# Report on Implementation of Fast Multiplication Routines in Fortran

B.Sc. Adam Kosiorek, M.Sc. Nils Zander

*Abstract*—**This report describes the derivation of performance models and development of fast matrix multiplication routines in Fortran for Intel Xeon E5-2689 Sandy Bridge CPU for the ADHOC++ project.**

## I. INTRODUCTION

**T**HE majority of numerically intensive software projects can be reduced on the lowest level to linear algebra problems. It is hardly surprising that the main suite of tools for them - Basic Linear Algebra Subroutines (BLAS) - is amongst the most heavily optimized software packages out there. Using advanced data distribution and load balancing techniques, it leverages data locality and various levels of parallelism to boost performance. Many BLAS implementations are, however, ill-suited for small sized problems. As shown later, Intel MKL – arguably the fastest implementation for x86 CPUs [1] – is easily outperformed by the most naïve general matrix-matrix multiplication code when inputs are not big enough. In many cases it is therefore possible and necessary to devise custom matrix multiplication routines optimized for smaller problems. We develop performance model and optimized code for the following cases of matrix-matrix multiplication:

$$C = Ax + C \qquad (1)$$

where $C \in \mathcal{R}^{N \times K}$, $A \in \mathcal{R}^{N \times M}$, $x \in \mathcal{R}^{M \times K}$ and $N \in \mathcal{N}^{+}$, $K \in \{1, 2, 3\}$

## II. PERFORMANCE MODEL

Theoretical performance models, while not perfectly accurate, can hugely benefit development of optimized numerical programming routines by providing maximum performance estimates as a goal to strive for. Maximum theoretical performance can be computed by:

$$P_{peak} = IPC * FPI * Frequency \qquad (2)$$

where $IPC$ - instructions per cycle, $FPI$ - floating point operations per instruction, $Frequency$ - CPU's frequency. Intel Xeon E5-2689 operates with the frequency of 2.9 GHz and is capable of executing 1, 2 or 4 FPI, depending on the mode (Scalar, SSE, AVX) [need ref]. IPC depends on the numerical algorithm or the average number of floating point instructions executed per cycle to compute a single element of the output matrix. To derive IPC and $P_{peak}$, we refer to table I, which lists instructions that can be executed in one cycle by the Sandy Bridge CPU and to table II, which depicts number of each low level instructions needed to be executed to

compute one entry of $C$ matrix. Finally, table III summarizes results.

$$K = 1: \quad C(i,1) = C(i,1) + A(i,1) * x(1,1) \qquad (3)$$

$$K = 2: \quad \begin{aligned} C(i,1) = C(i,1) &+ A(i,1) * x(1,1) \\ &+ A(i,2) * x(2,1) \end{aligned} \qquad (4)$$

$$K = 3: \quad \begin{aligned} C(i,1) = C(i,1) &+ A(i,1) * x(1,1) \\ &+ A(i,2) * x(2,1) \qquad (5) \\ &+ A(i,3) * x(3,1) \end{aligned}$$

TABLE I
LOW LEVEL INSTRUCTIONS PER CYCLE.

| Instruction | Scalar | AVX |
|---|---|---|
| load | 1 or 2 | 1 |
| store | 1 or 0 | 0.5 |
| multiply | 1 | 1 |
| add | 1 | 1 |

TABLE II
LOW LEVEL INSTRUCTIONS REQUIRED TO COMPUTE ONE ENTRY OF $C$ MATRIX.

| K | load | store | multiply | add |
|---|---|---|---|---|
| 1 | 2 | 1 | 1 | 1 |
| 2 | 3 | 1 | 2 | 2 |
| 3 | 4 | 1 | 3 | 3 |

TABLE III
FLOATING POINT INSTRUCTIONS PER CYCLE AND PEAK THEORETICAL PERFORMANCE FOR SCALAR AND AVX MODES.

| K | IPC | | $P_{peak}$ [GFlops] | |
|---|---|---|---|---|
| | Scalar | AVX | Scalar | AVX |
| 1 | 1 | 1 | 2.9 | 11.60 |
| 2 | 2 | 1.33 | 5.8 | 15.47 |
| 3 | 2 | 1.5 | 5.8 | 17.40 |

## III. IMPLEMENTATION AND RESULTS

Fortran is the programming language traditionally used to implement high performance numerical algorithms. One reason is that it makes weaker assumptions about memory layout of variables than e.g. C, thus making vectorization easier. We decided to use Intel's Fortran compiler ifort, which often produces faster code than GNU's gfortran [2], but also

provides easier to read and more complete optimization feedback. Performance maximization requires code vectorization or execution in AVX mode. The following subsections describe constraints and counter-measures that enabled successful vectorization.

## A. Memory Layout

Fortran assumes column-major memory layout, that is, a single column of a matrix is comprised of elements adjacent in memory. In our case the routines are executed from C/C++, which assumes row-major memory layout. Normally, that would require declaration of transposed matrices in Fortran: e.g. double A[sizeA1][sizeA2] in C/C++ would become DOUBLE PRECISION A(sizeA2, sizeA1) in Fortran. Efficient multiplication, however, requires iteration along rows of the original matrices, which transfers to impossible to vectorize iteration along the innermost dimension in Fortran. Thus, we resort to computing indices manually and declaring matrices as one dimensional vectors i.e. DOUBLE PRECISION A(sizeA1 * sizeA2).

Another limiting factor is memory alignment. High performance AVX instructions can be used only when data are aligned at appropriate memory boundaries. Even if that is the case, the compiler cannot be sure if every possible input will be memory-aligned. One way to ensure it is to use inline compiler directives, in this case: !DEC$ VECTOR ALIGNED just before a loop referencing vectors.

## B. Data Dependency - Read after Write

Listing 1 contains the Read after Write dependency. It prevents vectorization, since the same element in memory is referenced multiple times. One possible solution is to introduce a local variable that would accumulate the results and store it in C(i) only once as in listing 2. Compilers than unrolls the outer loop containing code from listing 2 and vectorizes the result.

Listing 1. Read after Write dependency.
```
DO i = 1, 3
   C(i) = C(i) + A(index + i) * x(i)
END DO
```

Listing 2. Read after Write dependency resolved.
```
sum = 0
DO i = 1, 3
   sum = sum + A(index + i) * x(i)
END DO
C(i) = C(i) + sum
```

## C. Results

Achieved results, summarized in table IV, are not perfect and are true only for matrices small enough to fit in L2 cache. Divergence of theoretical and achieved results might stem from the fact that we did not consider memory bandwidth or cache effects.

TABLE IV
ACHIEVED RESULTS AND COMPARISON WITH THEORETICAL PEAK PERFORMANCE

| K | $P_{achieved}$ [Gflops] | $P_{achieved}/P_{peak}$ |
|---|---|---|
| 1 | 9.17 | 0.79 |
| 2 | 11.89 | 0.77 |
| 3 | 13.96 | 0.8 |

## REFERENCES

[1] Intel Math Kernel Library: https://software.intel.com/en-us/intel-mkl
[2] Intel Parallel Studio XE 2015 Product Brief: https://goo.gl/vc5eYo