

# Roofline and Matrix Multiplication PAPI Analysis

Brito Rui  
Department of Informatics  
University of Minho  
ruibrito666@gmail.com

Alves José  
Department of Informatics  
University of Minho  
zealves.080@gmail.com

## Abstract

This paper focuses on the comparative analysis, based on the Roofline model[1], of two computers and on matrix dot product algorithm performance improvements.

## 1 Roofline

### 1.1 Machine Profiles

The machines used for this study were an Apple MacBook Pro late 2008 and a 2010 HP dv6-2190ep. Information about the machines was gathered from various \*NIX tools (e.g. */proc/cpuinfo*, */proc/meminfo*, *dmidecode* and *sysctl*) and from the web (e.g. *Intel Ark* and *Crucial*). To calculate cache and main memory bandwidth we used the tool *bandwidth*<sup>1</sup>.

#### 1.1.1 Performance Peaks

In order to calculate the rooflines, we needed the Floating-Point(FP) Performance Peak and the Memory Bandwidth's Peak. To attain the FP Performance Peak we solve the following formula:

$$\text{GFlop/s}_{\max} = \#_{\text{cores}} \times f_{\text{clock}} \times \#_{\text{SIMD}}$$

MacBook Pro FP Performance Peak:

$$\text{GFlop/s}_{\max} = 2 \times 2.8 \times 8 = 44.8 \text{GFLOPSs}$$

HP Pavillion FP Performance Peak:

$$\text{GFlop/s}_{\max} = 4 \times 1.6 \times 8 = 51.2 \text{GFLOPSs}$$

To calculate the Memory Bandwidth Peak we solve the following formula:

$$\text{BW}_{\max} = \#_{\text{channels}} \times \text{mem}_{\text{clock}} \times \text{bus}_{\text{bandwidth}}$$

MacBook Pro Memory Bandwidth Peak:

$$\text{GFlop/s}_{\max} = 2 \times 1067 \times 64 = 17.072 \text{GBbyte}$$

HP Pavillion Memory Bandwidth Peak:

$$\text{GFlop/s}_{\max} = 2 \times 1333 \times 64 = 21.328 \text{GBbyte}$$

#### 1.1.2 Specifications

The specifications for the MacBook Pro are displayed on Table 1.

<b>Manufacturer:</b>	Apple
<b>Model:</b>	MacBook Pro late 2008
<b>Processor</b>	
Manufacturer:	Intel
Arch:	Core
Model:	Core 2 Duo T9600
Cores:	2
Clock Frequency:	2.80 GHz
FP Performance's Peak:	44.8 GFlops/s
<b>Cache</b>	
Level:	1
Size:	32KB + 32KB
Line Size:	64 B
Associative:	8-way
Memory Access Bandwidth:	40 GB/s
Level:	2
Size:	6 MB
Line Size:	64 B
Associative:	24-way
<b>RAM</b>	
Type:	SDRAM DDR3 PC3-8500
Frequency:	1067 MHz
Size:	4 GB
Num. Channels:	2
Latency:	13.13 ns

**Table 1:** MacBook Pro late 2008 specifications

<sup>1</sup><http://zsmith.co/bandwidth.html>

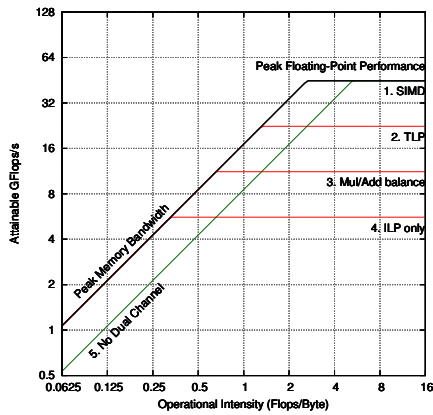
The specifications for the HP dv6-2190ep are displayed on Table 2.

<b>Manufacturer:</b>	HP
<b>Model:</b>	Pavillion dv6-2190ep
<b>Processor</b>	
Manufacturer:	Intel
Arch:	Nehalem
Model:	i7-720QM
Cores:	4
Clock Frequency:	1.60 GHz
FP Performance's Peak:	51.2 GFlops/s
<b>Cache</b>	
Level:	1
Size:	32KB + 32KB
Line Size:	64 B
Associative:	4/8-way
Memory Access Bandwidth:	22 GB/s
Level:	2
Size:	256 KB
Line Size:	64 B
Associative:	8-way
Level:	3
Size:	6 MB
Line Size:	64 B
Associative:	12-way
<b>RAM</b>	
Type:	SDRAM DDR3 PC3-10600
Frequency:	1333 MHz
Size:	4GB
Num. Channels:	2
Latency:	13.5 ns

**Table 2:** HP Pavillion dv6-2190ep specifications

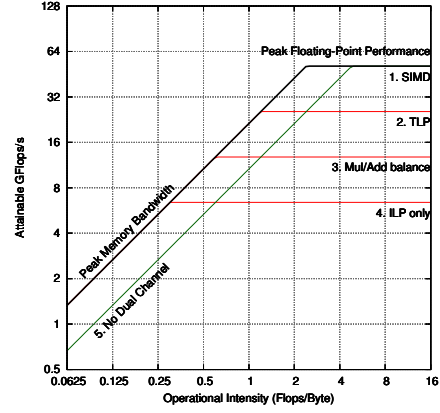
## 1.2 Roofline Model

### 1.2.1 MacBook Pro Roofline



**Figure 1:** Mackbook Pro late 2008

### 1.2.2 HP Pavillion Roofline



**Figure 2:** HP Pavillion dv6-2190ep Roofline

## 2 PAPI Case Study

### 2.1 Problem

The case study of this report, is to analyse the performance of a **matrix dot product** algorithm,

$$\text{Matrix}A * \text{Matrix}B = \text{Matrix}C \quad (1)$$

wich contains a triple nested loop with the indexes i,j and k(line,column and position). A naive implementation of this algorithm will provide, at best, very weak performance, since for every iteration of the j loop the whole current line of the matrix will be brought to cache (if it fits), however, since the j loop sweeps columns, this proves to be very inefficient. Our aim was to minimize this inefficiency.

### 2.2 Algorithm Analysis

The implementation produced to calculate the matrix multiplication was coded in C and compiled with Optimization level 3 (-O3, so the compiler can explore SIMD extensions). The naive algorithm of matrix multiplication is presented here, in order to better understand the problem at hand.

```
for (i = 0; i < size; i++) {
    for (j = 0; j < size; j++) {
        for(k = 0; k < size; k++) {
            acc += matrixA[i][k] * matrixB[k][j];
        }
        matrixC[i][j] = acc;
        acc = 0;
    }
}
```

And, for completeness' sake, here's the optimized version. Note that the matrix *tMatrix* is transposed so we can take advantage of a unit stride.

```

for (i = 0; i < size; i++) {
    for (j = 0; j < size; j++) {
        for(k = 0; k < size; k++) {
            acc += matrixA[i][k] *
                tMatrix[j][k];
        }
        matrixC[i][j] = acc;
        acc = 0;
    }
}

```

Two versions of the program were run, one with the original matrixB and other where matrixB is transposed. With this second version, it is expected for the algorithm to access a continuous memory space, thus leveraging unit stride, increasing overall performance while reducing memory accesses (due to the reduced miss rate).

## 2.3 Tests

### 2.3.1 Methodology

To measure the algorithm's performance, hardware counters were used. To gather information from these counters we used *PAPI*(Performance API). This tool allowed us to measure (among others) the following counters:

**PAPI\_TOT\_CYC** Total number of cycles;  
**PAPI\_TOT\_INS** Instructions completed;  
**PAPI\_LD\_INS** number of load instructions;  
**PAPI\_SR\_INS** number of store instructions;  
**PAPI\_FP\_OPS** Floating point operations;  
**PAPI\_FP\_INS** Floating point instructions;  
**PAPI\_L1\_DCA** L1 data cache accesses;  
**PAPI\_L1\_DCM** L1 data cache misses;  
**PAPI\_L2\_DCA** L2 data cache accesses;  
**PAPI\_L2\_DCM** L2 data cache misses;  
**PAPI\_L3\_DCA** L3 data cache accesses;

Moreover, all tests were run on a dedicated execution of the algorithm with process niceness set to -20 (*nice -n -20*) to ensure that essential machine time wasn't being spent on some other task. Also, to minimize overhead, OS Widgets and network were disabled. To decrease the chance of human error, the execution of all tests was "commanded" by a bash script. This script was also responsible for checking the 5% error margin, and, in case it wasn't satisfied, re-running the tests.

### 2.3.2 Test Cases

All four tests, presented below, were chosen to run in the two different version(normal and transpose). Each test was run four times, with the best execution time being selected within a margin of 5%.

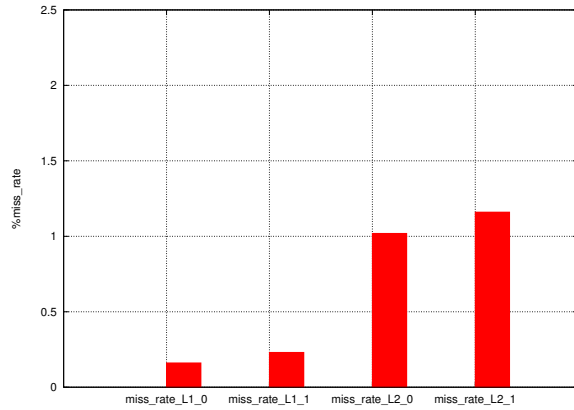
Memory	Size	Matrix Size
L1	30 KB	50
L2	255 KB	146
L3	3 MB	500
RAM	7.68 MB	800

**Table 3:** Test cases

## 2.4 Results

### 2.4.1 Analysis Chache

To measure the data cache misses the counters *PAPLL1\_DCM* and *PAPLL2\_DCM* were used. Each test fits in a different memory hierarchy level(L1, L2, L3 and RAM). (Note that test suffixed with a 0 mean unoptimized version, while those suffixed with a 1 mean, optimized version).



**Figure 3:** Percentage of Cache Misses

The above graphic shows an increase of percentage of misses with the optimized version, they are misleading. The next graph shows that although the percentage of misses increased, the total of misses didn't because the number of accesses also dropped.

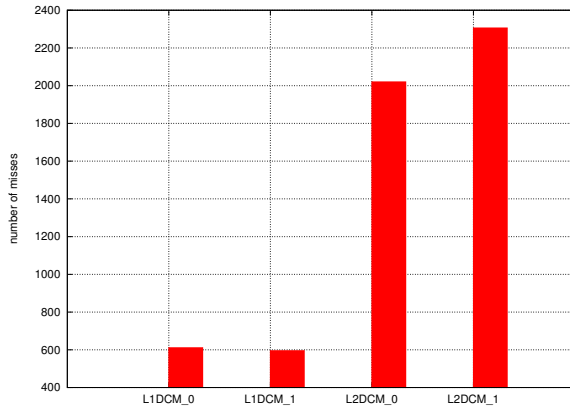
Usage of both levels of cache was estimated with specific counters. *PAPLL1\_DCA* and *PAPLL2\_DCA* provided the number of data accesses to the caches.

Before the results were out, it was expected a decrease of cache accesses from version one to version two of the algorithm.

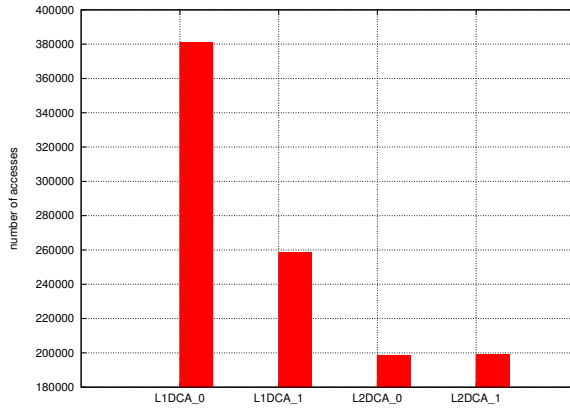
As we can see, the number of access to cache drops significantly from the first version to the second version while running with the L1 Cache Test. Though in the second test, the L2 Cache, the number of accesses slightly increased.

For more palpable values, the following table shows the execution times.

As it can be seen, the best improvements are seen in cache level 3 and in the main memory.



**Figure 4:** Number of Cache Misses



**Figure 5:** Number of Memory Accesses

Test	Time ( $\mu$ s)
L1_0	271.091
L1_1	254.000
L2_0	6983.450
L2_1	6629.000
L3_0	521429.0
L3_1	260535.0
RAM_0	6.60003e+06
RAM_1	3 1.08117e+06

**Table 4:** Execution Times

## References

- [1] **Roofline: An insightful Visual Performance Model for Floating-Point Programs and Multicore Architectures**  
*Samuel Webb Williams, Andre Waterman, David A. Patterson*  
23th November 2012
- [2] **Experiences and Lessons Learned with a Portable Interface to Hardware Performance Counters**  
*Jack Dongarra, Kevin London, Shirley Moore, Philip Mucci, Daniel Terpstra, Haihang You, Min Zhou*
- [3] **Computer Architecture: A Quantitative Approach, 5th Ed.**  
*John L. Hennessy, David A. Patterson*
- [4] <http://ark.intel.com/>
- [5] <http://www.crucial.com>
- [6] <http://www.wikipedia.com>

## Appendices

### A CUDA Implementation

Lorem ipsum dolor sit amet, consectetur adipisicing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.