
MATRIX MULTIPLICATION REPORT

January 17, 2020

Cheng Chen (605428367)

0.1 ABSTRACT

This report summarizes the results of parallel matrix multiplication and the blocked version of it. The results are compared against the sequential version.

0.2 RESULTS SUMMARY

0.2.1 Problem size 1024³

| Performance metric | Sequential | Parallel | Parallel-blocked |
|--------------------|------------|----------|------------------|
| GFlops | 0.649662 | 81.449 | 60.1233 |
| Speedup | 1 | 125.37 | 92.5455 |

0.2.2 Problem size 2048³

| Performance metric | Sequential | Parallel | Parallel-blocked |
|--------------------|------------|----------|------------------|
| GFlops | 0.292448 | 87.4917 | 77.5286 |
| Speedup | 1 | 299.17 | 265.102 |

0.2.3 Problem size 4096³

| Performance metric | Sequential | Parallel | Parallel-blocked |
|--------------------|------------|----------|------------------|
| GFlops | 0.19768 | 44.7062 | 117.761 |
| Speedup | 1 | 226.154 | 595.715 |

0.2.4 Summary

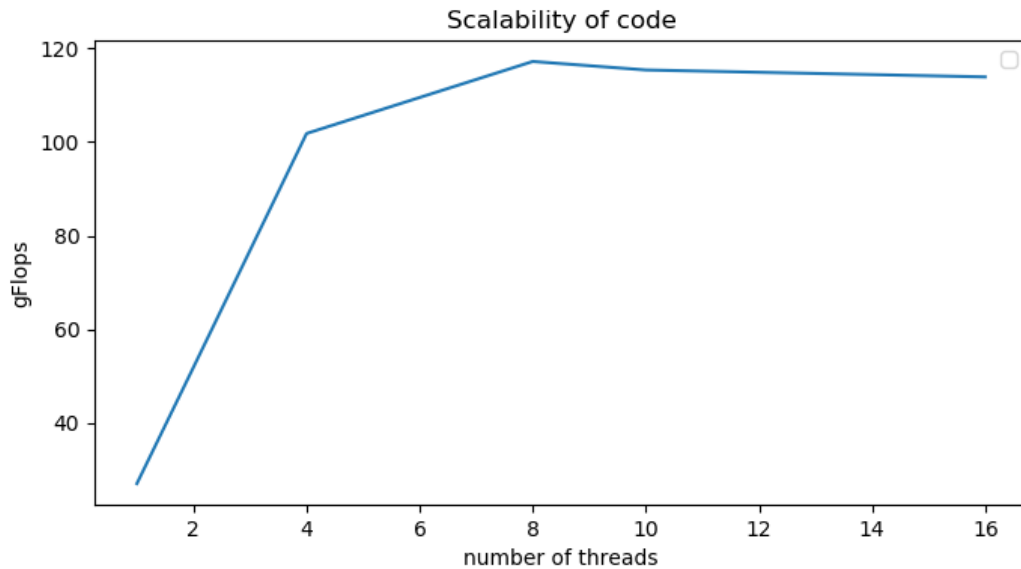
It's clear that in all of 3 problems sizes, paralleled versions are much better than sequential one. However, one may notice that only when problem size is 4096³, parallel-blocked one is better than not-blocked version. This is because when problem size is not that large, a whole line of data can be cached all together, so the blocked version cannot increase the spatial locality.

0.3 IMPACT OF EACH OPTIMIZATION

| Optimization | GFlops | Speedup |
|----------------------|---------|---------|
| not blocked | 44.7062 | 1 |
| block size optimized | 117.761 | 2.625 |

0.4 CODE SCALABILITY

The results are shown below as a plot:



This result is reasonable as m5.2xlarge supports 8 threads, so the performance is maximized when number of threads is set to be 8. Note that the performance does not drop much when number of threads is greater than 8. As only 8 threads can essentially proceed in parallel, more threads will need to wait for some other threads to finish. This will create some overhead, but not that much.

0.5 CONCLUSION

Several optimizations are made to speedup the code:

- omission of C array initialization: in C++, the array of numbers have value of 0 by default, so there is no need to initialize the C array
- loop permutation: make k-loop at the lowest level to utilize spatial locality
- static scheduling: as the workload more or less remains the same, static scheduling is better to reduce the overhead
- number of threads: performance is maximized when number of threads aligns with number of CPUs

- blocking: blocked version improves the performance to a large extent as spatial locality is increased

It is clear that the optimized parallelization speeds up the code to a great extent, and the blocked version only has advantage if the problem size is large enough.