# PHY905, section 4
# Homework Assignment 7

Håkon V. Treider
MSU id: treiderh

April 23, 2016

## Introduction

*The tests have been run on Mac OS X (10.11.4) with a 2.5 GHz Intel Core i7.*

In this assignment, I look into how we can use OpenMP to speed up the computation of the dense matrix-matrix multiply (algorithm) by means of parallelization.

I have run the serial code with and without OpenMP to check how much overhead it adds to the computation. Then I will use $t$ threads from 2 to 8, which is the number of logical cores on my chip. The number of physical cores is 4. I ran the code for square matrices of the order $n = 20, 200, 1000, 1300$ and the results can be seen below in table 1 or in figure 1.

Table 1: Average time of matrix-matrix multiply, in seconds

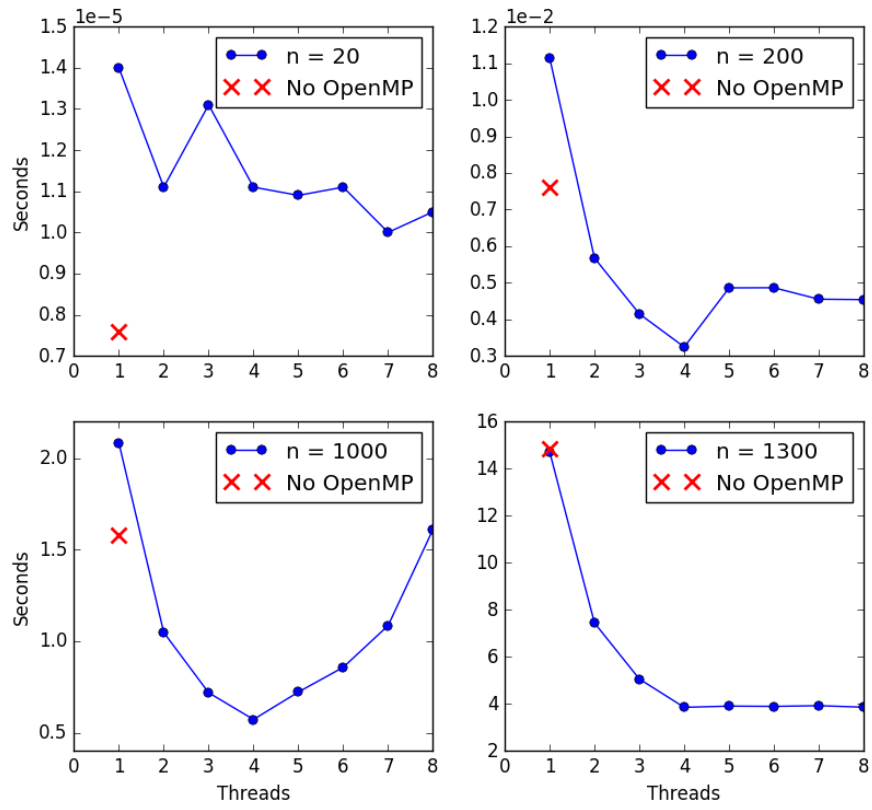| Nt | OpenMP? | $n = 20$ | $n = 200$ | $n = 1000$ | $n = 1300$ |
|----|---------|----------|-----------|------------|------------|
| 1 | No | 0.000008 | 0.0076 | 1.58 | 14.8 |
| 1 | Yes | 0.000014 | 0.011 | 2.08 | 14.7 |
| 2 | Yes | 0.000011 | 0.0057 | 1.05 | 7.44 |
| 3 | Yes | 0.000013 | 0.0041 | 0.72 | 5.06 |
| 4 | Yes | 0.000011 | 0.0032 | 0.57 | 3.85 |
| 5 | Yes | 0.000011 | 0.0049 | 0.72 | 3.90 |
| 6 | Yes | 0.000011 | 0.0049 | 0.86 | 3.89 |
| 7 | Yes | 0.000010 | 0.0045 | 1.09 | 3.92 |
| 8 | Yes | 0.000011 | 0.0045 | 1.61 | 3.86 |

Figure 1: Computation time as function of number of threads used. Comparison with serial code.

# Discussion

In the case of small matrix size (i.e. n=20), the overhead of using OpenMP is severe. So much in fact that using it, slows the computation by 20-100 %, depending on the number of threads. Going on step up on the magnitude ladder (n=200), we get a speed up when using more than one thread. This is of course because while the overhead stays roughly the same, the computation cost goes up by $\mathcal{O}(n^3)$ and since we can divide this work between multiple threads, the overall ratio goes down.

For reasonable matrix sizes, we have a distinct peak performance for *threads equal the number of processor cores*. This is best seen visually in the figure.

In homework assignment 6, we used a simple performance estimate where we assumed the total computation time scaled as $1/N_t$, up to a certain point where the computation is limited by memory bandwidth. At this point the computation time was almost constant for increasing number of threads. When we look at our result for the matrix $n = 1300$ we can see that this is also the case. The reason the three others doesn't seem to follow the same law, is because they are small enough to fit entirely in cache.

## Code

```
1   #include <libiomp/omp.h>
2   #include <armadillo>
3
4   int main(void){
5       const int n = 20; int Nt_max = 8; // Matrix size and max threads
6       arma::mat A(n,n),B(n,n),C(n,n);
7       double t0,t1,T,s;
8       double T_end_tests = 0.5; // Seconds passed before quitting loop
9       int test = 0;
10      int max_tests_default = 300000; int max_tests = max_tests_default;
11
12      // Fill A and B with numbers to check C is correctly computed
13      for (int i=0; i<n; i++){
14          for (int j=0; j<n; j++){
15              A(i,j) = i+j+1;
16              B(i,j) = 2*j+i+1;
17          }
18      }
19      // One thread, NO openMP
20      printf("Matrix %dx%d, No openMP,  ",n,n);
21
22      for (test = 0; test<max_tests; test++){
23          t0 = omp_get_wtime();          // Start timing
24          for (int i=0; i<n; i++){
25              for (int j=0; j<n; j++){
26                  s = 0.0;         // The sum for C(i,j)
27                  for (int k=0; k<n; k++){
28                      s += A(i,k)*B(k,j);
29                  }
30                  C(i,j) = s;
31              }
32          }
33          t1 = omp_get_wtime();          // End timing
34          T += t1-t0;
35          if (T>T_end_tests){max_tests = test+1;}
36      }
37      printf("Total t: %.2f, t/loop: %.7f, tests: %d, %f\n",T,T/max_tests,max_tests,C(1,1));
38
39      // 1-8 threads, with openMP
40      for (int Nt=1; Nt<Nt_max+1; Nt++){
41          T = 0;                         // Reset timer
42          max_tests = max_tests_default; // Reset counter
43          omp_set_num_threads(Nt);
44          printf("Matrix %dx%d, threads: %d, ",n,n,Nt);
45
46          for (test = 0; test<max_tests; test++){
47              t0 = omp_get_wtime();          // Start timing
48  #pragma omp parallel for
49              for (int i=0; i<n; i++){
50                  for (int j=0; j<n; j++){
51                      s = 0.0;         // The sum for C(i,j)
52                      for (int k=0; k<n; k++){
53                          s += A(i,k)*B(k,j);
54                      }
55                      C(i,j) = s;
56                  }
57              }
58              t1 = omp_get_wtime();          // End timing
59              T += t1-t0;                    4
60              if (T>T_end_tests){max_tests = test+1;}
61          }
62          printf("Total t: %.2f, t/loop: %.7f, tests: %d, %f\n",T,T/max_tests,max_tests,C(1,1));
63      }
64      return 0;
65  }
```