



INSTITUTO SUPERIOR TÉCNICO

DEPARTAMENTO DE ENGENHARIA INFORMÁTICA

## ORGANIZAÇÃO DE COMPUTADORES

LEIC

### **First Lab Assignment: System Modeling and Profiling**

Version 1.1.0

2023/2024

# 1 Introduction

The goal of this assignment is twofold: (i) to determine the characteristics of a computer's caches, and (ii) to leverage the obtained knowledge about the caches in order to optimize the performance of a given program. For this task, the students will make use of a performance analysis tool to have direct access to hardware performance counters available on most modern microprocessors. The tool that will be used is the standard Application Programming Interface (API): PAPI [1].

In the rest of this section, we make a brief introduction to PAPI, and describe the targeted computer platform and the development environment. In Section 3, we describe the procedure for modeling the L1 and L2 caches of the targeted platform (Subsection 3.1), and provide a guide for analyzing the performance of a matrix-multiply code segment and optimizing it based on the characteristics of the L2 cache of the target architecture (Subsection 3.2).

## 1.1 Targeted Platform and Development Environment

**IMPORTANT: This assignment must be performed on the computers of your lab classes room.** These computers have similar hardware characteristics, and any of them can be used as a target platform. Note that, since this work is hardware-dependent, conducting it on a computer with different hardware characteristics could produce unexpected results, and hence invalidating your work. This means you should always use the same lab. If you are an Alameda student, you can access the specific lab computer you want (see <https://welcome.rnl.tecnico.ulisboa.pt/#labs-access>).

To properly setup the development environment, it is necessary to obtain the PAPI library and a set of auxiliary program files. This material can be found in the package `lab1_kit.zip`, which can be downloaded from the course website. After downloading and uncompressed this package on any of the lab classes' computers, PAPI must be built. To this end, change directories to the location of the PAPI source code: folder `papi-X.X.X/src`. Compile the code by issuing the commands: `./configure`, and `make`. This operation will produce a set of helper tools located in directory `src/utils/` and create the PAPI library `papilib.a`. The tool `papi_avail`, in particular, is useful to determine the PAPI events supported on the target platform. The library will be linked to the auxiliary programs presented in the following sections.

## 2 Exercise

To help determining the characteristics of the labs computer's caches, the following exercises will help you estimate cache parameters from small C applications.

The first step to get acquainted with the procedure is to determine only the size of the cache using a small C application on a (known) machine, such as the code you have analyzed on lab exercise VI.3. This C code, is a simplified version of the following programs in this assignment. Basically, it iterates over an array to determine the cache size.

To guarantee that you measure the time accurately, please use the source code available in the lab kit (file spark.c).

In order to perform the evaluation you should go to your lab in order to access the cache size by running the application there. You may want to repeat the evaluation of the elapsed time a few times to achieve statistical significance. You should table the relevant results for different cache sizes on the response sheet and make a conclusion regarding the cache size. You can calculate more measures before the output, examine the final part of the source code file.

1. What is the cache capacity of the computer you tested? Please justify.

To discover the other cache parameters, you're going to modify the C application, so that it generates different data access patterns. Please spend a few minutes analyzing the modifications to the source code.

```
for(size_t cache_size = CACHE_MIN; cache_size < CACHE_MAX; cache_size = 2*cache_size) {  
    for(size_t stride = 1; stride <= cache_size/2; stride = 2*stride){  
        limit = cache_size - stride + 1;  
        for(ssize_t i = 10 * stride; i > 0; i--) {  
            for(index = 0; index < limit; index += stride) {  
                array[index] = array[index] + 1;  
            }  
        }  
    }  
}
```

The meaning of each variable is the following:

**array[]** an arbitrary large array that will be repeatedly accessed to measure the cache miss pattern;

**cache\_size** value of the cache size under test; all cache sizes given by integer powers of 2, between CACHE\_MIN = 8kB and CACHE\_MAX = 64kB should be considered;

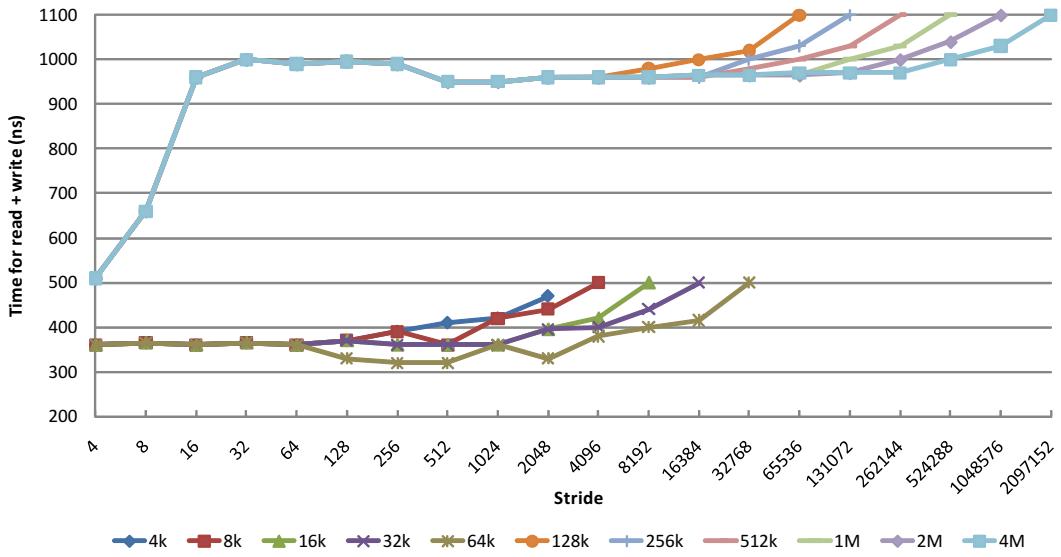
**stride** states how many entries are being skipped at each access; for example, if the stride is 4, entries 0, 4, 8, 12, ... in the array are being accessed, while entries 1, 2, 3, 5, 6, 7, 9, 10, 11, ... are skipped;

**limit** the largest address that will be accessed for the cache size and access pattern under test;

**repeat** denotes the number of times that each access pattern will be repeated in array.

The execution time for this code segment on this machine yield the chart depicted in Figure 1, by varying the adopted value for the *stride* parameter and for different array sizes, defined between ARRAY\_MIN = 4kB and ARRAY\_MAX = 4MB.

2. What is the cache capacity of the computer?
3. What is the size of each cache block?
4. What is the L1 cache miss penalty time?



**Figure 1:** Variation of the cache access time with the adopted *stride* value for different array sizes.

### 3 Procedure

#### 3.1 Modeling Computer Caches

In the first part of this assignment, the goal is to model the characteristics of the L1 data cache and L2 cache of the targeted computer platform. Next, we provide instructions for performing this analysis.

Use the forms at the end to answer the questions below.

##### 3.1.1 Modeling the L1 Data Cache

The methodology to experimentally model the L1 data cache consists in considering the total amount of data cache misses during the execution of the following code sequence of program `cml.c`, similar to the program in Section 2. This program can be found in the package `lab1_kit.zip`.

```

for(array_size=ARRAY_MIN; array_size < ARRAY_MAX; array_size=array_size*2)
    for(stride=1; stride <= array_size/2; stride=stride*2) {
        limit = array_size - stride + 1;
        for(repeat=0; repeat<=200*stride; repeat++)
            for(index=0; index<limit; index+=stride)
                x[index] = x[index] + 1;
    }
}

```

- a) Change to directory `cml/`, in the package `lab1_kit.zip`, and analyze de code of the program `cml.c`. Identify its source code with the program described above.  
What are the processor events that will be analyzed during its execution? Explain their meaning.
- b) Compile the program `cml.c` using the provided `Makefile` and execute `cml`. Plot the variation of the average number of misses (*Avg Misses*) with the `stride` size, for each considered dimension of the L1 data cache (8kB, 16kB, 32kB and 64kB).

NOTE: A fast sketch of these plots can be drawn in your computer by running the following commands:

```

./cml > cml.out
./cml_proc.sh

```

NOTE 2: You can draw these tables and plots on your computer, print, and attach to the report. You do not have to

fill them by hand on the printed report.

NOTE 3: You may need to mark the script as executable before being able to run it.

c) By analyzing the obtained results:

- Determine the **size** of the L1 data cache. Justify your answer.
- Determine the **block size** adopted in this cache. Justify your answer.
- Characterize the **associativity set size** adopted in this cache. Justify your answer.

### 3.1.2 Modeling the L2 Cache

In this part of the assignment, the goal is to experimentally model the characteristics of the L2 cache of the targeted computer platform. To analyze the computer's L2 cache, we will use the same methodology that was introduced in the previous section to model the L1 data cache.

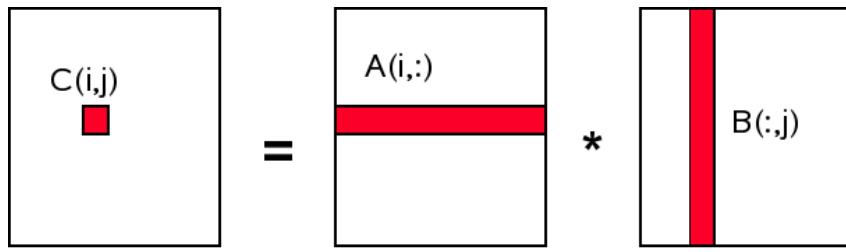
- a) Modify the program `cm1.c` in order to analyze the characteristics of the L2 cache. (Hint: use the event `PAPI_L2_DCM`.) Describe and justify the changes introduced in this program.
- b) Compile the program `cm1.c`, execute `cm1`, and plot the variation of the average number of misses (*Avg Misses*) with the `stride` size, for each considered dimension of the L2 cache.
- c) By analyzing the obtained results:
  - Determine the **size** of the L2 cache. Justify your answer.
  - Determine the **block size** adopted in this cache. Justify your answer.
  - Characterize the **associativity set size** adopted in this cache. Justify your answer.

## 3.2 Profiling and Optimizing Data Cache Accesses

Often, programmers wishing to improve their programs' performance focus their attention on how the programs affect the computer's caches. In the following, it will be analyzed how simple code changes can help to improve that performance for a matrix multiplication application.

Consider a simple matrix multiplication application, operating on two square matrices of  $N \times N$  16-bit integer elements, with  $N = 1024$ . From a mathematical point of view, given two matrices **A** and **B**, with elements  $a_{ij}$  and  $b_{ij}$  such that  $0 \leq i, j < N$ , the product matrix **C** is defined as:

$$c_{ij} = \sum_{k=0}^{N-1} a_{ik} b_{kj} = a_{i1} b_{1j} + a_{i2} b_{2j} + \dots + a_{i(N-1)} b_{(N-1)j} \quad (1)$$



**Figure 2:** Straightforward matrix multiplication.

### 3.2.1 Straightforward implementation

A straight-forward C implementation of Eq. 1 can look like this:

```
for (i = 0; i < N; ++i) {
    for (j = 0; j < N; ++j) {
        for (k = 0; k < N; ++k) {
            res[i][j] += mul1[i][k] * mul2[k][j];
        }
    }
}
```

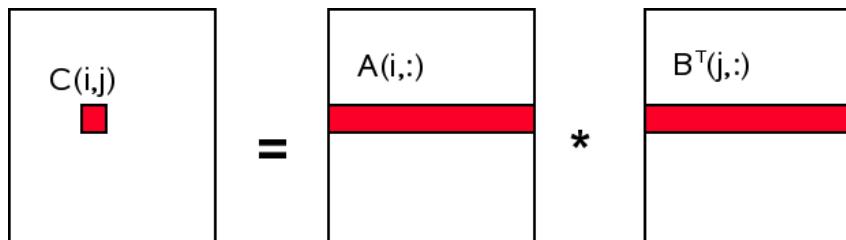
The two input matrices are `mul1` and `mul2`. The result matrix `res` is assumed to be initialized to all zeroes.

The provided program `mm1.c` includes this code sequence and all the necessary initialization steps, as well as the set of statements that are required in order to profile its execution using the PAPI toolbox.

- a) Change to directory `mm1/` and analyze de code of the program `mm1.c`. Identify its source code with the program described above.  
What is the total amount of memory that is required to accommodate each of these matrices?
- b) Compile the source file `mm1.c` using the provided `Makefile` and execute it. Fill the table with the obtained data.
- c) Evaluate the resulting L1 data cache *Hit-Rate*.

### 3.2.2 First Optimization: Matrix transpose before multiplication [2]

By analyzing the obtained results, it can be observed that such a straightforward implementation suffers from a severe penalty in what concerns the amount of L2 cache misses resulting from its access pattern. In fact, while `mul1` matrix is accessed sequentially, the inner loop advances the row number of `mul2` (see Fig. 2), meaning successive accesses to far away memory positions.



**Figure 3:** Transposed matrix multiplication.

One possible remedy to attenuate such problem is based on matrix transposition. In fact, since each matrix element is accessed multiple times, it might be worthwhile to rearrange (“transpose,” in mathematical terms) the second matrix `mul2` before using it (see Fig. 3):

$$c_{ij} = \sum_{k=0}^{N-1} a_{ik} b_{jk}^T = a_{i1} b_{j1}^T + a_{i2} b_{j2}^T + \dots + a_{i(N-1)} b_{j(N-1)}^T \quad (2)$$

After the preliminary transposition step, both matrices may be iterated sequentially. As far as the C code is concerned, it now looks like this:

```
int16_t tmp[N][N];

// transposition
for (i = 0; i < N; ++i) {
    for (j = 0; j < N; ++j)
        tmp[i][j] = mul2[j][i];
}

// multiplication
for (i = 0; i < N; ++i) {
    for (j = 0; j < N; ++j) {
        for (k = 0; k < N; ++k)
            res[i][j] += mul1[i][k] * tmp[j][k];
    }
}
```

Variable `tmp` is a temporary array to store the transposed matrix.

One direct consequence of this optimization is that it now requires additional accesses to the data memory. Hopefully, this extra cost can be easily recovered, since the 1024 non-sequential accesses per column are usually much more expensive.

- a) Change to directory `mm2/` and analyze the code of the program `mm2.c`. Identify its source code with the program described above. Compile this program using the provided `Makefile` and execute it.

Fill the table with the obtained data.

- b) Evaluate the resulting L1 data cache *Hit-Rate*.

- c) Change the code in the program `mm2.c` in order to include the matrix transposition in the execution time. Compile this program using the provided `Makefile` and execute it.

Fill the table with the obtained data.

Comment on the obtained results when including the matrix transposition in the execution time.

- d) Compare the obtained results with those that were obtained for the straightforward implementation, by calculating the difference of the resulting hit-rates ( $\Delta$ HitRate) and the obtained speedups.

### 3.2.3 Second Optimization: Blocked (tiled) matrix multiply [2]

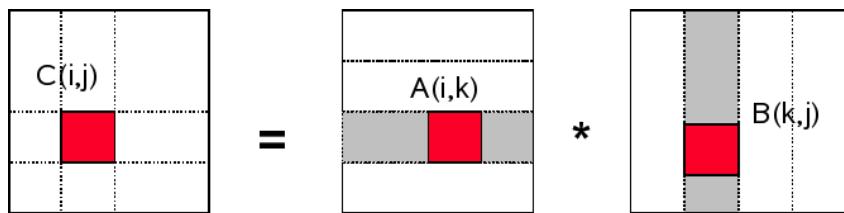
Despite the good results that may be obtained with the matrix transposition method, in many applications this approach can not be applied, either because the matrix is too large or the available memory is too small. Hence, other alternatives, which do not require the extra copy procedure, should be studied.

The search for an alternative processing scheme should start with a close examination of the involved math and the operations performed by the original implementation. Trivial math knowledge shows that the order of the several additions to obtain each element of the result matrix is irrelevant, as long as

each addend appears exactly once. This understanding will lead to solutions which reorder the additions performed in the inner loop of the original code.

According to the original algorithm, the adopted order to access the elements of matrix `mull` is: (0,0), (1,0), ... , (N -1,0), (0,1), (1,1), ... . Although the elements (0,0) and (0,1) are in the same cache line, by the time the inner loop completes one round, this cache line has long been evicted. For this example, each round of the inner loop requires, for each of the three matrices, 1024 cache lines, which is much more than what is available in most processors' caches.

One possible solution is to simultaneously handle more than one iteration of the middle loop, while executing the inner loop. In this case, several values which are guaranteed to be in cache will be used, thus contributing to a reduction of the L2 cache miss-rate. Hence, to maximize the speedup provided by this technique, it is necessary to adapt the dimension of the sub-matrix under processing to the cache block size, by taking into account the size of each matrix element. As a hypothetical example, considering that a `short` operand occupies 2-Bytes, this means that a 64-Byte cache block will accommodate 32 matrix elements, thus defining the optimal size for the sub-matrix line to be 32 (see Fig. 4).



**Figure 4:** Blocked matrix multiplication.

As far as the C code is concerned, it now looks like this:

```
#define SUB_MATRIX_SIZE (CACHE_LINE_SIZE / sizeof (short))

for (i = 0; i < N; i += SUB_MATRIX_SIZE) {
    for (j = 0; j < N; j += SUB_MATRIX_SIZE) {
        for (k = 0; k < N; k += SUB_MATRIX_SIZE) {
            for (i2 = 0, rres = &res[i][j], rmull1 = &mull[i][k];
                 i2 < SUB_MATRIX_SIZE;
                 ++i2, rres += N, rmull1 += N) {
                for (k2 = 0, rmull2 = &mull2[k][j]; k2 < SUB_MATRIX_SIZE; ++k2, rmull2 += N) {
                    for (j2 = 0; j2 < SUB_MATRIX_SIZE; ++j2) {
                        rres[j2] += rmull1[k2] * rmull2[j2];
                    }
                }
            }
        }
    }
}
```

The most visible change is that the code has six nested loops now. The outer loops iterate with intervals of `SUB_MATRIX_SIZE` (the cache line size `CACHE_LINE_SIZE` divided by `sizeof(short)`). This divides the matrix multiplication in several smaller problems which can be handled with more cache locality. The inner loops iterate over the missing indexes of the outer loops. There are, once again, three loops. The `k2` and `j2` loops are in a different order. This is done because, in the actual computation, only one expression depends on `k2` but two depend on `j2`.

- a) Change to directory `mm3/` and analyze the code of the program `mm3.c`. Identify its source code with the program described above.

Change the program source code in order to comply the algorithm parameterization (sub-matrix line size) with the block size (`CLS`) that was determined in Section 3.1.

How many matrix elements can be accommodated in each cache line?

- b) Compile this program using the provided `Makefile` and execute it. Fill the table with the obtained data.

- c) Evaluate the resulting L1 data cache *Hit-Rate*.
- d) Compare the obtained results with those that were obtained for the straightforward implementation, by calculating the difference of the resulting hit-rates ( $\Delta$ HitRate) and the obtained speedup.
- e) Compare the obtained results with those that were obtained for the matrix transpose implementation by calculating the difference of the resulting hit-rates ( $\Delta$ HitRate) and the obtained speedup. If the obtained speedup is positive, but the difference of the resulting hit-rates is negative, how do you explain the performance improvement? (Hint: study the hit-rates of the L2 cache for both implementations; You may use the following PAPI events PAPI\_L2\_DCH (or PAPI\_L2\_DCM) and PAPI\_L2\_DCA. Run papi\_avail to check for available events and understand their meaning.)

## References

- [1] Performance Application Programming Interface (PAPI). Webpage. "<http://icl.cs.utk.edu/papi>", December 2008.
- [2] Ulrich Drepper. What every programmer should know about memory. Technical report, Red Hat, Inc., November 2007.
- [3] *PAPI User's Guide*.
- [4] *PAPI Programmer's Reference*.

# First Lab Assignment: System Modeling and Profiling

STUDENTS IDENTIFICATION:

Number:	Name:
102695	Leonor Fortes
102703	Isabela Pereira
102823	Beatriz Raulo

## 2 Exercise

Please justify all your answers with values from the experiments.

- What is the cache capacity of the computer you used (please write the workstation name)?

Array Size	16 kB	32 kB	64 kB	128 kB	256 kB	512 kB
t2-t1	0.003439	0.006728	0.022584	0.060545	0.122237	0.159111
# accesses a[i]	320	640	1280	2560	5120	10240
# mean access time	$1.07 \times 10^5$	$1.05 \times 10^5$	$1.76 \times 10^5$	$2.37 \times 10^5$	$2.38 \times 10^5$	$1.55 \times 10^5$

Usamos o computador lab6pt1

A partir da tabela acima podemos ver que o tamanho da cache L1 é 32 kbytes devido a um subito aumento do tempo de acesso depois deste valor de cache, sendo que em 16 kbytes e 32 kbytes o tempo é praticamente constante.

Consider the data presented in Figure 1. Answer the following questions (2, 3, 4) about the machine used to generate that data.

- What is the cache capacity?

Consideramos que a capacidade da cache é de 64 kbytes devido ao subito aumento do tempo de leitura e escrita dos dados na cache, entre os tamanhos de 64 kbytes e 128 kbytes do array, devido a uma maior taxa de miss-rate pois excede o tamanho da cache tendo que despende mais tempo a ir buscar diretamente a memória. O que demonstra que arrays com tamanho menor ou igual a 64 kbytes não excedem o tamanho da cache, porém maior que esse valor tem uma miss-rate maior

- What is the size of each cache block?

O tamanho de cada bloco é de 16 bytes, pois o tamanho de um bloco pode ser determinado encontrando-se o primeiro stride onde o tempo de acesso, para o grupo de tamanho de array com tempos de acesso maiores, se estabiliza. Analizando o gráfico podemos entao concluir que o valor de stride em que estabiliza o tempo de acesso de arrays maiores, é efectivamente de 16 bytes.

- What is the L1 cache miss penalty time?

Poderemos calcular a miss penalty time do cache L1 comparando o tempo total de execução quando a miss rate é de aproximadamente 0% (hit-rate aproximadamente 100%, ou seja quando não excede o tamanho da cache,  $\leq 64\text{ kbytes}$ ) com quando é de aproximadamente 100% (quando excede totalmente o armazenamento no cache,  $\geq 128\text{ kbytes}$ ). Dito isto podemos, pelo gráfico, que os tempos respetivos são 300 ns < 1000 ns, por isso a miss penalty da cache L1 é de  $1000\text{ ns} - 300\text{ ns} = 690\text{ ns}$

### 3 Procedure

#### 3.1.1 Modeling the L1 Data Cache

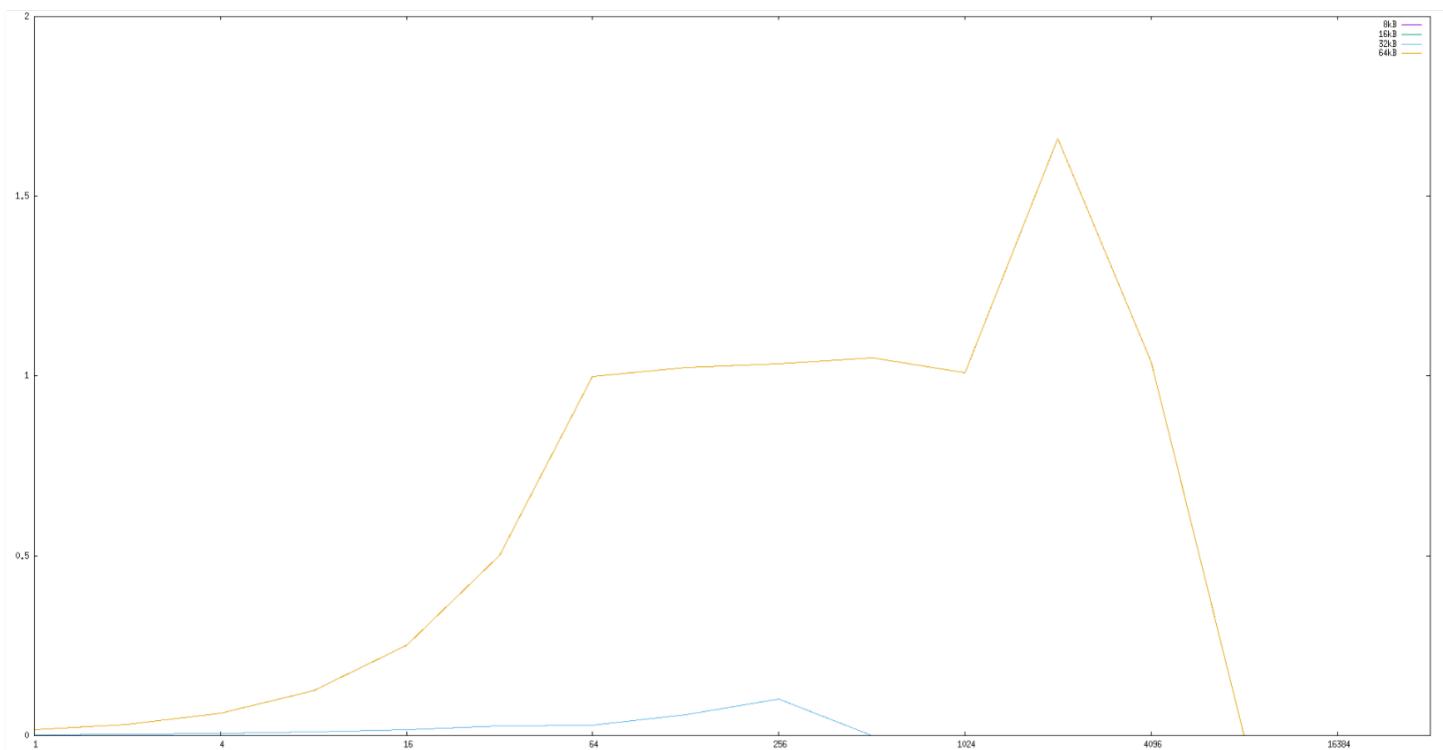
- a) What are the processor events that will be analyzed during its execution? Explain their meaning.

O evento adicionado que é processado no pedaço de código do ficheiro é o PAPI-L1-DCM que adiciona um evento que monitoriza os misses da cache L1. Ou seja, quando tentarmos dar fetch a dados que não estão na cache L1.

- b) Plot the variation of the average number of misses (*Avg Misses*) with the stride size, for each considered dimension of the L1 data cache (8kB, 16kB, 32kB and 64kB).

**Note that, you may fill these tables and graphics (as well as the following ones in this report) on your computer and submit the printed version.**

<b>Array Size</b>	<b>Stride</b>	<b>Avg Misses</b>	<b>Avg Cycl Time</b>	<b>Array Size</b>	<b>Stride</b>	<b>Avg Misses</b>	<b>Avg Cycl Time</b>
8kBytes	1	0.000174	0.002311	32kBytes	1	0.002070	0.002191
	2	0.000142	0.002318		2	0.003123	0.002193
	4	0.000155	0.002383		4	0.005137	0.002150
	8	0.000064	0.002214		8	0.008715	0.002134
	16	0.000043	0.002150		16	0.014548	0.002081
	32	0.000038	0.002152		32	0.037251	0.002012
	64	0.000051	0.002205		64	0.047779	0.002027
	128	0.000035	0.002056		128	0.048356	0.002133
	256	0.000021	0.001985		256	0.097663	0.002070
	512	0.000010	0.001937		512	0.000175	0.002049
	1024	0.000007	0.001925		1024	0.000069	0.001982
	2048	0.000007	0.001988		2048	0.000033	0.001996
	4096	0.000007	0.002028		4096	0.000016	0.002135
	8192	0.000030	0.002118		8192	0.000003	0.002121
	16384	0.000002	0.002034		16384	0.000002	0.002034
16kBytes	1	0.000180	0.002248	64kBytes	1	0.015639	0.001949
	2	0.000227	0.002242		2	0.031273	0.001849
	4	0.000178	0.002180		4	0.062539	0.002117
	8	0.000177	0.002244		8	0.125199	0.002179
	16	0.000193	0.002243		16	0.250361	0.002233
	32	0.000191	0.002196		32	0.500647	0.002030
	64	0.000188	0.002178		64	1.00036	0.001871
	128	0.000053	0.002060		128	1.027076	0.001835
	256	0.000019	0.001999		256	1.038051	0.001867
	512	0.000019	0.001985		512	1.053458	0.001937
	1024	0.000019	0.001985		1024	1.008145	0.001829
	2048	0.000011	0.002084		2048	1.901926	0.005348
	4096	0.000064	0.002167		4096	1.036992	0.005062
	8192	0.000004	0.002083		8192	0.000015	0.002159
	16384	0.000001	0.002119		16384	0.000001	0.002045
	32768	0.000001	0.002045		32768	0.000001	0.002045



c) By analyzing the obtained results:

- Determine the **size** of the L1 data cache. Justify your answer.

Consideramos que a capacidade da cache é de 32 Kibytes devido ao súbito aumento do tempo de leitura e escrita dos dados na cache, entre o tamanho de 32 Kibytes e 64 Kibytes do array, devido a uma maior taxa de miss-rate pois excede o tamanho da cache tendo que despende mais tempo a ir buscar diretamente à memória. O que demonstra que arrays com tamanho menor ou igual a 32 Kibytes não excedem o tamanho da cache, porém maior que esse valor tem uma miss-rate maior.

- Determine the **block size** adopted in this cache. Justify your answer.

O tamanho de cada bloco da cache é de 64 bytes, pois o tamanho de um bloco de cache pode ser determinado encontrando-se o primeiro stride entre o tempo de acesso, para o grupo de tamanho de array com tempos de aceso mais altos, se estabiliza. Analisando o gráfico podemos então concluir que o valor de stride em que estabiliza o tempo de acesso de arrays maiores, é efetivamente de 64bytes.

- Characterize the **associativity set size** adopted in this cache. Justify your answer.

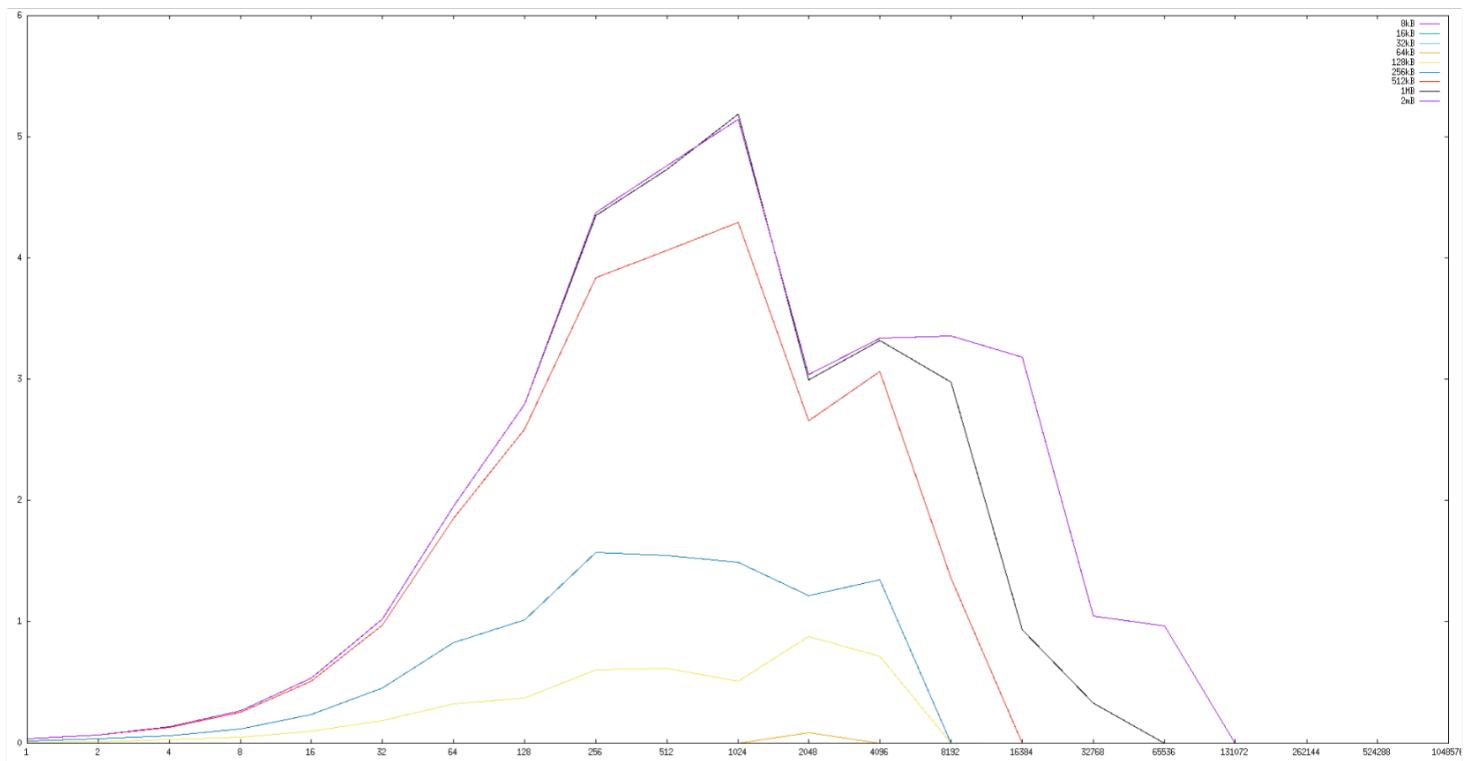
A associative set size é associada ao stride em que a miss rate diminui de volta a zero para a maior cache existente (neste caso, 64KiB). Como o stride determina quantos blocos distintos de dados do array iremos acessar, sabemos que para um stride de  $2^{15}$  (o máximo para este tamanho de array) apenas iremos acessar a 2 blocos diferentes desta cache. Se a cache for, pelo menos, associativa bidirecional, temos uma miss rate próxima de zero. Se repetirmos esse processo para strides mais baixas, notamos que para  $2^{14}$  e  $2^{13}$  a miss rate também é próxima a zero. Isso deve significar que a cache é, pelo menos, associativo de 4-way ou 8-way, respectivamente. Para um stride menor que  $2^{12}$ , a miss rate é maior do que 0, então devemos ter ultrapassando o número de vias da nossa cache. Daí concluimos que o associative set size é de 8.

### 3.1.2 Modeling the L2 Cache

- a) Describe and justify the changes introduced in this program.

As mudanças que efetuamos no código, foram a mudança do evento analizado, de PAPI\_L1-DCM para PAPI\_L2-DCM porque queríamos medir o miss rate da cache L2. Aumentamos também o CACHE\_MAX para 2 MiB devido ao aumento do tamanho de cache L2 relativamente à cache L1.

- b) Plot the variation of the average number of misses (*Avg Misses*) with the stride size, for each considered dimension of the L2 cache.



c) By analyzing the obtained results:

- Determine the **size** of the L2 cache. Justify your answer.

Consideramos que a capacidade da cache é de 256 Kibytes apesar de haver um "salto" entre os valores 128 Kibytes e 256 Kibytes, também houve um segundo "salto", mais saliente, entre os 256 Kibytes e os 512 Kibytes, dai termos considerado que o tamanho da cache é 256 Kibytes

- Determine the **block size** adopted in this cache. Justify your answer.

Consideramos que o block size da cache L2 é de 1024 bytes devido a ser o "pico" de dos arrays com maior size (2MB, 1MB, 512 Kibytes). Tendo a miso rate quase a 100% isto nos diz que o block está completamente preenchido com stride de 1024, dai concluirmos que o tamanho de cada bloco ser o mesmo.

- Characterize the **associativity set size** adopted in this cache. Justify your answer.

A associate set size é associada ao stride em que a miss rate diminui de volta a zero para a maior cache existente (neste caso, 2Mibytes). Como o stride determina quantos blocos distintos de dados do array iremos acessar, sabemos que para um stride de  $2^{20}$  (o máximo para este tamanho de array) operas iremos acessar a 2 blocos diferentes desta cache. Se a cache for, pelo menos, associativa bidirecional, temos uma miss rate próxima de zero. Se repetirmos esse processo para strides mais baixas, notamos que para  $2^{16}$  e  $2^{18}$  a miss rate também é próxima a zero. Isso deve significar que a cache é, pelo menos, associativo de 4-way ou 8-way, respectivamente. Para um stride menor que  $2^{17}$ , a miss rate é maior do que 0, porém para  $2^{17}$  ainda é 0, daí ter uma associatividade de 16-way.

## 3.2 Profiling and Optimizing Data Cache Accesses

### 3.2.1 Straightforward implementation

- a) What is the total amount of memory that is required to accommodate each of these matrices?

O tamanho total da matriz a guardar é  $512 \times 512$  e cada entrada da matriz é um número inteiro que ocupa 2 bytes de tamanho, dará o tamanho total de memória para guardar as matrizes é de  $512 \times 512 \times 2 = 5241280$  bytes = 512 kibytes.

- b) Fill the following table with the obtained data.

Total number of L1 data cache misses	$135.144818$	$\times 10^6$
Total number of load / store instructions completed	$536.871884$	$\times 10^6$
Total number of clock cycles	$579.192441$	$\times 10^6$
Elapsed time	$0.193065$	seconds

- c) Evaluate the resulting L1 data cache Hit-Rate:

$$\text{Hit-rate} = 1 - (\text{L1 data cache misses} / (\text{load/store instructions completed})) = 0.74827 = 74.83\%$$

Observámos 74.83% de hit rate logo as matrizes não cabem na totalidade na L1.

### 3.2.2 First Optimization: Matrix transpose before multiplication [2]

- a) Fill the following table with the obtained data.

Total number of L1 data cache misses	4.219230	$\times 10^6$
Total number of load / store instructions completed	536.871884	$\times 10^6$
Total number of clock cycles	525.633187	$\times 10^6$
Elapsed time	0.175212	seconds

- b) Evaluate the resulting L1 data cache *Hit-Rate*:

$$\text{Hit-rate} = 1 - \text{miss rate} = 1 - (\text{L1 data cache misses} / (\text{load/store instructions completed})) \\ = 0.9921 = 99.21\%$$

Observamos 99.21% de hit rate logo as matrizes não cabem na totalidade na L1, apesar do miss rate ser praticamente nulo

- c) Fill the following table with the obtained data.

Total number of L1 data cache misses	4.485384	$\times 10^6$
Total number of load / store instructions completed	537.396174	$\times 10^6$
Total number of clock cycles	527.401185	$\times 10^6$
Elapsed time	0.175801	seconds

Comment on the obtained results when including the matrix transposition in the execution time:

O tempo de execução aumentou, comparativamente à contagem do tempo após a transposição, como observado em a), uma vez que se realiza a transposição já em tempo de execução.

- d) Compare the obtained results with those that were obtained for the straightforward implementation, by calculating the difference of the resulting hit-rates ( $\Delta\text{HitRate}$ ) and the obtained speedups.

$\Delta\text{HitRate} = \text{HitRate}_{\text{mm2}} - \text{HitRate}_{\text{mm1}}: 0.9921 - 0.7483 = 0.2438$
$\text{Speedup}(\#Clocks) = \#\text{Clocks}_{\text{mm1}} / \#\text{Clocks}_{\text{mm2}}: 579.192441 / 525.633187 = 1.10189473444$
$\text{Speedup}(\text{Time}) = \text{Time}_{\text{mm1}} / \text{Time}_{\text{mm2}}: 0.193065 / 0.175212 = 1.10189370591$
Comment:
O speedup do método de transposta em relação ao método original foi de 1.1, tanto no #Clocks e no tempo de execução. Isto foi devido a um aproveitamento da localidade espacial da cache proporcionado pela transposição dos dados em linha de cache, com isto a multiplicação de matrizes é mais eficaz.

### 3.2.3 Second Optimization: Blocked (tiled) matrix multiply [2]

- a) How many matrix elements can be accommodated in each cache line?

$$\text{elements} = \text{block size} * \text{associative set size}/\text{element size} = 64 \text{ bytes} * 8 / 2 \text{ bytes} = 256$$

- b) Fill the following table with the obtained data.

Total number of L1 data cache misses	5.096860	$\times 10^6$
Total number of load / store instructions completed	536.888314	$\times 10^6$
Total number of clock cycles	215.610207	$\times 10^6$
Elapsed time	0.071871	seconds

- c) Evaluate the resulting L1 data cache *Hit-Rate*:

$$\text{Hit - Rate} = 1 - \text{miss rate} = 1 - (\text{L1 cache misses} / (\text{load/store instructions completed})) = 0.990506006159 = 99.05\%$$

Observamos 99.05% de hit rate logo as matrizes não cabem na totalidade na L1, apesar do miss rate ser praticamente nula.

- d) Compare the obtained results with those that were obtained for the straightforward implementation, by calculating the difference of the resulting hit-rates ( $\Delta\text{HitRate}$ ) and the obtained speedup.

$$\Delta\text{HitRate} = \text{HitRate}_{\text{mm3}} - \text{HitRate}_{\text{mm1}}: 0.99051 - 0.74827 = 0.24224 = 24.22\%$$

$$\text{Speedup}(\# \text{Clocks}) = \# \text{Clocks}_{\text{mm1}} / \# \text{Clocks}_{\text{mm3}}: 579.192441 / 215.610207 = 2.6862941651$$

Comment:

Comparando a implementação mm3 em relação à 3.2.1, mm1, observamos que houve uma melhoria na hit rate e no speedup do #Clock de 24.22 e 2.69, respectivamente. Isto deve-se à implementação por blocos (mm3) que nos permite aproveitar de forma mais eficiente o tamanho dos blocos da cache.

- e) Compare the obtained results with those that were obtained for the matrix transpose implementation by calculating the difference of the resulting hit-rates ( $\Delta\text{HitRate}$ ) and the obtained speedup. If the obtained speedup is positive, but the difference of the resulting hit-rates is negative, how do you explain the performance improvement? (Hint: study the hit-rates of the L2 cache for both implementations;)

$$\Delta \text{HitRate} = \text{HitRate}_{\text{mm3}} - \text{HitRate}_{\text{mm2}}: 0.99051 - 0.9921 = 0.00159 = -0.16\%.$$

$$\text{Speedup}(\# \text{Clocks}) = \# \text{Clocks}_{\text{mm2}} / \# \text{Clocks}_{\text{mm3}}: 525.633187 / 215.610207 = 2.43788637984$$

Comment:

Observámos que existe um speedup de 2.44 do mm3 em relação ao mm2, apesar da hit-rate do mm3 ser menos que a do mm2, apesar da hit-rate do mm3 ser menos que a do mm2. Isto deve-se ao facto de termos observado que a hit rate da cache L2 no mm3 é bastante superior ao hit rate dessa mesma cache no mm2, logo compensa o facto de a cache L1 ser pior, pois leva a uma melhor performance no over all.

### 3.2.3 Comparing results against the CPU specifications

Now that you have characterized the cache on your lab computer, you are going to compare it against the manufacturer's specification. For this you can check the device's datasheet, or make use of the command `lscpu`. Comment the results.

Usando o comando `lscpu -C` conseguimos concluir que os valores que obtivemos para a cache L1 (tamanho da cache, do bloco e a associatividade) estão de acordo com as especificações do processador. O tamanho da cache e do bloco da cache L2 também vão de encontro a esses dados, porém os gráficos que obtivemos não demonstram a associatividade real, de valor 4, mas sim a que comentamos anteriormente.

## A PAPI - Performance Application Programming Interface

The PAPI project [1] specifies a standard Application Programming Interface (API) for accessing hardware performance counters available in most modern microprocessors. These counters exist as a small set of registers that count *Events*, defined as occurrences of specific signals related to the processor's function (such as cache misses and floating point operations), while the program executes on the processor. Monitoring these events may have a variety of uses in the performance analysis and tuning of an application, since it facilitates the correlation between the source/object code structure and the efficiency of the actual mapping of such code to the underlying architecture. Besides performance analysis, and hand tuning, this information may also be used in compiler optimization, debugging, benchmarking, monitoring and performance modeling.

PAPI has been implemented on a number of different platforms, including: Alpha; MIPS R10K and R12K; AMD Athlon and Opteron; Intel Pentium II, Pentium III, Pentium M, Pentium IV, Itanium 1 and Itanium 2; IBM Power 3, 4 and 5; Cell; Sun UltraSparc I, II and II, etc.

Although each processor has a number of events that are native to that specific architecture, PAPI provides a software abstraction of these architecture-dependent *Native Events* into a collection of *Preset Events*, also known as *predefined events*, that define a common set of events deemed relevant and useful for application performance tuning. These events are typically found in many CPUs that provide performance counters. They give access to the memory hierarchy, cache coherence protocol events, cycle and instruction counts, functional unit, and pipeline status. Hence, preset events may be regarded as mappings from symbolic names (PAPI preset name) to machine specific definitions (native countable events) for a particular hardware resource. For example, Total Cycles (in user mode) is mapped into PAPI\_TOT\_CYC. Some presets are derived from the underlying hardware metrics. For example, Total L1 Cache Misses (PAPI\_L1\_TCM) is the sum of L1 Data Misses and L1 Instruction Misses on a given platform. The list of preset and native events that are available on a specific platform can be obtained by running the commands `papi_avail` and `papi_native_avail`, both provided by the papi source distribution.

Besides the standard set of events for application performance tuning, the PAPI specification also includes both a high-level and a low-level sets of routines for accessing the counters. The high level interface consists of eight functions that make it easy to get started with PAPI, by simply providing the ability to start, stop, and read sets of events. This interface is intended for the acquisition of simple but accurate measurement by application engineers [3, 4]:

- `PAPI_num_counters` – get the number of hardware counters available on the system;
- `PAPI_flops` – simplified call to get Mflops/s (floating point operation rate), real and processor time;
- `PAPI_ipc` – gets instructions per cycle, real and processor time;
- `PAPI_accum_counters` – add current counts to array and reset counters;
- `PAPI_read_counters` – copy current counts to array and reset counters;
- `PAPI_start_counters` – start counting hardware events;
- `PAPI_stop_counters` – stop counters and return current counts.

The following is a simple code example of using the high-level API [3, 4]:

```

#include <papi.h>

#define NUM_FLOPS 10000
#define NUM_EVENTS 1

int main() {
    int Events[NUM_EVENTS] = {PAPI_TOT_INS};
    long_long values[NUM_EVENTS];

    /* Start counting events */
    if (PAPI_start_counters(Events, NUM_EVENTS) != PAPI_OK)
        handle_error(1);

    do_some_work();

    /* Read the counters */
    if (PAPI_read_counters(values, NUM_EVENTS) != PAPI_OK)
        handle_error(1);

    printf("After reading the counters: %lld\n", values[0]);

    do_some_work();

    /* Add the counters */
    if (PAPI_accum_counters(values, NUM_EVENTS) != PAPI_OK)
        handle_error(1);

    printf("After adding the counters: %lld\n", values[0]);

    do_some_work();

    /* Stop counting events */
    if (PAPI_stop_counters(values, NUM_EVENTS) != PAPI_OK)
        handle_error(1);

    printf("After stopping the counters: %lld\n", values[0]);
}

```

#### Possible output:

```

After reading the counters: 441027
After adding the counters: 891959
After stopping the counters: 443994

```

The fully programmable low-level interface provides more sophisticated options for controlling the counters, such as setting thresholds for interrupt on overflow, as well as access to all native counting modes and events. Such interface is intended for third-party tool writers or users with more sophisticated needs.

The PAPI specification also provides access to the most accurate timers available on the platform in use. These timers can be used to obtain both real and virtual time on each supported platform: the real time clock runs all the time (e.g., a wall clock), while the virtual time clock runs only when the processor is running in user mode.

In the following code example, `PAPI_get_real_cyc()` and `PAPI_get_real_usec()` are used to obtain the real time it takes to create an event set in clock cycles and in microseconds [3, 4]:

```

#include <papi.h>

int main(){
    long long start_cycles, end_cycles, start_usec, end_usec;
    int EventSet = PAPI_NULL;

    if (PAPI_library_init(PAPI_VER_CURRENT) != PAPI_VER_CURRENT)
        exit(1);

    /*Create an EventSet */
    if (PAPI_create_eventset(&EventSet) != PAPI_OK)
        exit(1);

    /* Gets the starting time in clock cycles */
    start_cycles = PAPI_get_real_cyc();

    /* Gets the starting time in microseconds */
    start_usec = PAPI_get_real_usec();

    do_some_work();

    /* Gets the ending time in clock cycles */
    end_cycles = PAPI_get_real_cyc();

    /* Gets the ending time in microseconds */
    end_usec = PAPI_get_real_usec();

    printf("Wall clock cycles: %lld\n", end_cycles - start_cycles);
    printf("Wall clock time in microseconds: %lld\n", end_usec - start_usec);
}

```

**Possible output:**

```

Wall clock cycles: 100173
Wall clock time in microseconds: 136

```