

Assignment 3

Bjørn Christian Weinbach Thameez Bodhanya

October 5, 2021

Instructions

For compiling the code that was provided in the assignment as well as the code that was edited for solving the exercises, the makefile provided was used. In some cases, the programs would not compile with `-std=c++11` flag. This was mentioned in the lectures. If that is the case. Replace the flag with `-std=c++0x` flag.

We have used *siegnbahn.it.uu.se* for checking that the code runs. The code is structured in the following way:

- Q1
 - sieve.cpp
 - Makefile
 - results.text
 - test.sh
- Q2
 - Game_Of_Life.c
 - Makefile
 - results.text
 - test.sh
- Q3
 - matmul.cpp
 - Makefile
 - 1-loop-parallel.text
 - 2-loop-parallel.text
 - 3-loop-parallel.text
 - test.sh

- Q4
 - colBackSub.cpp
 - colBackSubInner.cpp
 - rowBackSub.cpp
 - rowBackSubInner.cpp
 - Makefile
 - test.sh

There is a makefile for every question and subquestion. To check a question just navigate to the corresponding code and run:

Make PROG=filename

Where *filename* is the name of the C/C++ file to compile.

Question 1: Sieve of Eratosthenes - OpenMP

In this task we were asked to rewrite our implementation of the *Sieve of Eratosthenes* using OpenMP instead of Posix Threads.

- Step one of this task was to remove any references to posix threads (this included thread creation, locks and joins).
- Step two was then to implement OpenMP by using specific OpenMP constructs such as:
 - `#pragma omp parallel num_threads(numThreads)` - with the number of requested threads as a command line argument
 - `#pragma omp single` - to run the serial initiation of the seeds
 - `#pragma omp for` - to split up the workload between the threads in an as even fashion as possible
- Step three was to compile and test the code, and finally compare results

Note that one could have set the number of requested threads using the environment variable `OMP_NUM_THREADS`, however to stay true to the implementation of the original code it was decided to set it explicitly using a command-line argument.

For the complete implementation please see the attached code file *sieve.cpp*

The OpenMP implementation was far easier to write when compared to the Posix version. Many functions were able to be dropped and the code became far more readable and concise. The code is also far safer as it is easier to understand. And the written implementation allows it to be run on a single-threaded system without any code changes.

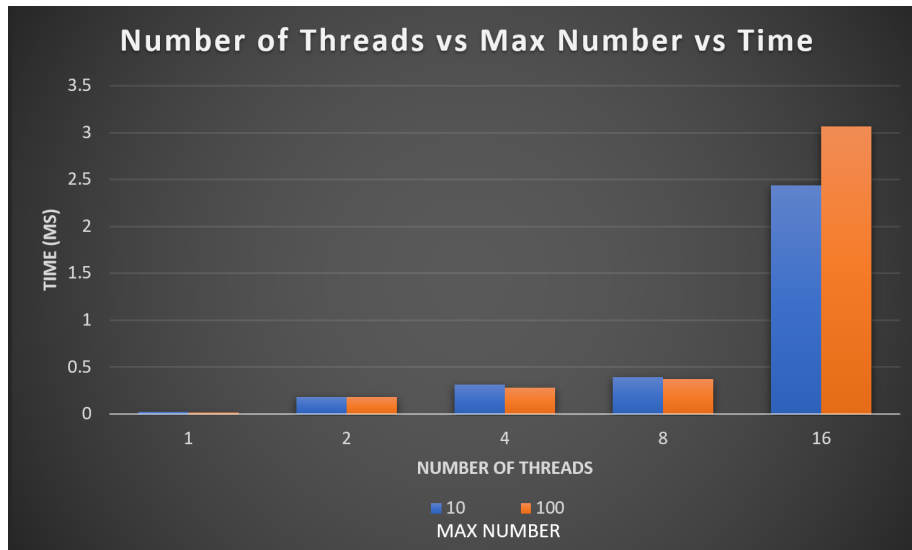


Figure 1: Sieve of Eratosthenes with max number 10 and 100.

Results

All results are from the *vitsippa.it.uu.se* system.

As per previous tests, we can see that for small max numbers (10 or 100), the overhead of threads actually increases the amount of time taken to compute the results. This is shown in figure 1. However as we hit medium max numbers (1000000 or 10000000), we already see an approx 50% reduction of time as we increase number of threads. This is shown in figure 2. And as we hit large max numbers we see that we have hit 50% reduction per addition of threads. This is shown in figure 3.

This time decrease hits substantially as we break into larger numbers as we overcome the overhead that parallelizing the application adds (and thus can take advantage of the threads) *Note: the difference between 8 and 16 threads are nominal, this indicates that the system was unable to give us a full 16 threads at testing time, however should we have gotten them, we would expect to see a full 50% reduction in time* Interestingly if we compare the runs of the Posix version to the OpenMP, we see a significant reduction in times (almost 75%). This is likely due to the less overhead OpenMP adds to parallelize the application vs our implementation Posix threads (with the additional functions and jumps). See figure 4.

Question 2: Conway's Game of Life

For this task, we want to parallelise Conway's Game of Life. This makes sense since the game works on 2D grid of arbitrary size. This means as the dimension

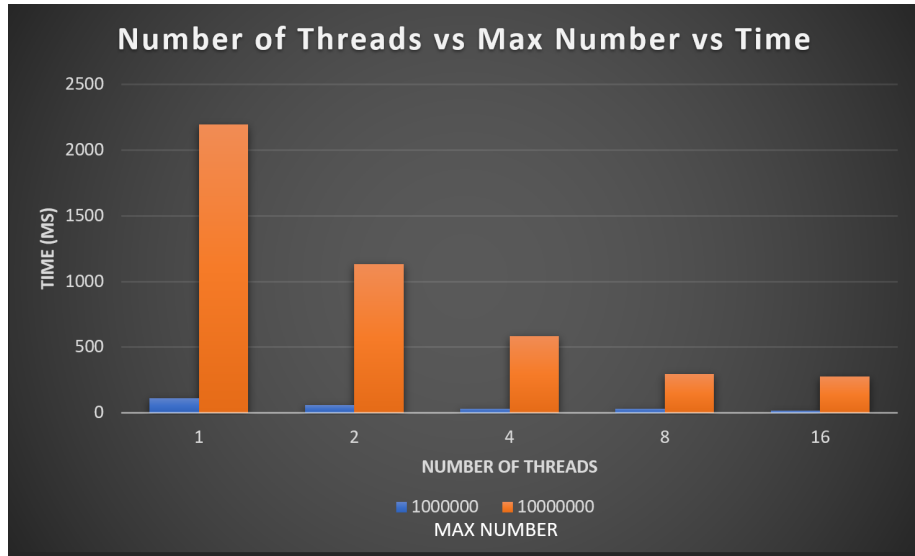


Figure 2: Sieve of Eratosthenes with max number 10^6 and 10^7 .

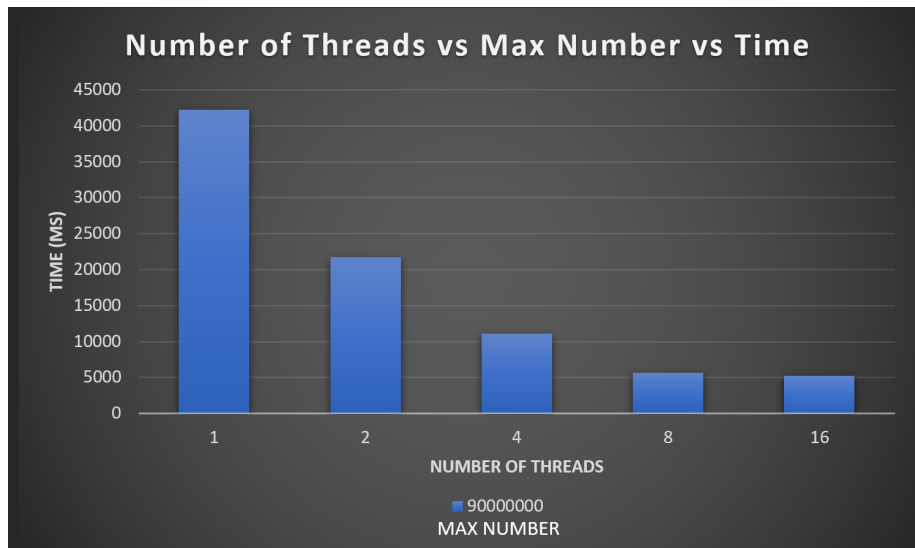


Figure 3: Sieve of Eratosthenes with max number 9×10^7 .

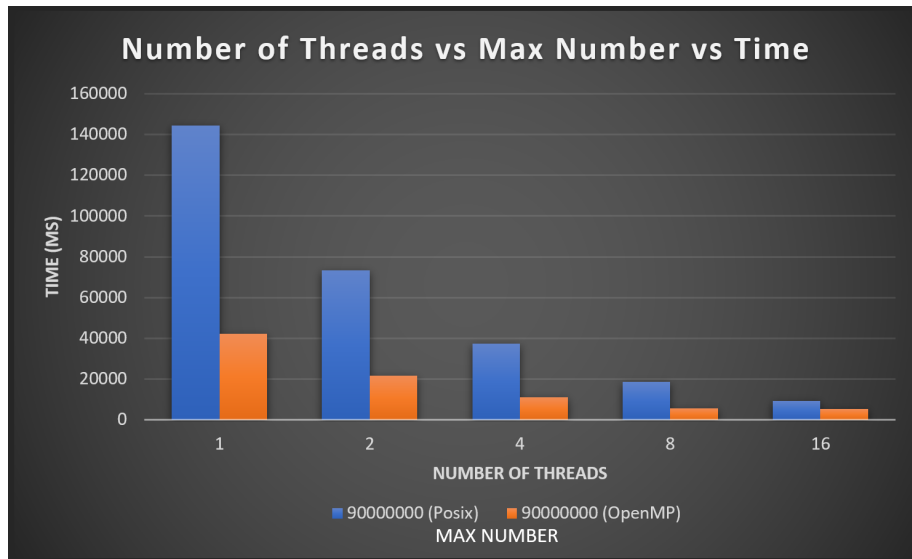


Figure 4: Posix Threads vs OpenMP.

of the grid increases by n , the overall runtime increases by $n * n = n^2$. This means that the time complexity of Conway's Game of Life is $O(n^2)$. As the new state is dependent on the previous state, the calculations at each gridpoint is independent of eachother and should be parallisable.

By using openMP this speedup did not require many changes. The changes that we did is shown in the listing below.

```

1 #pragma omp parallel num_threads(L) private(i, j, nbrs)
2 {
3     /*
4      * Pragma omp for. Collapse(2) since it is a nested for loop
5      */
6     #pragma omp for collapse(2)
7     for (i = 1 ; i < N-1 ; i++)
8         for (j = 1 ; j < N-1 ; j++) {
9             nbrs = previous[i+1][j+1] + previous[i+1][j] + previous
10                [i+1][j-1] \
11                + previous[i][j-1] + previous[i][j+1] \
12                + previous[i-1][j-1] + previous[i-1][j] + previous[
13                i-1][j+1];
14             if (nbrs == 3 || ( previous[i][j]+nbrs == 3))
15                 current[i][j] = 1;
16             else
17                 current[i][j] = 0;
18         }
19     }

```

This gave us a speedup in runtime shown in figure 5

We see that the improvement reduce drastically and slows down between 8

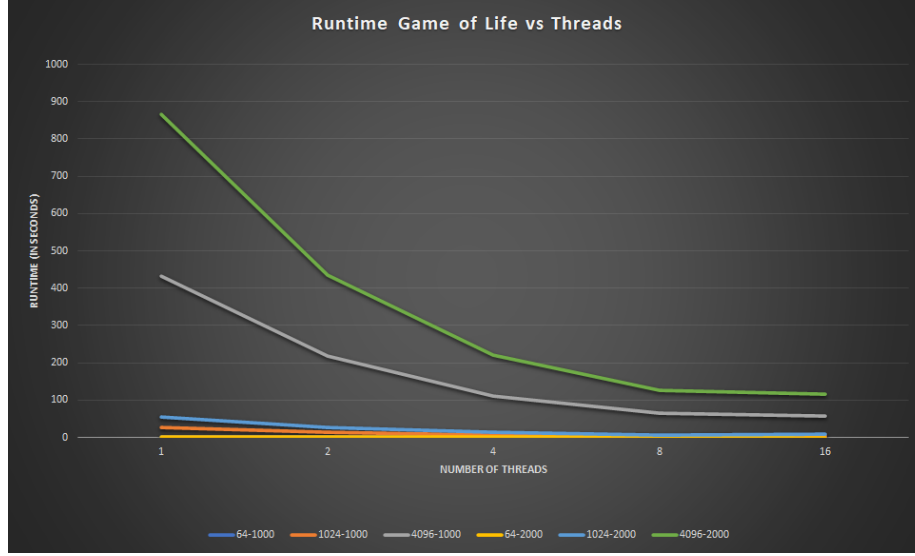


Figure 5: Runtime of Conway’s Game of Life with no of threads $t \in \{1, 2, 4, 8, 16\}$ dimensions $d \in \{64, 1024, 4096\}$ and timesteps $s \in \{1000, 2000\}$.

and 16 threads. This is due to parallel slowdown ([Wikipedia contributors, 2018]).

Q3: Matrix Multiplication

As in Conway’s Game of Life, in this task we have been provided code to parallelise. In this case, it is a matrix multiplication. Given a $m \times n$ matrix A and a $n \times p$ matrix B . $AB = C$ where C is an $m \times p$ matrix. Each value of c is calculated as follows:

$$c_{ij} = \sum_{k=1}^n a_{ik}b_{kj} \quad (1)$$

Where $i = 1, \dots, m$ and $j = 1, \dots, p$. This will require that our code has three loops to calculate all values of C . These three nested for loops is what we want to parallelise and we have done it as shown in the listing below.

```

1  #pragma omp parallel num_threads(numThreads) default(shared)
   private(i, j, k)
2  {
3      #pragma omp for schedule(static)
4      for (i = 0; i < dim; i++)
5          for (j = 0; j < dim; j++)
6              for (k = 0; k < dim; k++)
7                  c[i][j] += a[i][k] * b[k][j];
8  }
```

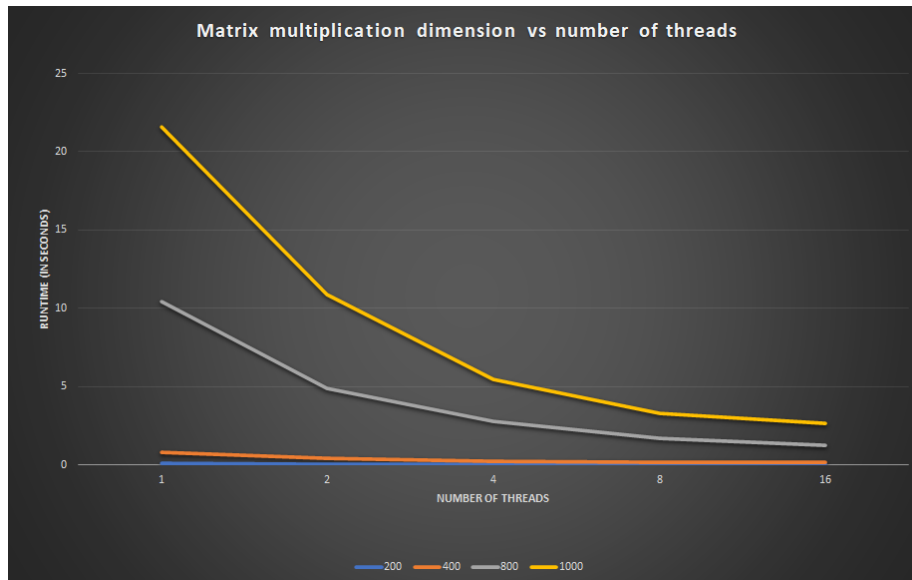


Figure 6: Runtime of matrix multiplication with no of threads $t \in \{1, 2, 4, 8, 16\}$ and dimensions $d \in \{200, 400, 800, 1000\}$.

We tried to collapse the three for loops but we saw little improvement. We are not sure exactly why this is. If only the outer loops is parallelised the threads gets assigned a chunk of that loop which has two nested for loops inside. If We collapse all three loops, the threads gets assigned a chunk of a single loop and iterates through that and calculates values of C . According to our experiments, both these cases is equally demanding for the threads and there is very little difference. Since there was no difference, only the outer loop parallelisation is plotted in figure 6.

Many factors can explain this. One is that the work done at each iteration is a simple sum that gets assigned and these are independent and a nested for loop with as many iterations as a longer single for loop should take about the same time. It is also unknown how the compiler handles this nested for loop and if it vectorises it. What we know at least, is that empirically from our little experiment, collapsing all loops does not improve the speed.

Question 4: Gaussian Elimination in OpenMP

In this task we were required to parallelize the back-substitution portion of a Gaussian Elimination algorithm. Two versions of the algorithm were presented, each containing two for-loops. Additionally once parallelized, we were asked to then gauge performance of our solutions with varying scheduling types.

Row-orientated algorithm

The algorithm to parallelize can be seen below:

```
1  for (row = n-1; row >= 0; row--)
2  {
3      x[row] = b[row];
4      for (col = row+1; col < n; col++)
5          x[row] -= A[row][col] * x[col];
6      x[row] /= A[row][row];
7  }
```

If we look at the algorithm, we see that it is not possible to parallelize the outer-loop as the inner loop relies on values generated in previous iterations of the outer-loop. Namely in line $x[\text{row}] -= A[\text{row}][\text{col}] * x[\text{col}]$, we see $x[\text{row}]$ is being updated to a value which contains $x[\text{col}]$. Should the outer-loop be parallelized, the value of $x[\text{col}]$ would depend on the order at which the threads execute, thus leading to a race condition. In addition, at the start of each thread, the value of $x[\text{row}]$ is being set to $x[\text{row}] = b[\text{row}]$ which would again influence the value of $x[\text{col}]$.

We can however parallelize the inner-loop safely, as we know that the value of col is set to $\text{row}+1$. And since the outer-loop is being run serially counting down from numThreads , the value of $x[\text{col}]$ is constant. In addition, by splitting the assignment of $x[\text{row}]$ from the heavy calculation $a[\text{row}][\text{col}] * x[\text{col}]$, we can ensure that there are no race conditions by wrapping $x[\text{row}] -= \text{diff}$ with an `#pragma omp atomic`

An implementation of this can be seen below:

```
1  for (int row = numUnknowns - 1; row >= 0; row--)
2  {
3      x[row] = b[row];
4      #pragma omp parallel for default(shared)
5      for (int col = row + 1; col < numUnknowns; col++)
6      {
7          double diff = a[row][col] * x[col];
8          #pragma omp atomic
9          x[row] -= diff;
10     }
11     x[row] /= a[row][row];
12 }
```

Column-orientated algorithm

The algorithm to parallelize can be seen below:

```
1  for (int row = 0; row < numUnknowns; row++)
2      x[row] = b[row];
3  for (int col = numUnknowns - 1; col >= 0; col--)
4  {
5      x[col] /= a[col][col];
6      for (int row = 0; row < col; row++)
7          x[row] -= a[row][col] * x[col];
```

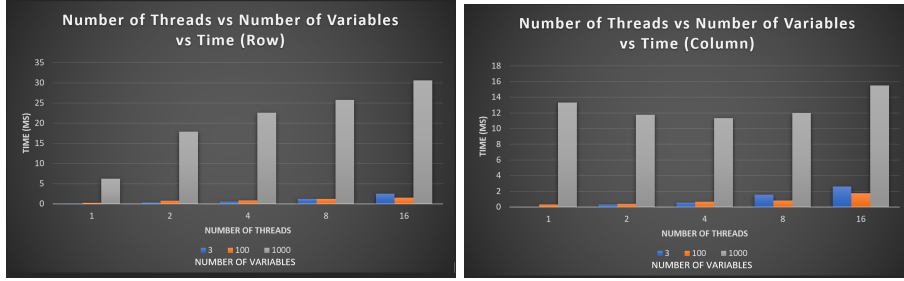



Figure 7: Column- and Row oriented performance vs no of threads.

8 }
9

If we look at the algorithm, we see that it is not possible to parallelize the outer-loop, as both the inner-loop and outer-loop update the same variable. Namely the first statement of the outer-loop $x[col] /= a[col][col]$ changes the value of $x[col]$, the current value of $x[col]$ is changed during the preceding runs of the inner-loop. This can be seen as the inner-loop counts up from 0 towards the value of col . Thus the value of $x[col]$ changes depending on the order of thread execution, since $x[col]$ is updated via division (where the order of operation does matter), this is not thread-safe and thus a data-race.

We can however parallelize the inner-loop safely, as we know that the value of col is constant as row is always less than col , thus $x[col]$ is never changed. In addition, since row is the inner-looper iterator, there are not data-races when updating $x[row]$, thus no atomic pragma is required.

An implementation of this can be seen below:

```

1  for (int row = 0; row < numUnknowns; row++)
2  x[row] = b[row];
3  for (int col = numUnknowns - 1; col >= 0; col--)
4  {
5      x[col] /= a[col][col];
6      #pragma omp parallel for default(shared)
7      for (int row = 0; row < col; row++)
8          x[row] -= a[row][col] * x[col];
9  }
10

```

Results

When compared to the serial implementation, we have vetted that the algorithm performs as expect and produces the expected results (setting u to 3 runs the application in test mode). At lower number of unknowns (3, 100 and 1000), we see an expected increase in time taken to run the program. This is true for both the row and column based approaches, see figure 7.

However, when we get to 10000 number of variables, we see that the row-based solution still increasing in time as the number of threads are added. This

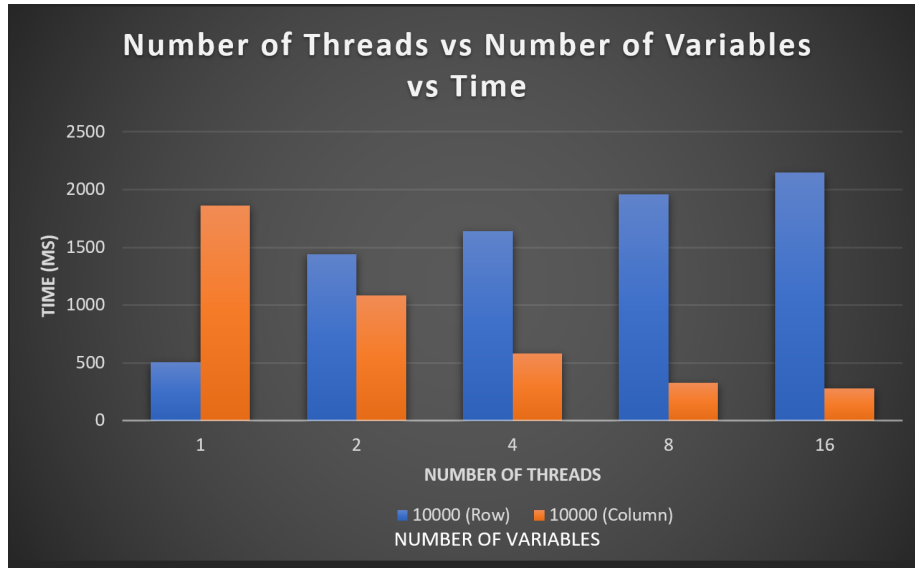


Figure 8: Column- and Row oriented comparison.

is in contrast to the column-based solution which practically halves in time as the number of threads double. A possible reason for this is that even though the inner-loop of the row-based solution has been parallelized, there exists a substantial calculation (namely $x[\text{row}] \neq a[\text{row}][\text{row}]$; is still being done serially in addition to the atomic portion of the inner-loop. See figure 8.

OpenMP Schedulers

This section documents our experimentation when running both the column-based solutions for 16 threads, with 42000 variables against the different OpenMP Schedulers

```

1  #!/bin/bash
2  set -eo pipefail
3
4  for s in {static,dynamic,guided,auto}; do
5      export OMP_SCHEDULE=$s
6      filename="results_col_${OMP_SCHEDULE}.txt"
7      set +e
8      rm $filename
9      set -e
10     OMP_NUM_THREADS=16 ./$PROG -u 42000 >> $filename
11 done
12

```

As we can see schedule type **auto** performed the best for this particular application resulting in an average increase of XXX per thread. With schedule type **dynamic** performing the worst. This would track as the amount of

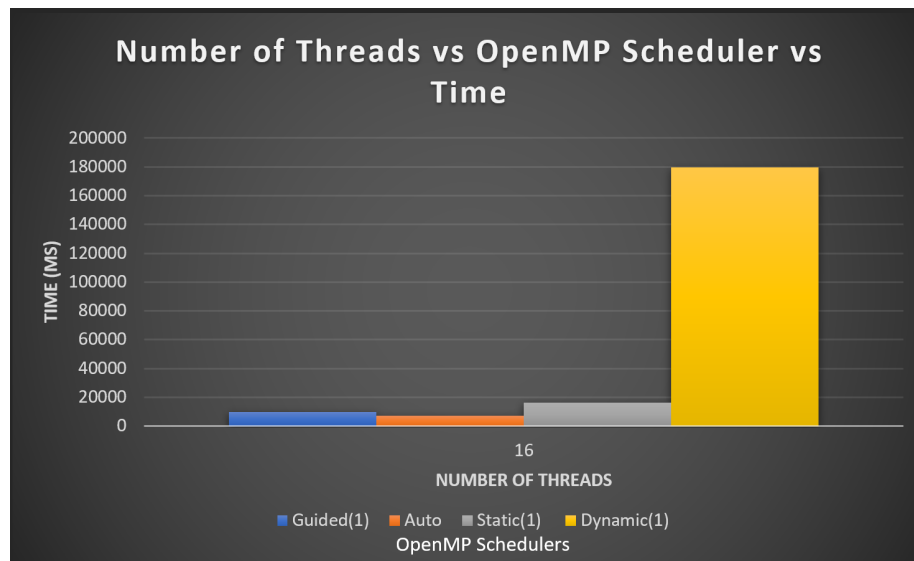


Figure 9: Column- and Row oriented performance vs no of threads.

work per iteration is vastly variable, and therefore a more generic slicing up of iteration per thread would on average perform the best. See figure 9.

References

- [Wikipedia contributors, 2018] Wikipedia contributors (2018). Parallel slowdown — Wikipedia, the free encyclopedia. https://en.wikipedia.org/w/index.php?title=Parallel_slowdown&oldid=876149288. [Online; accessed 3-October-2021].