

MCSC 6030G : High Performance Computing

Assignment 1: Matrix-Vector Product

Parikshit Bajpai
100693928

1 Introduction

Computing, in today's world, has radically changed since the days of Konrad Zuse who constructed the world's first fully automated, freely programmable computer with binary floating-point arithmetic in 1941 [1]. Zuse's great visions regarding the possible use of his device are now a reality, and computers, owing to their computational abilities at an incredible, ever-increasing speed, have become essential, not only in science and engineering but in all sectors of life [2].

As an essential tool for most research areas, both in academia and industry, computer-based simulations have become ever-present standard tools. With the rapid increase in the performance of computing systems, more and more problems of scientific interest come within our reach. Many of these problems have requirements for raw computational speed, storage, and/or main memory which can only be met by dedicated computational frameworks [3]. High Performance Computing, in general, refers to the practice of clustering computing power in a way that performs computations at significantly faster speeds than one could get out of a typical desktop computer or workstation in order to solve large problems in science, engineering, or business [4].

In times of stagnating single processor capabilities and increasing parallelism, there has been a growing focus on the performance and scalability, and the most sensible measure of performance in benchmarking is the *wallclock time*, also called *elapsed time* [3]. In order to reduce wallclock times, HPC systems rely on a few techniques such as parallelisation, compiler optimisation, and use of standard libraries.

Parallel processing refers to processing the program instructions by splitting them into multiple threads in order to reduce execution time. The different threads can perform different subtasks simultaneously thus reducing the time required for the execution. OpenMP is a shared memory parallel programming API that works on the fork-join model and consists of a set of compiler directives that a non-OpenMP-capable compiler would just regard as comments and ignore [5]. In the context of parallelisation, apart from the wallclock time, the following parameters are of interest:

$$Speed\ Up = \frac{Serial\ Wall\ Time}{Parallel\ Wall\ Time} = \frac{S}{P} \quad (1)$$

$$Efficiency = \frac{Speed\ Up}{Number\ of\ Threads} = \frac{S}{n \times P} \quad (2)$$

Apart from parallelisation, optimising compilers has become an essential component of modern high-performance computer systems. In addition to translating the input program into machine language, the compilers analyse it and apply various transformations to reduce its running time or its size [6]. Amongst the popular compilers, GNU Gfortran has developed an open-source *gfortran* compiler capable of working on multiple architectures and diverse environments, and, Intel has

developed the *ifort* compiler to deliver superior Fortran application performance and to boost Single Instruction Multiple Data (SIMD) vectorisation and threading capabilities [7].

2 Methodology

2.1 Objective

The present study is aimed at exploring the impact of code, compiler and execution on the wall time for matrix-vector multiplication problem. To this aim, compiler optimisation, parallelisation, and use of standard software library have been implemented and the wall times have been compared.

2.2 Machine Configuration

Manufacturer & Model: Lenovo ThinkPad Yoga 370

Processor: Intel Core i5 -7200U (4 processor cores)

Clock Rate: 2.50 GHz

RAM: 16 GB

Operating System: Ubuntu 18.04

2.3 Implementation

In order to obtain a reference performance benchmark, the standard double-loop code for matrix-vector multiplication was implemented without optimisation and parallelisation and compiled using the GNU compiler *gfortran*, and the matrix size for which the wall time obtained was of the order of 10 seconds was selected for further study. The matrix size corresponding to the aforementioned wall time was 8192×8192 . Once the matrix size was fixed, compiler optimisation, parallelisation and use of standard library was implemented and the obtained wall times were compared with that obtained for the double loop without optimisation case. For this purpose, the following primary cases and their combinations were considered:

1. Choice of compiler: GNU *gfortran* or Intel *ifort*
2. Compiler optimisation: no optimisation (-O0) or aggressive optimisation (-O3)
3. Parallelisation: 1, 2 or 4 OpenMP threads
4. Matrix multiplication approach: double loop or BLAS standard library

The BLAS (Basic Linear Algebra Subprograms) are routines that provide standard building blocks for performing basic vector and matrix operations [8] and LAPACK routines are written so that as much as possible of the computation is performed by calls to the Basic Linear Algebra Subprograms (BLAS) [9]. The wall times obtained based on the above strategies used individually and in conjunction were then used to calculate the speed up and efficiency of parallelisation based on equation 1 and equation 2 respectively.

Ibrahim Guiagoussou, Marcos Machado, Celina Desbiens and myself worked together on building the codes, executing and interpreting the results.

3 Results and Discussion

The source code was run for a number of different cases presented above and the obtained results have been presented in Table 1. The results highlight that the wall times obtained using the Intel compiler are significantly less than those obtained using the GNU compiler. This can be attributed to the optimisation of Intel compilers for systems using Intel processors and the high-level techniques used by the compiler to reduce stalls and produce codes that execute in the fewest possible number of cycles [7].

Table 1: Matrix vector multiplication wall times [s] for the different approaches under study.

			Open MP			BLAS
			1	2	4	
GNU	O0	12.757	12.731	07.669	05.652	00.190
	O3	12.950	12.887	07.862	05.824	00.198
INTEL	O0	17.028	17.019	08.808	07.475	00.190
	O3	06.870	00.167	00.143	00.142	00.191

Table 2: Speed-up and efficiency of parallelisation

			Open MP			
			2		4	
			Speed-up	Efficiency	Speed-up	Efficiency
GNU	O0	1.663	0.832		2.257	0.564
	O3	1.647	0.824		2.224	0.556
INTEL	O0	1.933	0.966		2.278	0.569
	O3					

Parallelisation was implemented using OpenMP directives and we observe a significant reduction in the wall times as the number of threads was increased. However, as presented in Table 2, we also notice that the speed up did not increase linearly when the number of threads was increased from 2 to 4. Such a non-linear increase can, most probably, be attributed to increase in overheads as the number of threads is increased. Use of the BLAS library with *gfortran* compiler resulted in the most significant decrease in wall time in absolute terms. However, the wall times using OpenMP and BLAS concurrently were not studied. This can be justified using the fact that BLAS code calls a single subroutine, *DGEMV*, for the matrix-vector multiplication and parallelising a single function call does not make any practical sense. Furthermore, BLAS is a highly optimised library which is itself a benchmark for performance tests and no valuable information will be achieved even if we use OpenMP directives in the BLAS code.

Interestingly, when using BLAS library, we observe that the wall times achieved using *gfortran* compiler were lower than those obtained using *ifort* compiler without optimisation. The minimum wall-times were, however, achieved using double-loop with OpenMP directives and compiled using Intel compiler with aggressive optimisation.

Note: The values obtained for the test case with *ifortran* compiler using aggressive optimisation seem erroneous. Specifically, the wallclock time obtained using the double-loop without parallelisation is much higher than expected and is incoherent with the wallclock times obtained with

parallelisation. The codes were compared to find the cause of this observation but such a solution can not be expected from a difference in the codes and the bug was not found at the level of code. Therefore, the error can, for now, be attributed to either the compilation process and the generated machine language code or to the machine itself. Therefore, the speed-up and efficiency values for the case in point have been omitted.

4 Conclusion

In general, the impact of compiler optimisation, parallelisation, and use of standard library on the performance is evident from the obtained results. The results of the different cases reveal that, as foreseen, the Intel *ifort* compiler is more efficient than the GNU *gfortran* compiler. Furthermore, we observe that both parallelisation and use of standard library significantly speeds up the code. The tests highlight the importance of the aforementioned ways of improving the performance of codes developed for HPC and show how these can be used individually and in conjunction with each other to achieve the best possible performance from the available computing device. In conclusion, we can say that, in this specific case, the double-loop code with Intel compiler and aggressive optimisation was the fastest. However, in general, BLAS library offers a much more reliable method of matrix-vector multiplication and performs equally well with both the compilers under consideration and for both the optimisation methods. Therefore, using BLAS and LAPACK libraries should be the preferable method for matrix-vector multiplication.

References

- [1] Raúl Rojas and Ulf Hashagen, editors. *The First Computers—History and Architectures*. History of Computing. The MIT Press, 2002.
- [2] Konrad Zuse. *The Computer – My Life*. Springer-Verlag Berlin Heidelberg, 1993.
- [3] Georg Hager and Gerhard Wellein. 'Modern Processors' in *Introduction to High Performance Computing for Scientists and Engineers*. Chapman & Hall/CRC Computational Science. CRC Press, 2010.
- [4] “What is high performance computing?”, insideHPC. [Online]. Available: <https://insidehpc.com/hpc-basic-training/what-is-hpc/>. [Accessed: 05-Nov-2018].
- [5] OpenMP Architecture Review Board. OpenMP Application Program Interface version 4.5, November 2015. [Online]. Available: <https://openmp.org/wp-content/uploads/openmp-4.5.pdf>. [Accessed: 05-Nov-2018]
- [6] David F. Bacon, Susan L. Graham, and Oliver J. Sharp. Compiler transformations for high-performance computing. *ACM Computing Survey*, 26(4):345–420, December 1994.
- [7] Intel Corporation. Intel Fortran Compiler 18.0 Developer Guide and Reference. Technical report, 2018. [Online]. Available: <https://software.intel.com/en-us/download/190-fortran-developer-guide-and-reference>. [Accessed: 05-Nov-2018]
- [8] L. S. Blackford, J. Demmel, J. Dongarra, I. Duff, S. Hammarling, G. Henry, M. Heroux, L. Kaufman, A. Lumsdaine, A. Petitet, R. Pozo, K. Remington, and R. C. Whaley. An updated set of basic linear algebra subprograms (blas). *ACM Transactions on Mathematical Softwares*, 28(2):135–151, June 2002.

- [9] C. L. Lawson, R. J. Hanson, D. Kincaid, and F. T. Krogh. Basic linear algebra subprograms for fortran usage. *ACM Transactions on Mathematical Softwares*, 5:308–323, 1979.