallel across 4 threads, for $N + 4$ adds total. So the serial fraction is $f = \frac{(\frac{N}{4})+4}{N+4}$, so $1 - f = \left(\frac{3}{4}\right) - \frac{3}{N+4}$, I think.

The code will *not* run without OpenMP, since we index by thread ID – this means we require the program be compiled with OpenMP, with no graceful degradation in its absence.



If we use more than 4 threads, then we need to restrict the loop to only use the first 4 or else we'll do bad things.

In general, with M threads, I think the serial fraction is just $f = \frac{\frac{N}{M} + M}{N + M}$, so $1 - f = (M - 1) \frac{N}{M(M + N)}$.

To adapt it to M threads, you just create the array of partial sums of length M , i.e. `omp_get_max_threads()`. Find a listing of the implementation in the `racecondition-solutions.c` file below:

`racecondition-solutions.c`

```
23 /* Manual - More Threads */
24 int manualMore(int N){
25     int partialsums[omp_get_max_threads()];
26     for (int n=0; n<omp_get_max_threads(); n++){
27         partialsums[n] = 0;
28     }
29     #pragma omp parallel for
30     for (int n=1; n<N+1; n++){
31         partialsums[omp_get_thread_num()] += n;
32     }
33     int sum = 0;
34     for (int n=0; n<omp_get_max_threads(); n++){
35         sum += partialsums[n];
36     }
37     return sum;
38 }
```

c

We can see that handling this pattern manually has some additional memory management which may be considered burdensome in a more interesting program.

A `#pragma omp critical` section

An arguably simpler method is to avoid the temporary shared array `int partialsums[M]` in favor of a private variable `int partialsum` on each thread. However, once you have these private `int partialsums`, you still need to recover their values to accumulate into `int sum`. In the program snippet below, you have a hint — we have `#pragma omp parallel` generating a whole parallel block (delineated by `{...}`) — to aid you in implementing a manual reduction using OpenMP.

`critical.c`

```
#include <stdio.h>
#include <omp.h>
int main(){
    const int N = 40000;
    int sum = 0;
    #pragma omp parallel
    {
        // ??
        // ??
        for (int n=1; n<N+1; n++){
            // ??
        }
        // ??
    }
    printf("Result is %d. It should be %d.\n", sum, N*(N+1)/2);
    return 0;
}
```

c



Exercise

Replace the commented lines (e.g. `// ??`) to implement a parallel reduction using a private `int partialsum` variable and a `#pragma omp critical` region. Use `default(none)` with `private()` and `shared()` to make the data availability of each variable explicit. Estimate the parallel fraction of the implementation. Is it different from the manual approach above? Is it different from the version you implemented with arbitrary number of threads? What would your result be if you replaced `#pragma omp critical` with `#pragma omp single`?

Solution

See solution code for this implementation:

racecondition-solutions.c

```
40 /* critical region */
41 int criticalRegion(int N){
42     int sum = 0;
43     #pragma omp parallel
44     {
45         int partialsum = 0;
46         #pragma omp for
47         for (int n=1; n<N+1; n++){
48             partialsum+=n;
49         }
50         #pragma omp critical
51         sum += partialsum;
52     }
53     return sum;
54 }
```

c

For M threads, your parallel fraction should be the same as the extended manual implementation. If you replaced `#pragma omp critical` with `#pragma omp single`, then your result would be... *wrong* – your program would sum only a subset of the whole range, returning only a single partial sum corresponding to whichever thread made it to the `#pragma omp single` region first.

An `#pragma omp atomic update`

Using a `#pragma omp critical` section is not the only way to reduce the intermediate calculations. While the `#pragma omp critical` blocks protects a section of code from multiple thread accesses, `#pragma omp atomic` protects a *variable* from multiple thread accesses.

Exercise

Adapt your implementation with the `#pragma omp critical` region to use an `#pragma omp atomic update` of `int sum` from the partial sums.

Where else could the `#pragma omp atomic` pragma be placed? Why is that placement better or worse than how you structured it?

Are the data availability tags `#pragma omp [...] private(...)` and `#pragma omp [...] shared(...)` changed? Why?

Solution

You could `#pragma omp atomic` the `int partialsum` or the `int sum`, and the latter is a more performant option because it *does less serial work*.

racecondition-solutions.c

```
56 /* atomic update */
57 int atomicUpdate(int N){
58     int sum = 0;
59     #pragma omp parallel
60     {
61         int partialsum = 0;
62         #pragma omp for
63         for (int n=1; n<N+1; n++){
64             partialsum+=n;
65         }
66         #pragma omp atomic
67         sum += partialsum;
68     }
69     return sum;
70 }
```

c

A `#pragma omp reduction` operator

This sort of pattern is called a reduction or fold, and OpenMP includes a specific pragma, `#pragma omp reduction(operator:variable)` to make this pattern easy to implement. In the pragma, `operator` should be a binary associative operator (like `+`) and `variable` should be the variable into which the reduction is stored. The reduction clause makes a private copy of the variable for per-thread calculations, and takes precedence over the data access clauses – therefore, one must be careful how the data access clauses change when using a reduction.

Exercise

Adapt the serial implementation to use the OpenMP `#pragma omp reduction` pragma to perform the reduction. Use `#pragma omp [...] default(none)` and data availability tags to make the memory sharing explicit. Are they different from what you expected based on your other implementations?

Investigate the *strong scaling* of your implementations. Can you make any conclusions about bottlenecks for this calculation? How might you change it to make the parallel performance scaling better?

Solution

The `#pragma omp reduction` takes primacy for data sharing clauses, and will error if you try to use `#pragma omp [...] private(sum)` or `#pragma omp [...] shared(sum)`.

The `#pragma omp reduction` should be much more performant for large numbers of threads and sufficiently large N , but for smaller problems you may find the others finish faster. Because the `#pragma omp reduction` does the reduction in parallel (at least partially), it should also scale better than any of the other options. Additionally, the

reduction clause produces the simplest and shortest implementation. A complete listing of the reduction clause can be seen here:

racecondition-solutions.c

```
72 /* reduction operator */
73 int reductionOperator(int N){
74     int sum = 0;
75     #pragma omp parallel for reduction(+:sum)
76     for (int n=1; n<N+1; n++){
77         sum+=n;
78     }
79     return sum;
80 }
```

The speedup of the implementations on my laptop are shown in Figure 1, normalized by the fastest single-threaded run, and one can immediately tell that the calculation does not scale well. Indeed, the fastest implementation is the serial one. The reduction clause is (faint praise!) the least terrible of the options, as I said. It is not, however, a good representation of parallel performance scaling because we *are not doing enough work, per thread, for the additional threads to improve things*. On a per-core slower machine, you would find this scales considerably better; alternatively, you might consider a problem where we are calculating something substantially more intensive and interesting.

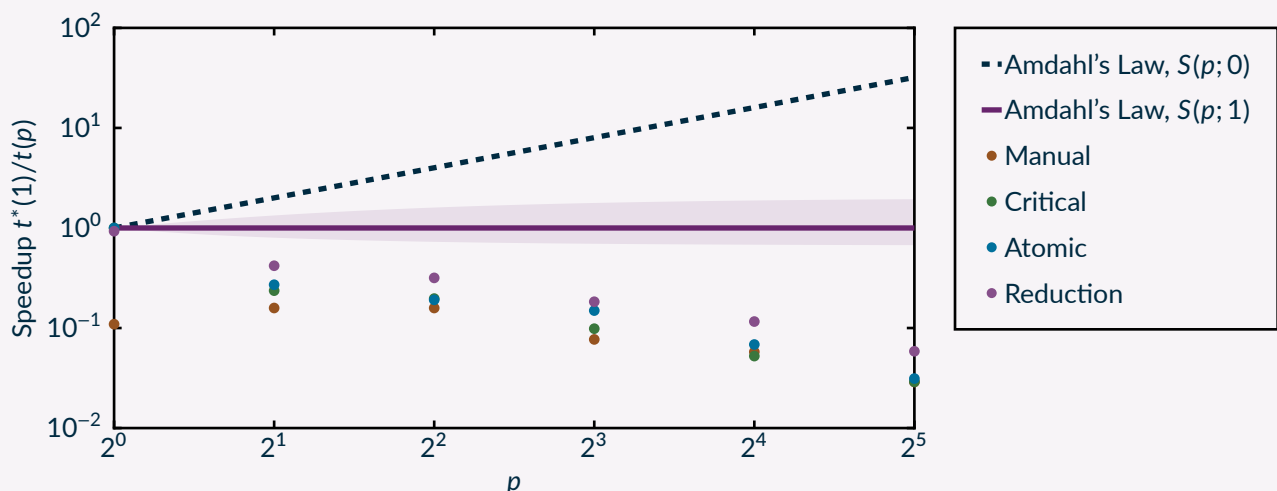


Figure 1: Speedup of reductions against the fastest serial code.

Aims

- Review of race conditions in OpenMP threaded code.
- Introduction of several ways to resolve this simple race condition.
- A gentle introduction to localisation environment variables for OpenMP threading.