

# ASIGNATURA Computación de altas prestaciones

## Práctica1 Vector processing and SIMD

You must write a report answering the questions proposed in each exercise, plus the requested files. Submit a zip file through Moodle. Check submission date in Moodle (deadline is until 11:59 pm of that date).

- Exercise 1:
  - o Identify your CPU model and list the supported SIMD instructions.

### Intel(R) Core(TM) i7-7500U CPU @ 2.70GHz

**flags:** fpu vme de pse tsc msr pae mce cx8 apic sep mtrr pge mca cmov pat pse36 clflush dts acpi mmx fxsr sse sse2 ss ht tm pbe syscall nx pdpe1gb rdtscp lm constant\_tsc art arch\_perfmon pebs bts rep\_good nopl xtopology nonstop\_tsc cpuid aperfmperf pni pclmulqdq dtes64 monitor ds\_cpl vmx est tm2 ssse3 sdbg fma cx16 xtpr pdcm pcid sse4\_1 sse4\_2 x2apic movbe popcnt tsc\_deadline\_timer aes xsave avx f16c rdrand lahf\_lm abm 3dnowprefetch cpuid\_fault epb invpcid\_single pti ssbd ibrs ibpb stibp tpr\_shadow vnmi flexpriority ept vpid ept\_ad fsgsbase tsc\_adjust bmi1 avx2 smep bmi2 erms invpcid mpx rdseed adx smap clflushopt intel\_pt xsaveopt xsavec xgetbv1 xsaves dtherm ida arat pln pts hwp hwp\_notify hwp\_act\_window hwp\_epp md\_clear flush\_l1d arch\_capabilities

- o Explain the main differences between both assembly codes (vectorized and non-vectorized) focused on the SIMD instructions generated by the compiler.

› diff --color simple2\_o3.s simple2\_o3\_native.s

### Bloque 1:

9a10

```
> xorl    %eax, %eax
```

12,17d12

```
< movdqa    .LC0(%rip), %xmm2
< movdqa    .LC1(%rip), %xmm4
< movdqa    .LC2(%rip), %xmm3
< movq     %rcx, %rax
< movq     %rdx, %rsi
< leaq     16384(%rcx), %rdi
```

### Bloque 2:

19,34c14,22

```
< movdqa    %xmm2, %xmm0
< addq    $32, %rax
< padd    %xmm4, %xmm2
< addq    $32, %rsi
< cvtdq2pd    %xmm0, %xmm1
< movaps %xmm1, -32(%rax)
< pshufd $238, %xmm0, %xmm1
< padd    %xmm3, %xmm0
< cvtdq2pd    %xmm1, %xmm1
< movaps %xmm1, -16(%rax)
< cvtdq2pd    %xmm0, %xmm1
< pshufd $238, %xmm0, %xmm0
< cvtdq2pd    %xmm0, %xmm0
< movaps %xmm1, -32(%rsi)
< movaps %xmm0, -16(%rsi)
< cmpq    %rdi, %rax
```

---

```
> pxor    %xmm0, %xmm0
> leal    1(%rax), %esi
> cvtsi2sdl    %eax, %xmm0
> movsd    %xmm0, (%rcx,%rax,8)
> pxor    %xmm0, %xmm0
> cvtsi2sdl    %esi, %xmm0
> movsd    %xmm0, (%rdx,%rax,8)
> addq    $1, %rax
> cmpq    $2048, %rax
```

### Bloque 3:

37c25

```
< movapd    .LC3(%rip), %xmm3
```

---

```
> movsd    .LC0(%rip), %xmm2
```

### Bloque 4:

44,49c32,35

```
< movapd    (%rdx,%rax), %xmm0
< mulpd    %xmm3, %xmm0
< addpd    (%rcx,%rax), %xmm0
< addq    $16, %rax
< addsd    %xmm0, %xmm1
< unpckhpd    %xmm0, %xmm0
```

---

## Computación de altas prestaciones HPC

```
> movsd (%rdx,%rax), %xmm0
> mulsd %xmm2, %xmm0
> addsd (%rcx,%rax), %xmm0
> addq $8, %rax
```

### Bloque 5:

67,68c53,54

```
< .section .rodata.cst16,"aM",@progbits,16
< .align 16
---
> .section .rodata.cst8,"aM",@progbits,8
> .align 8
```

### Bloque 6:

70,89d55

```
< .long 0
< .long 1
< .long 2
< .long 3
< .align 16
< .LC1:
< .long 4
< .long 4
< .long 4
< .long 4
< .align 16
< .LC2:
< .long 1
< .long 1
< .long 1
< .long 1
< .align 16
< .LC3:
< .long -611603343
< .long 1072693352
```

- Exercise 2:
  - o Provide the source code of *simple2\_intrinsics.c* after the vectorization of the loops. Explain how you have carried out the vectorization of the code.

```
> ./simple
4194513.817600
> ./intrinsics
4194513.817600
```

El primer bucle, queda de la siguiente forma:

```
1  __m256d vb = {0, 1, 2, 3};
2  __m256d va = {1, 2, 3, 4};
3
4  __m256d cons = {4, 4, 4, 4};
5
6  /* Populate A and B arrays */
7  for (i = 0; i < ARRAY_SIZE; i += 4)
8  {
9      _mm256_store_pd(&b[i], vb);
10     vb = _mm256_add_pd(vb, cons);
11
12     _mm256_store_pd(&a[i], va);
13     va = _mm256_add_pd(va, cons);
14
15 }
```

El segundo bloque queda de la siguiente forma:

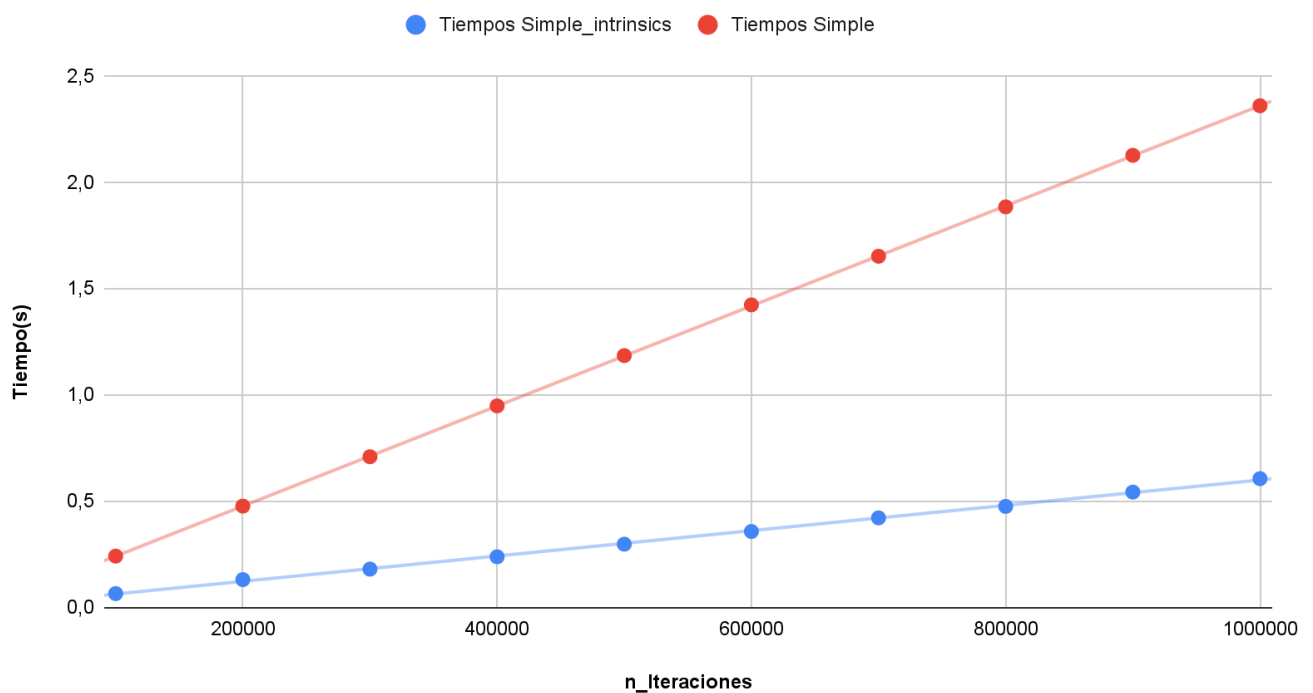
```
1  __m256d mm = {1.0001, 1.0001, 1.0001, 1.0001};
2  __m256d sum = {0.0, 0.0, 0.0, 0.0}; // to hold partial sums
3
4  /* Perform an operation a number of times */
5  for (t=0; t < NUMBER_OF_TRIALS; t++) {
6      for (i=0; i < ARRAY_SIZE; i += 4) {
7
8          // Load arrays
9          __m256d va = _mm256_load_pd(&a[i]);
10         __m256d vb = _mm256_load_pd(&b[i]);
11         // Compute m*a+b
12         __m256d tmp = _mm256_fmadd_pd(mm, va, vb);
13         // Accumulate results
14         sum = _mm256_add_pd(tmp, sum);
15
16         // c += m*a[i] + b[i];
17     }
18 }
19
20 // Get sum[2], sum[3]
21 __m128d xmm = _mm256_extractf128_pd(sum, 1);
22 // Extend to 256 bits: sum[2], sum[3], 0, 0
23 __m256d ymm = _mm256_castpd128_pd256(xmm);
24 // Perform sum[0]+sum[1], sum[2]+sum[3], sum[2]+sum[3], 0+0
25 sum = _mm256_hadd_pd(sum, ymm);
26 // Perform sum[0]+sum[1]+sum[2]+sum[3]...
27 sum = _mm256_hadd_pd(sum, sum);
28 c = sum[0];
```

- o Compare the execution time for different values of NUMBER\_OF\_TRIALS: from 100.000 to 1.000.000 in steps of 100.000. Plot the results in a graph. Discuss the results.

Tiempos Simple_intrinsics	Tiempos Simple
0,06320	0,239581
0,128949	0,474872
0,178317	0,707135
0,236446	0,946552
0,295623	1,183117
0,355387	1,421803
0,419435	1,650559
0,473498	1,882969
0,540748	2,125369
0,604040	2,358790

Tal y como observamos, los tiempos con intrínsecos son claramente menores.

### Tiempos Simple\_intrinsics y Tiempos Simple



- Exercise 3:
  - o The program includes two loops. The first loop (indicated as Loop 0) iterates over the arguments applying the algorithm to each of them. The second loop (indicated as Loop 1) computes the grayscale algorithm. Is this loop optimal to be vectorized? Why?

Consideramos que no es óptimo vectorizar el segundo bucle, ya que el segundo trata pixel a pixel las imágenes y tiene un tamaño demasiado grande como para que podamos vectorizarlo de manera óptima.

- o Provide the source code of the auto-vectorized version of the code. Explain the changes in the code to help the compiler to vectorize the loop.

Debido a que el orden de acceso a los datos es importante para la velocidad a la hora de recorrer matrices, hemos cambiado el orden de acceso para que recorra las columnas primero para maximizar la eficiencia de acceso a la memoria. Antes estaba al revés, con el acceso a filas primero, y gracias a eso podemos observar como el algoritmo mejora su tiempo de ejecución:

```
> ./greyScale ../images/8k.jpg
[info] Processing ../images/8k.jpg
[info] ../images/8k: width=7680, height=4320, nchannels=3
Tiempo: 0.956727
> make
gcc -O3 -march=native -fwhole-program -Wall -D_GNU_SOURCE -s
rm -f *.o *~
> ./greyScale ../images/8k.jpg
[info] Processing ../images/8k.jpg
[info] ../images/8k: width=7680, height=4320, nchannels=3
Tiempo: 0.319321
```

- o Provide the source code after manually vectorizing the code. Explain your solution.

```
1  for(int i = 0, j = 0; j < imageSize; i += 16, j += 4) {
2
3      __m128i* data_ptr_one = (__m128i*)(rgb_image + i);
4      __m128i* data_ptr_two = (__m128i*)(rgb_image + i + 8);
5      __m128i filas = _mm_loadl_epi64(data_ptr_one);
6      __m128i columnas = _mm_loadl_epi64(data_ptr_two);
7
8      // Extendemos los vectores y los convertimos a floats
9      __m256i extendedFilasInt = _mm256_cvtepu8_epi32(filas);
10     __m256 extendedFilasFloat = _mm256_cvtepi32_ps(extendedFilasInt);
11     __m256i extendedColumnasInt = _mm256_cvtepu8_epi32(columnas);
12     __m256 extendedColumnasFloat = _mm256_cvtepi32_ps(extendedColumnasInt);
13
14     // generamos el vector de coeficientes y lo usamos para multiplicar los vectores por pares de pixeles.
15     __m256 coeficientes = _mm256_set_ps(0.0, 0.1140, 0.5870, 0.2989, 0.0, 0.1140, 0.5870, 0.2989);
16     __m256 par1 = _mm256_mul_ps(extendedFilasFloat, coeficientes);
17     __m256 par2 = _mm256_mul_ps(extendedColumnasFloat, coeficientes);
18
19     // hacemos un horizontal add, dos veces porque no se completa en una sola.
20     __m256 h_add = _mm256_hadd_ps(par1, par2);
21     h_add = _mm256_hadd_ps(h_add, h_add);
22
23     // permutamos el vector final y lo extraemos para el outcome.
24     __m256 permutado = _mm256_permutevar8x32_ps(h_add, _mm256_set_epi32(0, 0, 0, 0, 5, 1, 4, 0));
25     __m128 outcome = _mm256_extractf128_ps(permutado, 0);
26
27     for (int k = 0; k < 4; k++) {
28         grey_image[j + k] = (int) outcome[k];
29     }
30 }
```

Tal y como hemos indicado en el código con los comentarios, separamos los pixeles de dos en dos, extendemos los vectores y los convertimos a floats, después generamos el vector de coeficientes y lo usamos para multiplicar los vectores por pares de pixeles, siendo par1, el pixel 1 y 2, y par2 el 3 y 4. Después haremos horizontal add, dos veces debido a que no acababa de hacerlo de una sola, y permutamos el vector y lo extraemos para sacarlo en la imagen.



- o Fill in a table with time and speedup results compared to the original version and auto-vectorized version for images of different resolutions (SD, HD, FHD, UHD-4k, UHD-8k). You must include a column with the fps at which the program would process. Discuss the results.

	greyscale auto (s)	greyscale manual (s)
8k	0,385949	0,356616
4k	0,124572	0,108843
FHD	0,055421	0,040849
HD	0,015936	0,016629
SD	0,005566	0,005998
	greyscale auto (FPS)	greyscale manual (FPS)
8k	3	3
4k	8	9
FHD	18	24
HD	63	60
SD	180	167

Debido a el tiempo que tardan en procesar la imagen, podemos observar cómo sería inviable tener este proceso para tamaños de 4k y 8k, con el manual podríamos acercarnos a una experiencia más cinematográfica con 24 fps, y si quisiéramos mantener ese procesamiento para ordenadores o móviles nos iríamos a un mínimo de HD para los 60 fps que vendrían bien debido a la habitual velocidad de refresco de 60hz.