

Matrix multiplication by using multi-threading

I. METHODOLOGY

The main point to improve the performance is not only increasing the number of threads but also making sure tasks are equality assigned into each thread. With mapping equally, it prevents main thread from waiting long time for one of thread still executing when others have completed their own tasks. Therefore, it leads in a lower CPU utilization and higher executing time.

The work allocates tasks to each thread by the total number of elements which depend on the size of matrix. In addition, to observe the multi-threading performance more flexible, dynamic memory allocation is adopted in my program to create the matrixes and threads by difference value of input (Fig 1).

```

C:\Users\Bread\source\repos\thread\Debug\thread.exe
Enter the size of matrix: 3
Enter the number of thread: 4
-----dynamic_array_1-----
0      1      2
3      4      5
6      7      8
-----dynamic_array_2-----
9      8      7
6      5      4
3      2      1
Thread 0 start from: 0 to 2 of elements Execute 3 elements
Thread 1 start from: 3 to 4 of elements Execute 2 elements
Thread 2 start from: 5 to 6 of elements Execute 2 elements
Thread 3 start from: 7 to 8 of elements Execute 2 elements
-----result-----
12      9      6
66      54      42
120     99      78
Press any key to continue . . .
    
```

Fig. 1

```

C:\Users\Bread\source\repos\thread\Debug\thread.exe
Enter the size of matrix: 3
Enter the number of thread: 6
Thread 0 start from: 0 to 1 of elements Execute 2 elements
Thread 1 start from: 2 to 3 of elements Execute 2 elements
Thread 2 start from: 4 to 5 of elements Execute 2 elements
Thread 3 start from: 6 to 6 of elements Execute 1 elements
Thread 4 start from: 7 to 7 of elements Execute 1 elements
Thread 5 start from: 8 to 8 of elements Execute 1 elements
-----TLB-----
thread_num  start point  stop point
0           0           1
1           2           3
2           4           5
3           6           6
4           7           7
5           8           8
Press any key to continue . . .
    
```

Fig. 2

Firstly, the program reads input parameters including matrix size and the number of threads to create global dynamic arrays and threads by using “new” instruction. Following that, the main program calculates the total number of elements of result array and assigns them to each thread. To minimize the executing time in each thread, main program calculates the start point and end point for every thread before thread start processing matrix multiplication, and put them into global array called TLB (Fig. 2). In this program, the work uses CreateThread() defined in windows.h to build threads. When creating threads, all of them are pointed to the same executing function but passed a unique number for each thread by using a parameter in CreateThread(). Therefore, in the thread function, by receiving the different thread number, each thread reads a different start point and stop point from TLB to get which part of element the thread is allocated to achieve computing separately and parallelly. In addition, the program can create multiple thread more flexible due to point threads to the same function. Furthermore, by using this method, the program will not happen critical section anymore due to every data in each thread is independent.

In this program, Clock() which is defined in Time.h, and adopted to calculate the time executing matrix multiplication. To estimate the running more accurately, the timer is only triggered during thread executing matrix multiplication. After checking threads and assigning tasks, main program starts the timer while all of threads are set to enable, until all the threads have been marked into finish state. As the result, the executing time can be more accurate and not affected by the other processes in main program.

II. PERFORMANCE ANALYSIS

To eliminate the tolerance of performance estimation, the work takes tests executing time for 100 times and take average value as the result.

According to the trend (Fig. 3), when executing nxn matrix multiplication by single thread, the using time remains stable with only approximately 1.5ms when matrix size is below 300. However, when the matrix size is over 300x300, the time increases significantly by growing around 1.7 time for per 100 n, due to the calculating becomes heavier to affect the performance. Therefore, it is inefficient for single thread to execute a heavy loading processing.

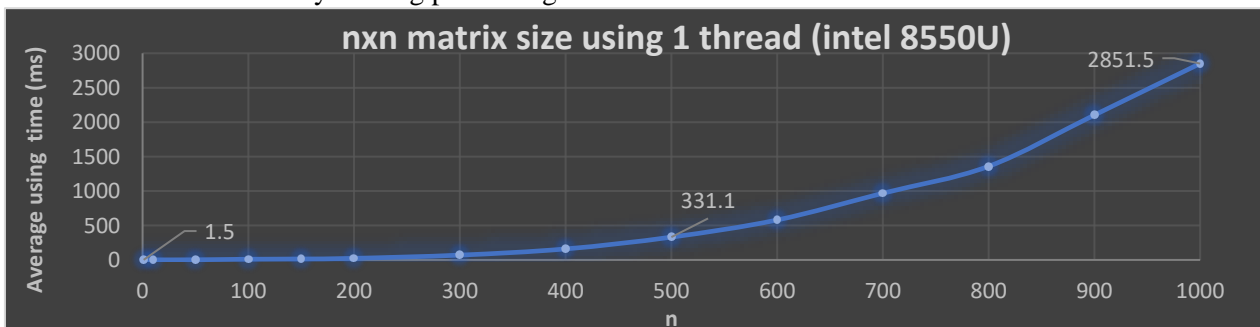


Fig. 3

- Observe the $O(n)$ performance by using threads

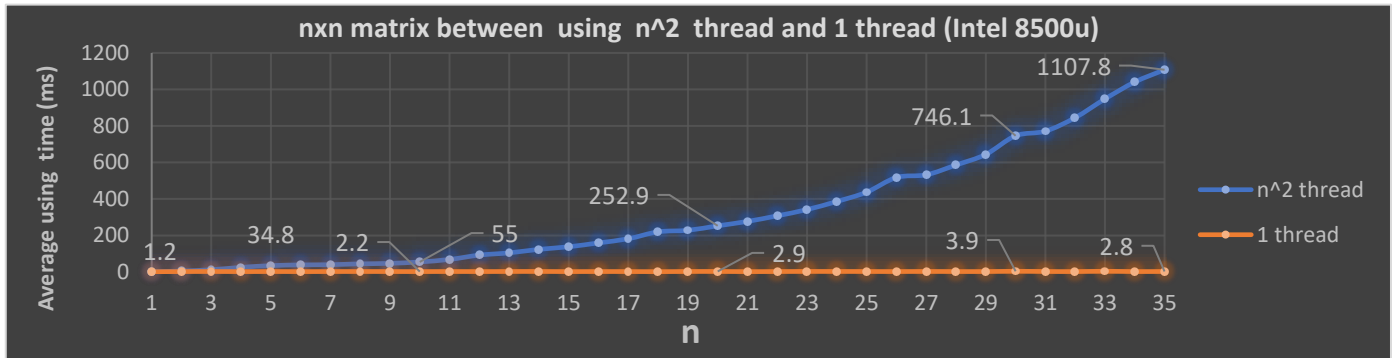


Fig. 4

The $O(n)$ notation means the limiting behavior of the function which is dominated by n , and it implies that the equation is linear function. When executing 10×10 matrix, the program creates 100 threads to execute 100 elements, and each thread spends a significant time to sleeping and synchronize by scheduler. Based on the phenomenon, the time that executes matrix multiplication by n^2 threads increases as a polynomial function and is approximately described as $f(n) = 1.2403n^2 - n + 65.654$ by Excel (Fig. 4). As a result, it is evident that the performance is not a linear function and limiting behavior becomes $O(n^2)$.

- The implication for performance when varying the number of threads

Compare to another trend in Fig. 4, when increasing the matrix size, the executing time by single thread is significantly lower than using multi-thread. The reason for this phenomenon is that, if the total number of threads is lower than 8 of the logical cores (using Intel 8550u), all threads complete their own work parallelly without any thread waiting. However, if the number of threads is over than 8, which means some of threads need to spend time waiting to be executed, every thread consumes a considerable executing time on synchronization as well (Fig. 5). Therefore, when the program executing a lower loading processing, the efficiency of using multi-thread is significantly lower than using single thread.

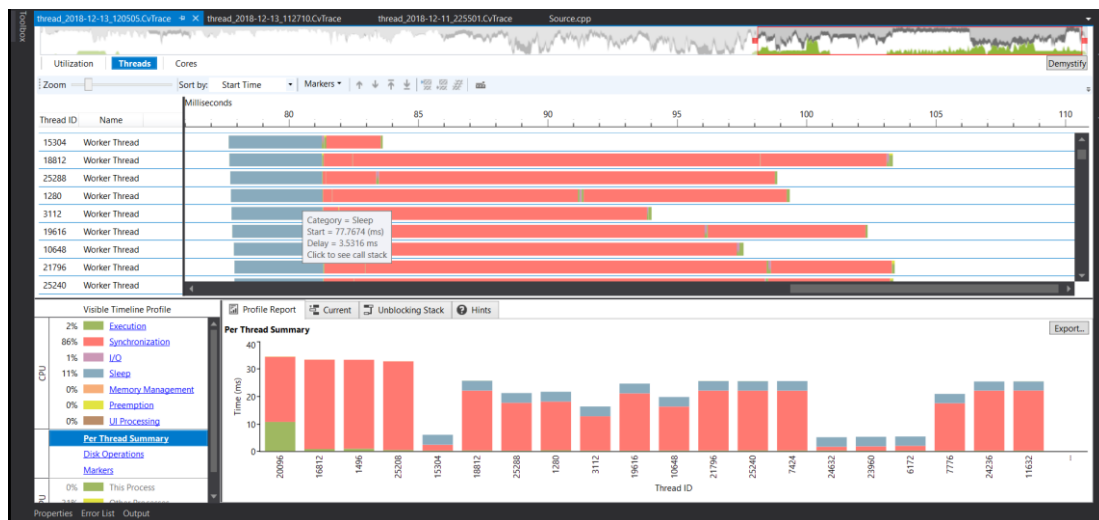


Fig. 5 multi-thread working distribution
(40x40 matrix multiplication by using 1600 thread)

- Is it possible to achieve roughly same performance by using smaller threads ($< n^2$)

As the work mention before, the using time increasing dramatically when executing heavy loading by single thread. To alleviate the issue, multi-thread plays an imperative role to enhance the performance. According to fig. 6, in intel 8550u, it cost around 2.85 second for single thread to calculate 1000×1000 matrix multiplication. However, the time drops to 1.4 second when adopting 2 thread to

calculate and stays constant with only 1.1sec by using 50 threads. As the result, due to heavy loading, the characteristic of thread system can be observed easily that with multi-threading, tasks are executed parallelly and each thread calculates their own task in the same machine cycle. It implies that the running time can be reduce efficiently by using multi-thread. However, due to the synchronization, even adopting 50 threads, the time remains at a constant and grows up again while each thread starts to consume time on synchronization.

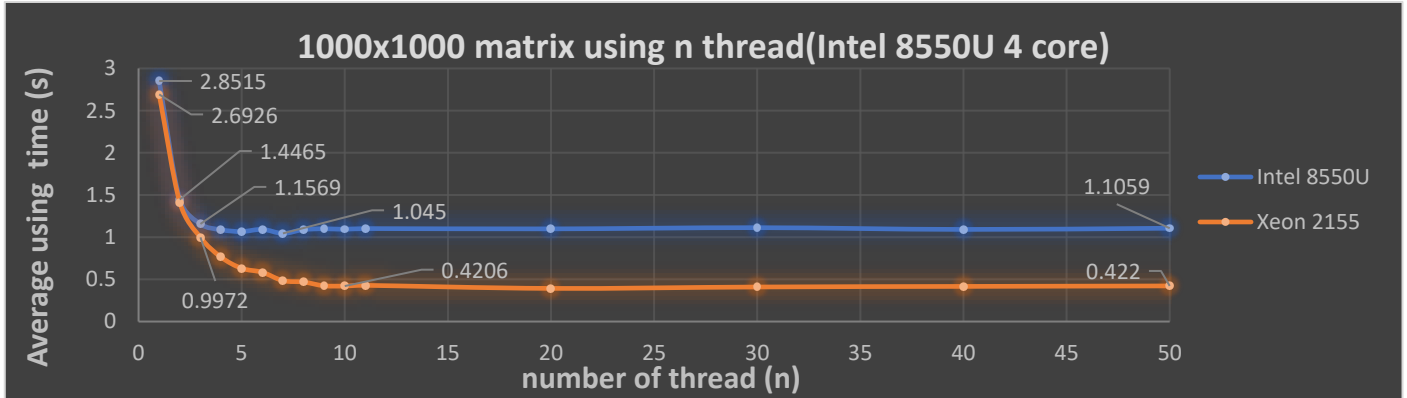


Fig. 6

• Implication by physical cores and cache size

According to fig. 6, the orange figure shows the using time when executing 1000x1000 matrix using intel Xeon 2155 CPU (10 cores 2 physical threads in each core). It is evident that due to the number of physical cores is more than intel 8550u (4 cores 2 physical threads in each core), the executing times can be decreased considerably with only 0.42 second by using 50 threads. When both of computer using double threads, the time depends on the frequency of CPU. With a large number of logical cores, it implies more and more threads are allowed to be allocated to each core in the same machine cycle. In addition, if cache memory is bigger, it suggests each thread would not spends too much time on I/O due to more and more data can be stored in cache at a time.

• Discussion

Initially, in my program, the work assigns the task to threads by dividing and conquering matrix to calculate multiplication such as Strassen's method. However, mapping tasks with this method will become more complicated when increases the number of threads. Alternatively, to make algorithm more intuitive and understandable, the work allocates tasks by dividing total number of elements to threads.

Due to compiler only can allocate 2GB address spaces for program, there is a limitation about number of creating multi-thread. Even a thread had no code and data and the entire memory space are occupied stacks, the program can create at most 2048 threads. Base on the restriction, the work uses Check_error() to check whether a thread can be created successfully after executing CreateThread(). If the check function occurs error, it returns 1 to let main program know it hits the limitation number of threads. After confirming all threads, main program reassigned the tasks to prevent some of elements from missing or pointing to an empty data address.

III. CONCLUSION

In conclusion, to evaluate whether thread system is a good technology, the answer depends on the process loading. When executing nonheavy load processing, it is worth to adopt single thread. Nevertheless, for multi-thread, due to each thread spending a significant time on synchronization, it is inefficient to use nxn threads to execute matrix multiplication. However, if the number of threads is below number of logical cores, it contributes the highest performance and utilization due to each thread execute parallelly without consuming lots of time on synchronization.