

Implementing BLAS Operation using OpenMP

Amit Singh Cem Basoy

March 25, 2021

Contents

1	Introduction	2
1.1	Hand-tuned Assembly	2
1.2	Compiler Intrinsics	3
1.3	Compiler Dependent Optimization	3
1.4	BLAS Routines	4
1.5	Machine Model	4
1.6	Performance Metrics	4
1.6.1	FLOPS	4
1.6.2	Speedup	5
1.6.3	Speed-down	5
1.6.4	Peak Utilization	5
1.7	System Information	5
1.8	Compiler Information	6
1.9	library Version	6
2	Vector-Vector Inner Product	7
2.1	Calculating Number of Operations	7
2.2	Algorithm	8
2.2.1	The First Section	10
2.2.2	The Second Section	13
2.2.3	The Third Section	14
3	Outer Product	16
4	Matrix-Vector Product	17
5	Matrix-Matrix Product	18

Chapter 1

Introduction

Linear algebra is a vital tool in the toolbox for various applications, from solving a simple equation to the art of Deep Learning algorithm or Genomics. The impact can be felt across modern-day inventions or day to day life. Many tried to optimize these routines by hand-coding them in the assembly or the compiler intrinsics to squeeze every bit of performance out of the CPU; some chip manufacturers provide library specific to their chip. A few high-quality libraries, such as **Intel's MKL**, **OpenBLAS**, **Flame's Blis**, **Eigen**, and more, each one has one common problem, they are architecture-specific. They need to be hand-tuned for each different architecture.

There are three ways to implement the routines and each one has its shortcomings:

1. Hand code them in assembly and try to optimize them for each architecture.
2. Use the compiler intrinsics.
3. Simply code them and let the compiler optimize.

1.1 Hand-tuned Assembly

Hand-coded assembly gives high performance due to more control over registers, caches, and instructions used, but that comes with its issues, which we exchange for performance:

- Need a deep understanding of the architecture
- Maintenance of the code
- It violates the DRY principle because we need to implement the same algorithm for a different architecture
- Development time is high

- Debugging is hard
- Unreadable code
- Sometimes lead to micro-optimization or worst performance than the compiler

Intel's MKL, **OpenBLAS** and **Flame's Blis** comes under this category

1.2 Compiler Ininsics

The compiler intrinsics is one layer above the assembly, and we lose control over which register to use. If an intrinsic has multiple representations, then we do not know which instruction will emit. There are the following issues with this approach:

- Need the knowledge of intrinsic
- Maintenance of the code much better than assembly but not great
- DRY principle achieved with little abstraction
- Development time is better than assembly but not great
- Debugging is still hard
- Unreadable code if not careful
- May lead to micro-optimization or worst performance than the compiler

Eigen uses the compiler intrinsics, and they fixed and avoided some of the above problems with the right abstraction.

1.3 Compiler Dependent Optimization

The compiler has many tools for optimizing code: loop-unrolling, auto-vectorization, inlining functions, etc. We will rely on code vectorization heavily, but auto-vectorization may or may not be applied if the compiler can infer enough information from the code. To avoid unreliable auto-vectorization, we will use **OpenMP** for explicit vectorization, an open standard and supported by powerful compilers. The main issues are:

- Performance depends on the compiler
- No control over vector instructions or registers
- Auto-vectorization may fail

Boost.uBLAS depends on the auto-vectorization, which does not guarantee the code will vectorize.

1.4 BLAS Routines

There are four BLAS routines that we will implement using **OpenMP**, which uses explicit vectorization and parallelization using threads. Each routine has its chapter, and there we go much deeper with performance metrics.

1. Vector-Vector Inner Product (**?dot**)
2. Vector-Vector Outer Product (**?ger**)
3. Matrix-Vector Product or Vector-Matrix Product (**?gemv**)
4. Matrix-Matrix Product (**?gemm**)

1.5 Machine Model

The machine model that we will follow is similar to the model defined in the [Low et al., 2016], which takes modern hardware into mind. Such as vector registers and a memory hierarchy with multiple levels of set-associative data caches. However, we will ignore vector registers because we let the **OpenMP** handle the registers, and we do not have any control over them. However, we will add multiple cores where each core has at least one cache level that not shared among the other cores. The only parameter we need to put all our energy in is the cache hierarchy and how we can optimize the cache misses.

All the data caches are set-associative and we can characterize them based on the four parameter defined bellow:

- C_{L_i} : cache line of the i^{th} level
- W_{L_i} : associative degree of the i^{th} level
- N_{L_i} : Number of sets in the i^{th} level
- S_{L_i} : size of the i^{th} level in Bytes

$$S_{L_i} = C_{L_i} W_{L_i} N_{L_i} \quad (1.1)$$

Let the S_{data} be the width of the type in Bytes.

We are assuming that the cache replacement policy for all cache levels is **LRU**, which also assumed in the [Low et al., 2016] and the cache line is same for all the cache levels. For most of the case, we will try to avoid the associative so that we could derive a simple equation containing the cache size only from the equation 1.1.

1.6 Performance Metrics

1.6.1 FLOPS

It represents the number of floating-point operations that a processor can perform per second. The higher the Flops are, the faster it achieves the floating-point specific operations, but we should not depend on the flops all the time because it might be deceiving. Moreover, it does not paint the whole picture.

$$FLOPS = \frac{Number\ of\ Operation}{Time\ taken} \quad (1.2)$$

1.6.2 Speedup

The speedup tells us how much performance we were able to get when compared to the existing implementation. If it is more significant than one, then reference implementation performs better than the existing implementation; otherwise, if it is less than one, reference implementation performs worse than the existing implementation.

$$Speedup = \frac{Flops_{reference}}{Flops_{existing}} \quad (1.3)$$

1.6.3 Speed-down

The speed down is the inverse of the speedup, and if it is below one, then we performing better than the existing implementation; otherwise, we are performing worse.

$$Speeddown = \frac{Flops_{existing}}{Flops_{reference}} \quad (1.4)$$

1.6.4 Peak Utilization

This tells us how much CPU we are utilizing for floating-point operations when the CPU can compute X amount of floating-point operations.

$$Peak\ Utilization = \frac{Flops}{Peak\ Performance} \times 100 \quad (1.5)$$

1.7 System Information

Processor	2.3 GHz 8-Core Intel Core i9
Architecture	x86
L1 Cache	8-way, 32KiB
L2 Cache	4-way, 256KiB
L3 Cache	16-way, 16MiB
Cache Line	64B
Single-Precision(FP32)	32 FLOPs per cycle per core
Double-Precision(FP64)	16 FLOPs per cycle per core

Peak Utilization found using the formula defined on the [FLOPS, 2021]

$$Peak\ Utilization = Frequency \times Cores \times \frac{FLOPS}{Cycle} \quad (1.6)$$

Peak Utilization_{FP32} = 588.8 GFLOPS

Peak Utilization_{FP64} = 294.4 GFLOPS

1.8 Compiler Information

Compiler	Apple clang version 12.0.0 (clang-1200.0.32.29)
Compiler Flags	-march=native -ffast-math -Xclang -fopenmp -O3
C++ Standard	20

1.9 library Version

Intel MKL	2020.0.1
Eigen	3.3.9
OpenBLAS	0.3.13
BLIS	0.8.0

Chapter 2

Vector-Vector Inner Product

This operation takes the equal length of the sequence and gives the algebraic sum of the product of the corresponding entries.

Let x and y be the vectors of length n , then

$$x \cdot y = \sum_{i=0}^n (x_i \times y_i) \quad (2.1)$$

$$y \cdot x = x \cdot y \quad (2.2)$$

Equation 2.2 can be easily proven using equation 2.1.

$$y \cdot x = \sum_{i=0}^n (y_i \times x_i)$$

We know that multiplication on a scalar is commutative. Using this fact, we can say

$$y \cdot x = \sum_{i=0}^n (x_i \times y_i)$$

From the equation 2.1

$$y \cdot x = x \cdot y$$

2.1 Calculating Number of Operations

Using equation 2.1, we fill the below table

Name	Number
Multiplication	n
Addition	$n - 1$

Total Number of Operations = Number of Multiplication + Number of Addition

Total Number of Operations = $n + (n - 1)$

Total Number of Operations = $2n - 1$

2.2 Algorithm

Algorithm 1: Inner Product SIMD Function

```
// a is the pointer to the first vector
// b is the pointer to the second vector
// n is the length of the vectors
Function simd_loop (a, b, n):
    sum  $\leftarrow$  0
    #pragma omp simd reduction(+:sum)
    for i  $\leftarrow$  0 to n by 1 do
        | sum  $\leftarrow$  sum + a[i]  $\times$  b[i]
    end
    return sum
end
```

Algorithm 2: Dot Product

Input: $c, a, b, \text{max_threads}, n$
// a and b are pointer to the vectors
// c is the pointer to the output
// n is the size of the vectors
// max_threads is the user provided thread count
begin
 $\text{number_el_}L_1 \leftarrow \lfloor \frac{S_{L_1}}{S_{data}} \rfloor$
 $\text{number_el_}L_2 \leftarrow \lfloor \frac{S_{L_2}}{S_{data}} \rfloor$
 $\text{block}_1 \leftarrow 2 \times \text{number_el_}L_1$
 $\text{block}_2 \leftarrow \lfloor \frac{\text{number_el_}L_1}{2} \rfloor$
 $\text{block}_3 \leftarrow \lfloor \frac{\text{number_el_}L_2}{2} \rfloor$
 $\text{sum} \leftarrow 0$
 $\text{omp_set_num_threads}(\text{max_threads})$
 if $n < \text{block}_1$ **then**
 $\text{sum} \leftarrow \text{simd_loop}(a, b, n)$
 end
 else if $n < \text{block}_3$ **then**
 $\text{min_threads} \leftarrow \lfloor \frac{n}{\text{block}_2} \rfloor$
 $\text{num_threads} \leftarrow \max(1, \max(\text{min_threads}, \text{max_threads}))$
 $\text{omp_set_num_threads}(\text{num_threads})$
 **#pragma omp parallel for schedule(static) **
 reduction(+:sum)
 for $i \leftarrow 0$ **to** n **by** block_2 **do**
 $\text{ib} \leftarrow \min(\text{block}_2, n - i)$
 $\text{sum} \leftarrow \text{sum} + \text{simd_loop}(a + i, b + i, \text{ib})$
 end
 end
 else
 #pragma omp parallel reduction(+:sum)
 begin
 for $i \leftarrow 0$ **to** n **by** block_3 **do**
 $\text{ib} \leftarrow \min(\text{block}_3, n - i)$
 $\text{ai} \leftarrow a + i$
 $\text{bi} \leftarrow b + i$
 #pragma omp for schedule(dynamic)
 for $j \leftarrow 0$ **to** ib **by** block_2 **do**
 $\text{jb} \leftarrow \min(\text{block}_2, \text{ib} - j)$
 $\text{sum} \leftarrow \text{sum} + \text{simd_loop}(\text{ai} + j, \text{bi} + j, \text{jb})$
 end
 end
 end
 end
 $c \leftarrow \text{sum}$
end

The algorithm fragmented into three sections, and each part represents cache hierarchy.

2.2.1 The First Section

This section handles the data when the vector's length is less than the S_{L_1} . Here, we do not touch the threads because the overhead incurred is significant to paralysis the performance, and we get the worst performance compared with no thread or other libraries. When the whole data can fit in the L_1 cache, it is best not to use threads. We move to the next section when the cache misses large enough to compensate for the thread overhead.

According to the **Amdahl's Law**, we know

$$Speedup_N = \frac{1}{p + \frac{1-p}{N}} - O_N \quad (2.3)$$

p is the proportion of execution time for code that is not parallelizable, where $p \in [0, 1]$

$1 - p$ is the proportion of execution time for code that is parallelizable

N is the number of processor cores or number of threads, where $N \in \mathbb{Z}_{\neq 0}$

O_N is the overhead incurred due to the N thread

Using the equation 2.3 for $N = 1$, we know the $Speedup_1 = 1$ since the program is already running on the thread, and we do not need to spawn another thread. Therefore the thread overhead incurred is none ($O_1 = 0$).

In order to perform better than the single-threaded program, we need the $Speedup_N$ to exceed the $Speedup_1$ and $N > 1$. We can now get the lower bound using the fact and the relationship between the vectors' length and the thread overhead.

$$Speedup_N > Speedup_1$$

From the equation 2.3, we get

$$\frac{1}{p + \frac{1-p}{N}} - O_N > 1$$

$$\begin{aligned}
O_N &< \frac{1}{p + \frac{1-p}{N}} - 1 \\
O_N &< \frac{N}{Np + 1 - p} - 1 \\
O_N &< \frac{N}{1 + p(N-1)} - 1 \\
O_N &< \frac{(N-1) - p(N-1)}{1 + p(N-1)} \\
O_N &< \frac{(N-1)(1-p)}{1 + p(N-1)} \\
O_N &< \frac{1-p}{\frac{1}{(N-1)} + p}
\end{aligned}$$

Let $\frac{1}{N-1}$ be C_N and $C_N \leq 1$, where $N > 1$

$$O_N < \frac{1-p}{C_N + p} \quad (2.4)$$

We know $p \in [0, 1]$. Now, we get

$$\begin{aligned}
0 &\leq \frac{1-p}{C_N + p} \leq \frac{1}{C_N} \\
0 &\leq O_N < \frac{1}{C_N} \\
0 &\leq O_N < N-1
\end{aligned}$$

Therefore, $O_N \in [0, N-1)$. We can theorise that p must be a function of Block (B_1).

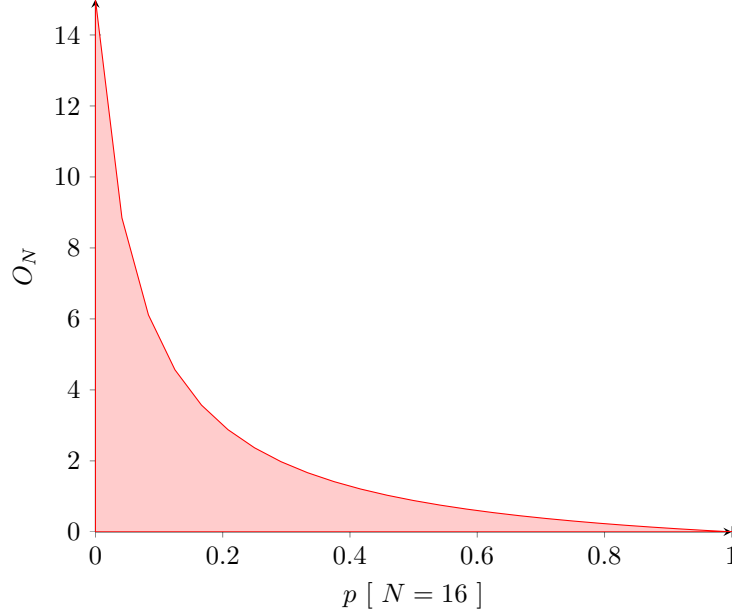
$$\begin{aligned}
p &\propto B_1 \\
p &= kB_1
\end{aligned}$$

where k is the proportionality constant. Now, replace p in the equation 2.4

$$O_N < \frac{1}{\frac{C_N + kB_1}{1 - kB_1}}$$

As we B_1 increase then $1 - kB_1$ decrease and $C_N + kB_1$ also increase, which over all decrease the overhead (O_N)

Thread Overhead



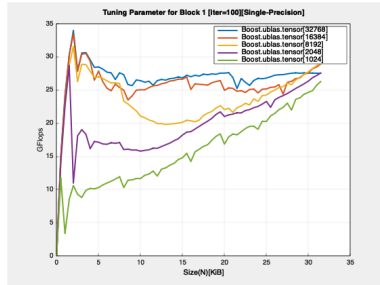
The B_1 's lower bound found to be the number of elements that can be fit inside the L_1 cache through the experimental data. This phenomenon happens because the vectors' elements cannot fit inside the cache, and cache misses increase, dominating the thread overhead.

$$B_1 \geq \frac{S_{L_1}}{S_{data}}$$

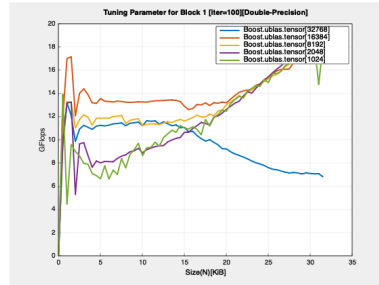
However, the optimal block size for this section found to be twice the vectors' elements that can put inside the cache.

$$B_1 = \frac{2S_{L_1}}{S_{data}} \quad (2.5)$$

Experimental Data for B_1



(a) Single-Precision



(b) Double-Precision

	Block[Single]	Block[Double]
Lower Bound	8192 (8K)	4096 (4K)
Optimal	16384 (16K)	8192 (8K)

2.2.2 The Second Section

Here, the cache misses is large enough to dominate the thread overhead, and the thread can provide speedup more significant than a single thread. To quench the data demand for each thread, we fit both vectors' elements inside the L_1 cache, and the accumulator always is inside the register—the whole L_1 cache used for vectors. The advantage of the threads come in handy because most of the CPU architecture provides a private L_1 cache, which owned by each core, and as the inner product is an independent operation, no two elements are dependent on each other for the result. Each private L_1 cache can fetch a different part of the sequence and run in parallel without waiting for the result to come from another thread, because of which there is no need to invalidate the cache. Once the data comes inside the L_1 cache, it will compute a fragment of the sequence and keeps the accumulated result in the register and then fetch another part. In the end, it will combine all the partial results into one final result.

$$S_{L_1} = S_{data}(\text{Block of the first vector} + \text{Block of the second vector})$$

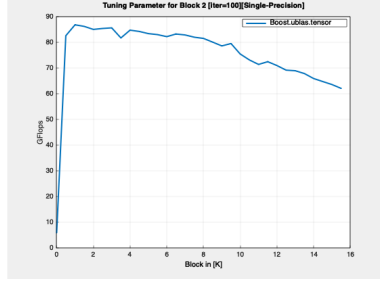
From the equation 2.1, we can infer that the block of both vectors must be equal to give the optimal block size.

$$B_2 = \text{Block of the first vector} = \text{Block of the second vector}$$

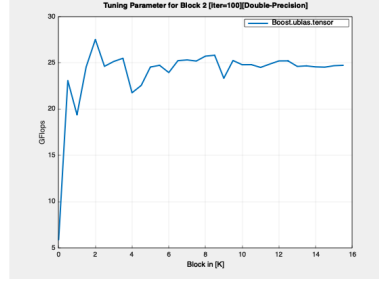
$$\begin{aligned} S_{L_1} &= S_{data}(B_2 + B_2) \\ S_{L_1} &= 2S_{data}B_2 \end{aligned}$$

$$B_2 = \frac{S_{L_1}}{2S_{data}} \quad (2.6)$$

Experimental Data for B_2



(a) Single-Precision



(b) Double-Precision

	Block[Single]	Block[Double]
Theoretical Optimal	4096 (4K)	2048 (2K)
Experimental Optimal	4096 (4K)	2048 (2K)

2.2.3 The Third Section

Here, once the data starts to exceed the L_2 cache or even L_3 cache, we use the CPU architecture's prefetch feature. This section may or may not affect some architecture; if the L_2 cache is also private for different cores, it will prefetch the data and fill it, and when it goes inside the loop, which utilizes the threads will fetch the data into the L1 cache from the L_2 cache. If the L_2 shared among all the cores, the gain might be minimal or see no gain in speedup. With the same logic we used to derive equation 2.6, we can use it here also.

$$S_{L_2} = S_{data}(\text{Block of the first vector} + \text{Block of the second vector})$$

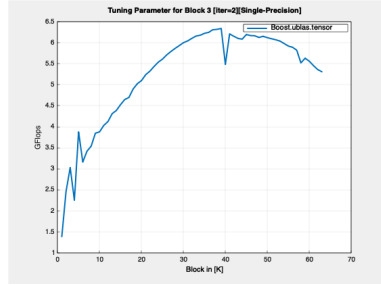
From the equation 2.1, we can infer that the block of both vectors must be equal to give the optimal block size.

$$B_3 = \text{Block of the first vector} = \text{Block of the second vector}$$

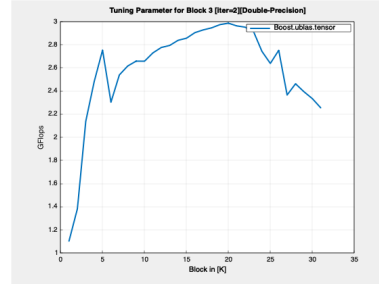
$$\begin{aligned} S_{L_2} &= S_{data}(B_3 + B_3) \\ S_{L_2} &= 2S_{data}B_3 \end{aligned}$$

$$B_3 = \frac{S_{L_2}}{2S_{data}} \quad (2.7)$$

Experimental Data for B_2



(a) Single-Precision



(b) Double-Precision

	Block[Single]	Block[Double]
Theoretical Optimal	32768 (32K)	16384 (16K)
Experimental Optimal	$\geq 30K$ and $\leq 40K$	$\geq 16K$ and $\leq 20K$

Chapter 3

Outer Product

Chapter 4

Matrix-Vector Product

Chapter 5

Matrix-Matrix Product

Bibliography

FLOPS. Flops — Wikipedia, the free encyclopedia. =
<https://en.wikipedia.org/wiki/FLOPS>,, 2021. [Online; accessed 24-March-2021].

Tze Meng Low, Francisco D. Igual, Tyler M. Smith, and Enrique S. Quintana-Ortí. Analytical modeling is enough for high-performance BLIS. *ACM Transactions on Mathematical Software*, 43(2):12:1–12:18, August 2016. URL <http://doi.acm.org/10.1145/2925987>.