

# Organización del Computador II

Departamento de Computación  
Facultad de Ciencias Exactas y Naturales  
Universidad de Buenos Aires

## Trabajo Práctico Número 2

Procesamiento de imágenes

Integrante	LU	Correo electrónico
Alcain Pablo	485/07	pabloalcain@gmail.com
Gamarra Ignacio	792/11	gamarra_mi32@yahoo.com.ar
Santos Diego	874/03	diego.h.santos@gmail.com

### Grupo

My Generation

# Índice

<b>1. Introducción</b>	<b>2</b>
<b>2. Desarrollo</b>	<b>4</b>
2.1. Diferencia de Imágenes . . . . .	4
2.1.1. Iteración del Ciclo Principal . . . . .	5
2.1.2. Sección .data . . . . .	5
2.1.3. Carga de datos . . . . .	7
2.1.4. Ciclo Principal . . . . .	7
2.2. Blur Gaussiano . . . . .	7
<b>3. Resultados</b>	<b>11</b>
3.1. SpeedUp . . . . .	11
3.2. Comparaciones entre C y ASM: Dependencia en tamaño . . . . .	11
3.2.1. Diferencia de Imágenes . . . . .	12
3.2.2. Blur Gaussiano . . . . .	12
3.3. Comparaciones entre C y ASM: Dependencia en aspect ratio . . . . .	14
3.4. Comparaciones entre C y ASM: Dependencia con el radio de blur . . . . .	14
3.5. Comparaciones entre C -O0 y C -O3: Filtro de diferencias . . . . .	14
3.6. Algunas conclusiones . . . . .	14
<b>4. Conclusión</b>	<b>19</b>

## 1. Introducción

En este trabajo hemos realizado la implementación de funciones de procesamiento de imágenes en dos lenguajes de programación, *Assembler* y *C*. Además realizamos un análisis de la performance de las mismas a fin de determinar cual es mas eficiente respecto al tiempo de ejecución. Las funciones implementadas son:

- **Diferencia de Imágenes**, a partir de dos imágenes genera una tercera imagen que resalta los pixeles donde las dos imagenes fuentes difieren. Para obtener la tercer imagen se toma la norma infinito la norma infinito de la resta vectorial entre los píxel, ignorando el canal alfa. La fórmula a aplicar para obtener la tercer imagen es;

$$\forall k \in \{r, g, b\} \quad O[i, j, k] = \|I_1[i, j] - I_2[i, j]\|_{\infty}$$

$$O[i, j, \alpha] = 255$$

- **Blur Gaussiano**, tomando una imagen de fuente genera una segunda similar a la fuente pero con un aspecto desenfocado. La manera de lograr este efecto es calculando cada componente de la imagen de salida como un promedio ponderado con una gaussiana 2D de los píxeles que la circundan. Las formulas a aplicar para realizar el filtro son las siguientes:

$$O[i, j, k] = (I * K)[i, j, k] = \sum_{x=-r}^r \sum_{y=-r}^r I[i + x, j + y, x] K[r - x, r - y]$$

Donde la matriz  $K$  está generada por los parámetros  $\sigma$  y  $r$ :

$$K_{\sigma,r}[i, j] = G_{\sigma}(r - i, r - j) \quad G_{\sigma}(x, y) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}}$$

En el caso de las implementaciones hechas en *Assembler* hemos usado el modelo de programación *SIMD*, a través del set de instrucciones *SSE*, pues el objetivo de este trabajo práctico es estudiar las ventajas y desventajas de usar ese modelo contra uno *SISD*.

En las siguientes secciones se explicará como se implementaron los filtros, se presentarán gráficos mostrando los tiempos de ejecución de cada filtro comparando ambas implementaciones y luego se presentará una conclusión acerca de los resultados obtenidos y el costo de usar un modelo *SIMD* para programar.

## 2. Desarrollo

### 2.1. Diferencia de Imagenes

Hicimos el filtro `diff` tanto en C como en ASM. Su implementación es bastante inmediata, al menos sin tener en cuenta los registros **XMM** de **SSE**. De esta manera, el filtro `diff` recorre todos los píxeles de ambas imágenes, de forma secuencial, resta los tres componentes de colores y toma el valor absoluto. Así, en cada píxel  $(i, j)$  [de 4 bytes cada uno], obtenemos tres bytes  $(\Delta R, \Delta G, \Delta B)$ . En este filtro `diff`, convenimos que la imagen diferencia fuera la norma infinito del vector en las tres posiciones (es decir, el máximo entre los tres). Así, una vez obtenidos  $(\Delta R, \Delta G, \Delta B)$ , calculamos el máximo entre los tres,  $\Delta C$ , y en el píxel  $(i, j)$  de la imagen de salida anotamos los tres bytes  $(\Delta C, \Delta C, \Delta C)$  y 0xFF en el canal alpha. Repetimos este procedimiento para todos los píxeles de la imagen.

En la implementación de ASM, aprovechamos los registros **SSE** disponibles en el procesador<sup>1</sup>. En cada registro entran 16 bytes: 4 píxeles que pueden procesarse a la vez. Ahora procesarlos no es tan trivial como antes, pero podemos ver que pensar algunas relaciones básicas de álgebra, ayuda a ver el camino a seguir. Recordemos los dos pasos:

**Restar los tres componentes y tomar valor absoluto** La clave para reducir este problema es considerando la siguiente igualdad:

$$\text{abs}(a - b) = \text{máx}(a, b) - \text{mín}(a, b) \quad (1)$$

Así, simplemente tenemos que asegurarnos de estar restando el máximo (byte a byte) menos el mínimo. Para eso, supongamos que tenemos las tiras de 16 bytes de ambas imágenes en **XMM0** y **XMM3**. Ahora, copiamos cualquiera de los dos (por caso, **XMM0**) a un registro extra **XMM1**, y guardamos **máx** (**XMM0**, **XMM3**) en **XMM0** [`pmaxub`]. Luego, guardamos **mín** (**XMM1**, **XMM3**) en **XMM1** [`pminub`]. Ahora, **XMM0** tiene los máximos y **XMM1** los mínimos, byte a byte, de los 4 píxeles. Los restamos byte a byte y guardamos el resultado en **XMM0** [`psubb`]. Hacemos aquí unas cuentas extras, ya que también estamos procesando el canal alpha. Esto no significa un detrimento en performance (ya que la operación está paralelizada), sino que si la información estuviera guardada en RGB sin canal alpha se podrían procesar más píxeles a la vez<sup>2</sup>.

**Obtener el máximo** Ahora el desafío es encontrar, con operaciones **SSE**, una forma de encontrar el máximo de cada bloque de 4 bytes. Pensemos cada bloque de 4 bytes por separado primero. Tenemos tres valores  $R, G, B$  y queremos que esto se transforme en  $C, C, C, X$  [ $C$  es el  $\text{máx}(R, G, B)$  y  $X$  es el canal alpha, que por ahora puede tomar cualquier valor]. Si  $\mathbf{t}$  es el vector de salida, queremos

$$\mathbf{t} = (\text{máx}(R, G, B), \text{máx}(R, G, B), \text{máx}(R, G, B)) \quad (2)$$

Para vectorizar esta expresión, tenemos que usar la siguiente identidad:

$$\text{máx}(R, G, B) = \text{máx}(R, \text{máx}(G, B)) = \text{máx}(G, \text{máx}(B, R)) = \text{máx}(B, \text{máx}(R, G)) \quad (3)$$

Definimos al vector  $\mathbf{p} = (R, G, B)$  y al operador `rot` tal que:

<sup>1</sup>El procesador utilizado, intel i7, tiene disponible **AVX**, pero no utilizamos estas instrucciones

<sup>2</sup>Claro que ahora no podríamos garantizarnos esta hermosura de que entren *justo* 4 píxeles en un registro **SSE**

$$\text{rot}(R, G, B) = (G, B, R) \quad (4)$$

Las sucesivas aplicaciones de **rot** las definimos también:

$$\mathbf{q} = \text{rot}(R, G, B) = (G, B, R) \quad (5)$$

$$\mathbf{r} = \text{rot}(\text{rot}(R, G, B)) = (B, R, G) \quad (6)$$

Con esta notación, la identidad 3 queda

$$\text{máx}(R, G, B) = \text{máx}(p_1 \text{máx}(q_1, r_1)) = \text{máx}(p_2, \text{máx}(q_2, r_2)) = \text{máx}(p_3, \text{máx}(q_3, r_3)) \quad (7)$$

y, en consecuencia, **t** es

$$\mathbf{t} = \text{máx}(\mathbf{p}, \text{máx}(\mathbf{q}, \mathbf{r})) \quad (8)$$

Esta ecuación ahora está vectorizada y podemos usar **pmaxub** para calcular los máximos. Sólo queda crear, a partir de **p** los vectores rotados **q** y **r**; es decir, implementar el operador **rot**. Para esto, utilizamos una máscara que rota a la izquierda los tres bytes, es decir:  $(R, G, B) \rightarrow (G, B, R)$ , utilizando **pshufb** y una máscara del tipo (1, 2, 0).

Como, en rigor, en cada píxel tenemos también el canal alpha, la máscara para cada byte va a ser (1, 2, 0, 0xFF), poniendo un 0 en el canal alpha [esto también podría ser la posición 3 o cualquier otra, porque el canal alpha *tiene que* ser reescrito].

Esto explica el procedimiento general. Aumentarlo para procesar  $n$  píxeles en paralelo es directo, simplemente la máscara de rotación va a ser **m** = (1, 2, 0, 0xFF, 4 + 1, 4 + 2, 4 + 0, 0xFF, 8 + 1, 8 + 2, 8 + 0, 0xFF, ...). En general,

$$m_i = \begin{cases} 4 * \text{floor}(i/4) + \text{mod}(i + 1, 3) & \text{mod}(i, 4) \neq 3 \\ 0xFF & \text{mod}(i, 4) = 3 \end{cases} \quad (9)$$

La implementación entonces es directamente una traducción de esto: Copiamos el registro **XMM0** (**p**) a un registro **XMM1** y mantenemos la máscara de rotación en **XMM2**. Ahora rotamos **XMM1** con **pshufb** y obtenemos **q**. Obtenemos el máximo de **XMM0** (**p**) y **XMM1** (**q**) con **pmaxub** y lo guardamos en **XMM0**. Rotamos nuevamente **XMM1** y resulta **r**. Ahora calculamos el máximo de **XMM0** (**máx(p, q)**) y **XMM1** (**r**) y lo guardamos en **XMM0** (el objetivo final, **t**).

En la figura 1 se puede observar un esquema de este procedimiento.

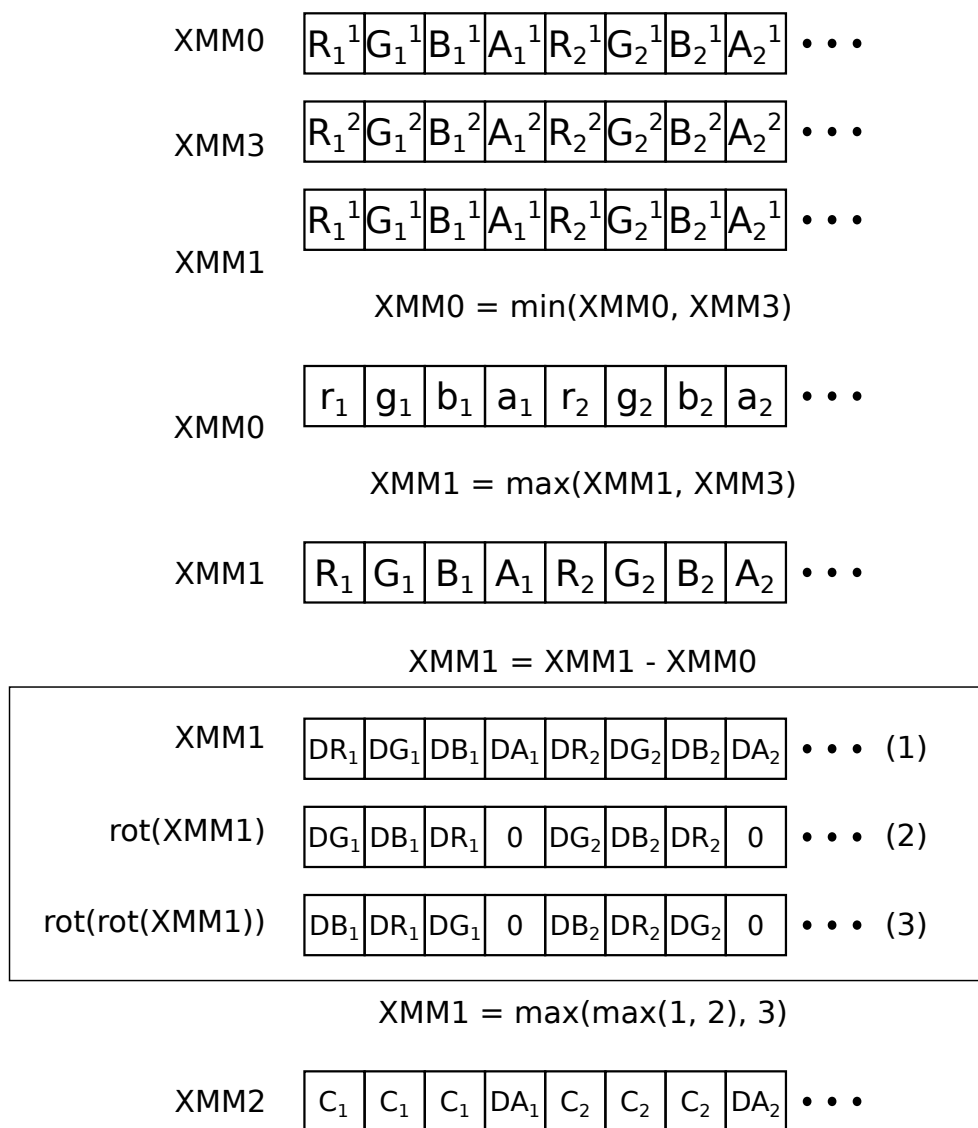
### 2.1.1. Iteración del Ciclo Principal

Vamos a mostrar como se cargan y modifican los registros en una iteración de nuestra implementación

### 2.1.2. Sección .data

Definimos una constantes en la sección .data (por ende es una constantes inicializada)

- **mask'rot**: DB 1, 2, 0, 0xFF, 5, 6, 4, 0xFF, 9, 10, 8, 0xFF, 13, 14, 12, 0xFF : es una máscara que utilizaremos para acomodar los canales



**Figura 1:** Esquema que caracteriza el algoritmo de diff con instrucciones SIMD

### 2.1.3. Carga de datos

La carga de datos inicial para el ciclo principal del algoritmo es:

- Copiamos en `XMM2` la máscara `mask_rot`
- Inicializamos `R9` en 0 con la instrucción `XOR` para utilizarlo de contador de columnas
- Inicializamos `R10` en 0 con la instrucción `XOR` para utilizarlo de contador de filas
- Preservamos en `R8` `m` y lo usamos para saber cuándo terminamos una columna
- Preservamos en `RCX` `n` y lo usamos para saber cuándo terminamos una fila

### 2.1.4. Ciclo Principal

Con los datos cargados en los registros vamos a mostrar una iteración del ciclo principal del filtro. Para cada línea de la matriz que representa la imagen fuente se realizan las siguientes acciones.

- Leemos desde la primer imagen fuente en `RDI` 16 bytes y los cargamos en `XMM0`, cargando aquí la fila  $i - 1$  con la instrucción `MOVDQU`
- Leemos desde la segunda imagen fuente en `RSI` 16 bytes y los cargamos en `XMM3`, cargando aquí la fila  $i - 1$  con la instrucción `MOVDQU`
- Obteniendo en cada uno de los registros

	$R3$	$G3$	$B3$	$A3$	$R2$	$G2$	$B2$	$A2$
<sup>127</sup>	$R1$	$G1$	$B1$	$A1$	$R0$	$G0$	$B0$	$A0$
	<sub>0</sub>							

## 2.2. Blur Gaussiano

El desarrollo del filtro de `blur` también lo realizamos en C y en ASM. Este filtro está caracterizado por el radio  $r$  y por el parámetro de difusión  $\sigma$ . Vamos a utilizar los siguientes términos

**Marco:** El recuadro de tamaño  $r$  que bordea a la figura, donde los píxeles no se afectan.

**Centro:** La parte de la figura que no es el marco.

**Filtro:** La matriz con la que se convoluciona cada píxel.

**Vecinos:** Los píxeles cercanos que afectan el valor de salida de un píxel en la imagen procesado. Es decir, todos los elementos que están a una distancia (en norma 1) menor a  $r$  de un píxel dado.

En este filtro, tenemos que aplicar la siguiente relación que mencionamos en la introducción:

$$O[i, j, k] = \sum_{x=-r}^{y=r} I[i + x, j + y, k] K_{\sigma, r}[r - x, r - y] \quad (10)$$

para todos los píxeles en el *centro*. Así, la implementación en C es, para cada componente  $(R, G, B)$  de cada píxel  $(i, j)$  del *centro*, hacer un ciclo en  $(x, y)$  entre  $(-r, -r)$  y  $(r, r)$ . Un detalle es la identificación del vecino  $(x, y)$  del píxel  $(i, j)$  con la posición absoluta en la imagen,  $(i+x, j+y) = (i+x)*n_{\text{cols}} + (j+y)$ . Así, con los valores de  $(R, G, B)$  del píxel  $(i, j)$  de la imagen salida inicializados a 0, acumulamos en cada uno  $I[i+x, j+y]K_{\sigma,r}[r-x, r-y]$ . Luego, redondeamos este resultado (que es de punto flotante de 4 bytes) a un char y lo guardamos en la imagen de salida.

Unas consideraciones hay que realizar en el cálculo del filtro:

1. Normalización: a pesar de que la función gaussiana está normalizada a 1, no sucede en este caso (ya que está discretizada y además truncada). De este modo, el coeficiente de normalización del filtro tiene que ser calculado para cada uno, de modo tal que la suma de todos los elementos de la matriz sea 1.
2. Cálculo: para calcular los elementos de una gaussiana hay que realizar muchos cálculos de punto flotante. Para evitarnos repetir estos cálculos una y otra vez, y como el filtro es siempre el mismo, lo implementamos como una *look up table* y lo guardamos en la memoria<sup>3</sup>. La creación de la LUT es en C [de cualquier modo, esta cuenta se realiza *una sola vez por imagen*, así que no es costoso su cálculo].

La implementación en ASM fue realizada nuevamente utilizando las instrucciones SSE. Para eso, tenemos que vectorizar el filtro. Llamando  $\mathbf{T}$  a la matriz de píxeles de salida (es decir,  $\mathbf{T}[i, j] = \mathbf{t}_{i,j} = (R_{i,j}^s, G_{i,j}^s, B_{i,j}^s)$ ) y  $\mathbf{P}$  a la matriz de píxeles de entrada ( $\mathbf{P}[i, j] = \mathbf{p}_{i,j} = (R_{i,j}^0, G_{i,j}^0, B_{i,j}^0)$ ) la ecuación 10 anterior queda

$$\mathbf{T}[i, j] = \sum_{x=-r}^{y=r} \mathbf{P}[i+x, j+y] K_{\sigma,r}[r-x, r-y] \quad (11)$$

$$\mathbf{t}_{i,j} = \sum_{x=-r}^{y=r} \mathbf{p}_{i+x, j+y} K_{\sigma,r}[r-x, r-y] \quad (12)$$

$$(13)$$

Escribimos cada término de la suma:

$$\hat{\mathbf{p}} = \mathbf{p}_{i+x, j+y} K_{\sigma,r}[r-x, r-y] \hat{\mathbf{p}} = \mathbf{p}_{i+x, j+y} \mathbf{k} \quad (14)$$

Aquí,  $\mathbf{k} = (K_{\sigma,r}[r-x, r-y], K_{\sigma,r}[r-x, r-y], K_{\sigma,r}[r-x, r-y], K_{\sigma,r}[r-x, r-y])$  es el filtro expandido hacia un vector, y el último producto es un producto elemento a elemento. Así, llegamos a una ecuación vectorizada en  $(R, G, B, A)$ , elementos contiguos en memoria que, al ser floats de 4 bytes cada uno, es el tamaño de los registros XMMn. En la implementación de ASM, precalculamos algunos valores que quedan en los registros:

1. Primer elemento del cuadrado de vecinos: La distancia, en memoria, entre el píxel  $(i, j)$  y el píxel  $(i-r, j-r)$ :

<sup>3</sup>Como esta matriz es utilizada muy seguido, intuimos que se mantiene siempre en cache, de forma tal que su acceso sería bastante rápido.



$$d_{1v} = \text{pos}[(i - r), (j - r)] - \text{pos}[i, j] \quad (15)$$

$$= [(i - r)n_{\text{cols}} + (j - r)] - [i n_{\text{cols}} + j] \text{ 4bytes} \quad (16)$$

$$= (-r n_{\text{cols}} - r) \text{ 4bytes} \quad (17)$$

$$= -r (n_{\text{cols}} + 1) \text{ 4bytes} \quad (18)$$

2. Nueva fila en el filtro: La distancia, en memoria, entre el píxel  $(i + x, j + r)$  y el píxel en la siguiente fila  $(i + (x + 1), j - r)$ :

$$\Delta = \text{pos}[(i + x), (j + r)] - \text{pos}[i + (x + 1), (j - r)] \quad (19)$$

$$= [(i + x)n_{\text{cols}} + (j + r)] - [(i + x + 1)n_{\text{cols}} + (j - r)] \text{ 4bytes} \quad (20)$$

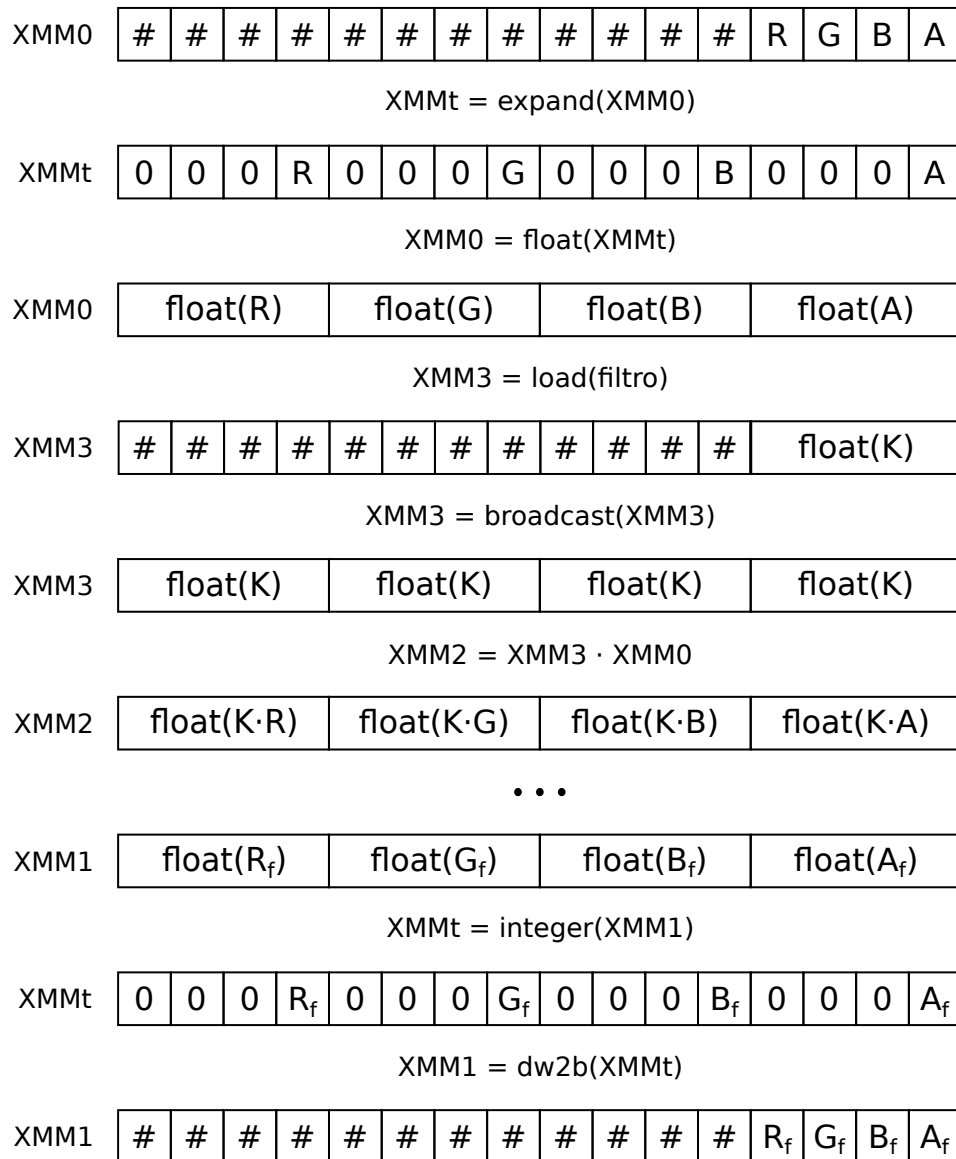
$$= (-n_{\text{cols}} + 2r) \text{ 4bytes} \quad (21)$$

Los ciclos ahora también van a ser sobre el *centro* de la imagen, dentro de cada cual calculamos la salida para el primer píxel  $(i, j)$ . Destinamos un registro, por ejemplo **XMM1**, al acumulador. Dentro de él vamos sumando (con **addps**) cada uno de los términos de acuerdo a la expresión 14. Vamos a explicar con detalle cómo se calcula esta expresión para el primer píxel vecino en  $(i - r, j - r)$ , que es la parte vectorizada.

Primero, cargamos los 4 bytes del píxel en un registro, llamémoslo **XMM2**, con **movd**. Ahora queremos convertirlos a punto flotante. Para eso, tenemos que considerar que la carga anterior llena la parte baja del registro. Para convertirlo, primero vamos a extender estos 4 bytes a 4 doublewords en un registro temporal **XMMt** (mediante **pmovzxbd**) y luego los convertimos a números de punto flotante de 4 bytes con **cvtdq2ps** y los guardamos de nuevo en **XMM2**. En este punto, **XMM2** tiene guardados los 4 “colores” en una representación de punto flotante ( $\mathbf{p}_{i+x, j+y}$ ). Ahora, en la parte baja del registro **XMM3** cargamos los 4 bytes (con **movd**) correspondientes al valor del filtro para  $(-r, -r)$  (usando  $d_{1v}$ ). Ahora, para crear el vector **k** hay que expandirlo (hacer un broadcast del elemento 0, viendo al registro como 4 números de punto flotante de 4 bytes) a todos los elementos, a través de **suhfps XMM3, XMM3, 0x00**. Ya tenemos los dos vectores, y los multiplicamos y lo guardamos en **XMM2** con **mulps**. Ya estamos listos para acumularlo en **XMM1** (**addps**). Ahora pasamos al siguiente elemento en el filtro y en la imagen (agregando 4 bytes). Este procedimiento lo realizamos para toda una fila, y luego, para pasar a la siguiente columna sumamos  $\Delta$  a la posición de memoria de la imagen. Una vez acumulado todo el filtro para un píxel en **XMM1**, lo convertimos a enteros de 4 bytes con **cvtps2dq**. Como el filtro está normalizado a 1, estos enteros no pueden ocupar más de 1 byte [ya que los datos originales son de 1 byte]. De este modo, ahora condensamos en **XMM1** el byte más bajo de cada entero de 4 byte, con una máscara **mask\_dw2b[0, 4, 8, 12]** y **psuhfb xmm3, mask\_d2wb**). Ya podemos guardar **XMM1**, nuevamente con **movd**, en la posición de memoria destino.

Este procedimiento se puede observar en el esquema 2.

Terminado esto, ¡ya tenemos procesado el píxel! Hacemos esto para toda la fila [sumando de a 4 bytes para pasar al siguiente elemento] y, para evitar el *marco*, en el último elemento de la fila sumamos  $2r$  4bytes para pasar al siguiente elemento del *centro*.



**Figura 2:** Esquema que caracteriza el algoritmo de blur con instrucciones SIMD

### 3. Resultados

En esta sección se muestran los tiempos de ejecución de cada algoritmo comparandose las implementaciones hechas en **C** contra las implementaciones hechas en **ASM**.

Dichos resultados se corresponden con la cantidad de ticks del procesador que cada algoritmo tomó, usando para ello el parámetro **-t** del programa principal, cuyo valor devuelto es justamente la cantidad de ticks para ejecutar el filtro de entrada una cantidad de veces equivalente a la pasada como parámetro.

Como medimos los ticks del procesador? Es algo que se hace internamente en el código que la cátedra nos proporcionó pero de todas formas sabemos que esto se lleva a cabo usando la instrucción **rdtsc**, la cual obtiene el Time Stamp Counter (TSC). Dicho registro se incrementa en uno con cada ciclo del procesador, de modo que la cantidad de ciclos total equivale a la diferencia del valor después y antes de ejecutar cada filtro.

Notar que este registro es global y por ende cuenta ticks que todos los procesos del sistema estan consumiendo, no solo el nuestro, de modo que sería incorrecto hacer solo una medición, en vez de eso hacemos 1000 y tomamos el promedio, para suavizar outliers (observación numéricamente muy distante al resto de los valores).

#### 3.1. SpeedUp

En computación paralela el *speedup* refiere a cuánto más rápido es un algoritmo paralelo (en nuestro caso refieren a las implementaciones en Asm que hacen uso de las instrucciones SSE) que el correspondiente algoritmo secuencial (lo que sería cada implementación en C).

Porqué consideramos importante medir que tanto más rápido es la versión paralela que la secuencial? Si bien este tipo de análisis excede originalmente lo pedido por la cátedra, nos pareció razonable dar una idea de la magnitud de que tantas veces es mejor una implementación paralela que secuencial. Si bien este valor se desprende de la cantidad de ticks insumidos por cada implementación, nos pareció adecuado formalizarlo usando el concepto de *speedup* entre las dos implementaciones. Se calcula con la fórmula:

$$S_{vsj} = \frac{T_i}{T_j} \quad (22)$$

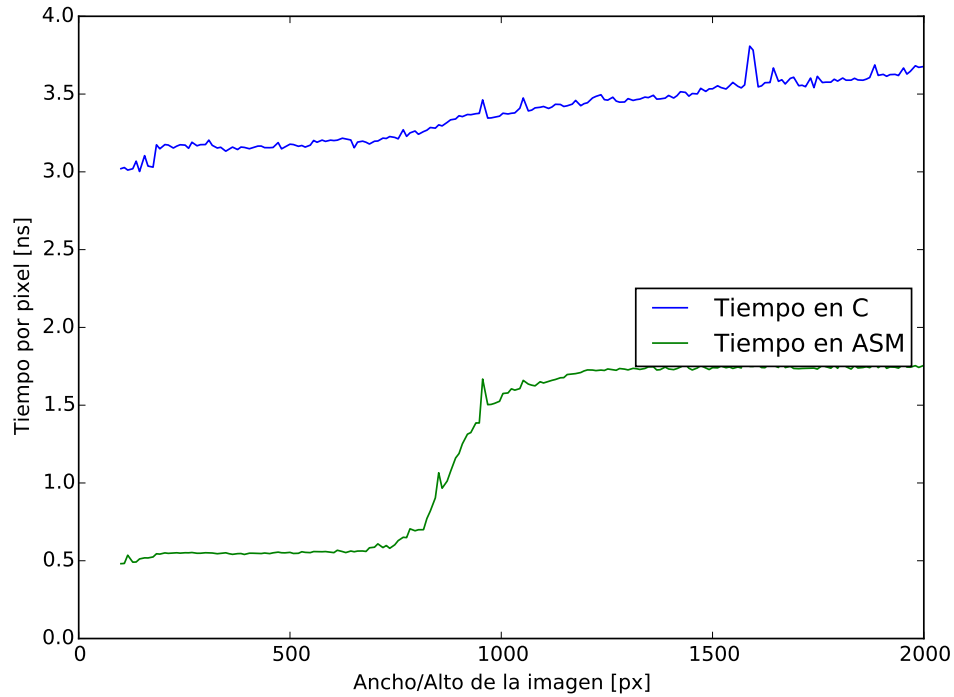
En nuestro caso, será:

- $T_i$ : cantidad de ticks del algoritmo en C
- $T_j$ : cantidad de ticks del algoritmo en ASM

Para el análisis de benchmark generamos un script en **python** que genera imagenes al azar, de distintas dimensiones.

#### 3.2. Comparaciones entre C y ASM: Dependencia en tamaño

Para un *aspect ratio* 1:1, comparamos el tiempo que tarda C con optimización -O3 y ASM en procesar cada imagen.



**Figura 3:** Tiempo por píxel para imágenes cuadradas. Observamos un salto, para ASM, entre 800 y 1000 píxeles.

### 3.2.1. Diferencia de Imágenes

En la figura 3 observamos el tiempo de procesamiento, por píxel, para varias imágenes entre 100x100 y 2000x2000 píxeles.

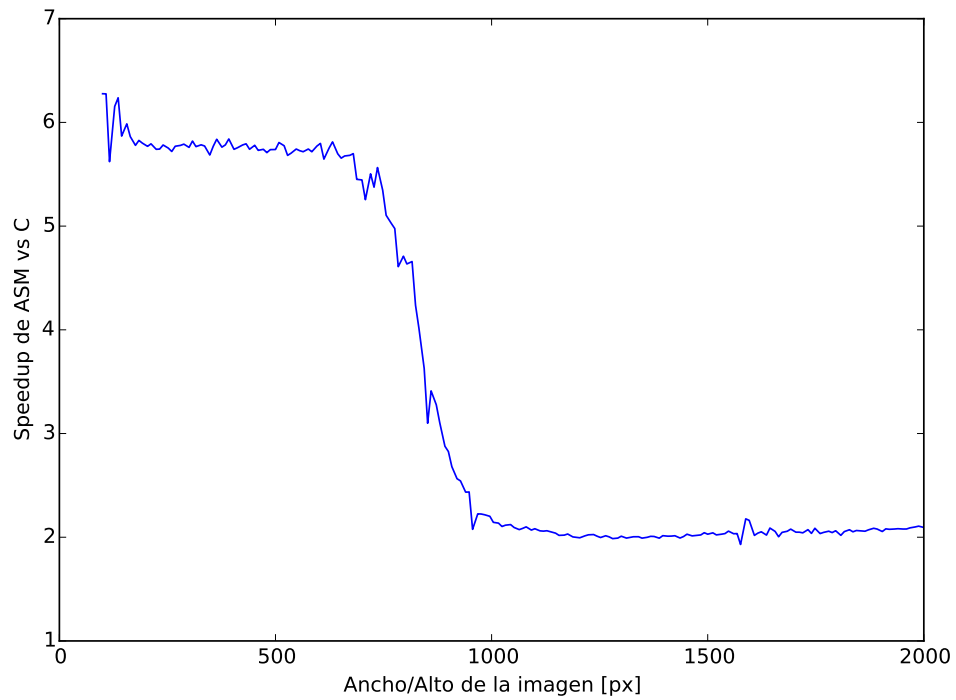
Observamos una muy pequeña tendencia creciente en el tiempo de procesamiento por píxel para la implementación en C (20% más lento en imágenes de 100 veces la cantidad de píxeles), mientras que en el caso de ASM el tiempo parece ser constante entre 200 y 800 píxeles y entre 1000 y 2000, con una bastante abrupta transición en la que se triplica el tiempo de procesamiento entre 800 y 1000 píxeles. Para 800x800 píxeles de imagen, tenemos en total  $0,64 \text{ MP} \cdot 4\text{bytes/pixel} = 2,5\text{MB}$ . Como son en total tres imágenes, estamos usando 7,5MB de memoria, que es el tamaño de la cache L3 del procesador utilizado. Se vuelve a estabilizar al llegar a 1000x1000, 4MB por imagen (dos imágenes entrarían en cache).

En la implementación de C no se observa ese comportamiento, probablemente debido a que la optimización O3 realiza un algoritmo *cache aware*, mediante el flag `-freorder-blocks-and-partition`. Este mismo salto, y el hecho de que suceda sólo en la implementación en ASM, implica un salto también en el *speedup*, como se ve en 4.

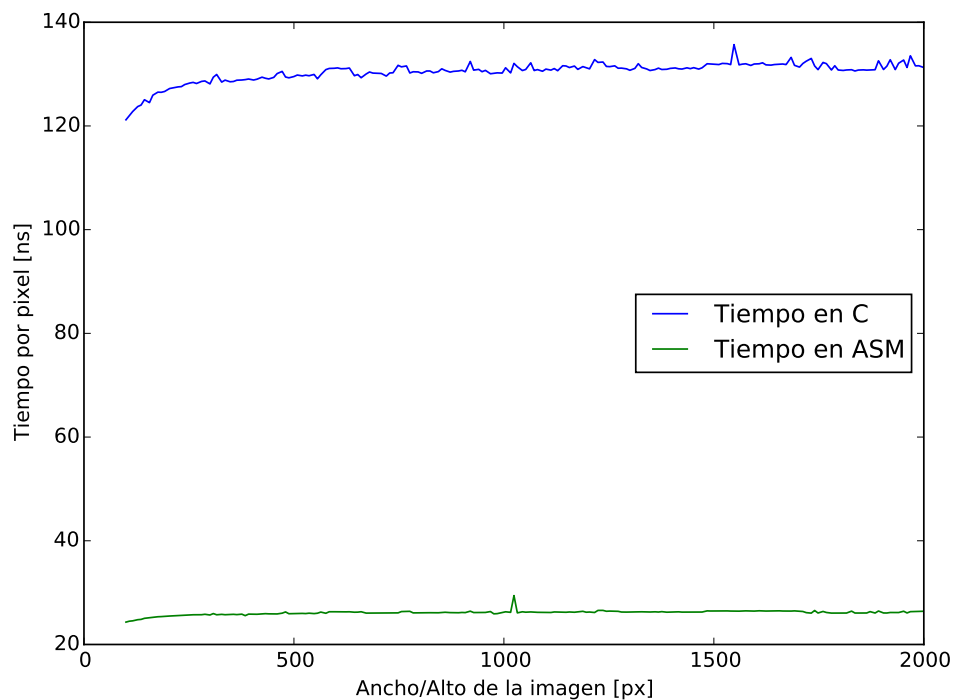
### 3.2.2. Blur Gaussiano

En la figura 5 observamos el tiempo de procesamiento, por píxel, para varias imágenes entre 100x100 y 2000x2000 píxeles.

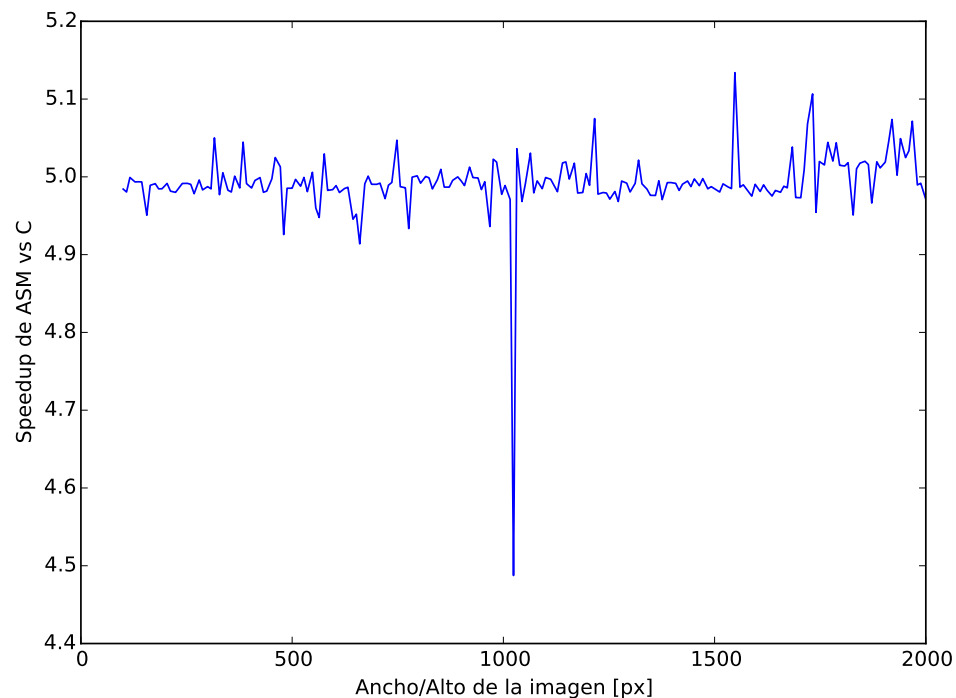
En este gráfico vemos que ambos algoritmos arrancan con una tendencia creciente hasta los 300 px, a partir del cual el tiempo por píxel se mantiene constante. Que ambos sean de la misma tendencia se corrobora en el estudio del *speedup*, en el que se ve que se mantiene uniforme para todos los tamaños



**Figura 4:** *Speedup* para imágenes cuadradas. Observamos un salto entre 800 y 1000 píxeles, que se traslada de la baja de performance en el código de ASM



**Figura 5:** Tiempo por píxel para imágenes cuadradas. Observamos un comportamiento creciente, que se estabiliza a partir de  $\approx 300$  px.



**Figura 6:** *Speedup* para imágenes cuadradas. Observamos un comportamiento uniforme de 5x.

de la imagen.

### 3.3. Comparaciones entre C y ASM: Dependencia en aspect ratio

En este caso, mantuvimos el tamaño de la imagen constante y variamos el aspect ratio. No se observó ninguna dependencia en ninguno de los casos, como se puede ver en las figuras 7, 8, 9, 10.

### 3.4. Comparaciones entre C y ASM: Dependencia con el radio de blur

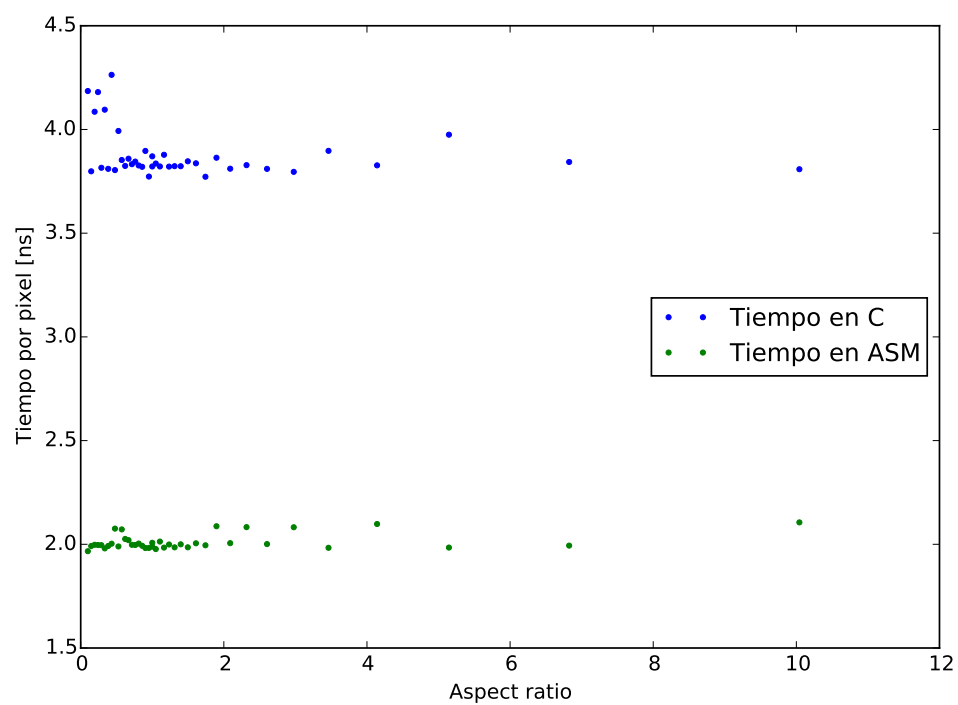
Finalmente, observamos en la figura 11 el tiempo de procesado por píxel para distintos radios. Observamos una tendencia, como era de esperar, creciente. Sin embargo, hacen falta más estudios para poder caracterizar la forma funcional de la dependencia con el tamaño del filtro.

### 3.5. Comparaciones entre C -O0 y C -O3: Filtro de diferencias

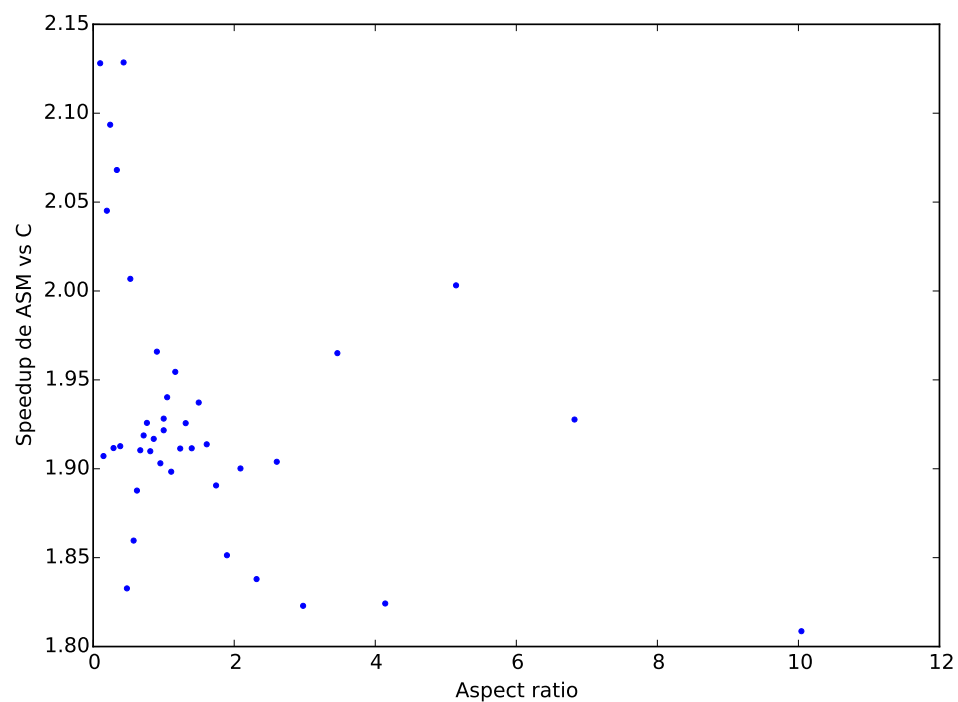
Observamos aquí que -O3 es más rápido que -O0 por un factor 5x, y ese factor va decreciendo con el tiempo

### 3.6. Algunas conclusiones

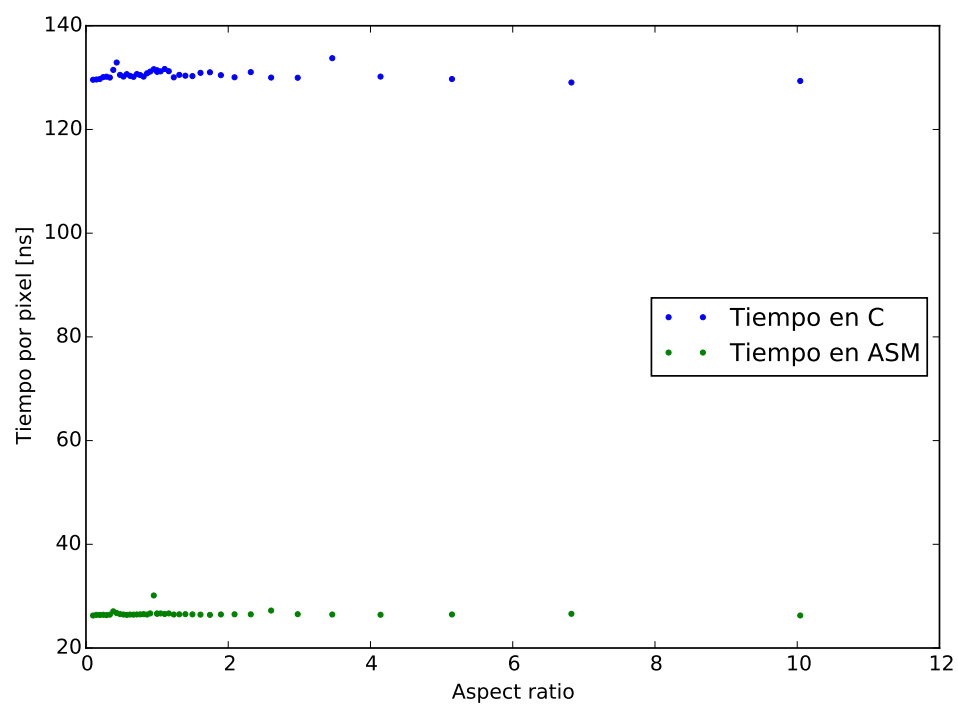
A partir de los resultados y tablas vistas podemos mencionar las siguientes conclusiones generales que aplican a todos los filtros:



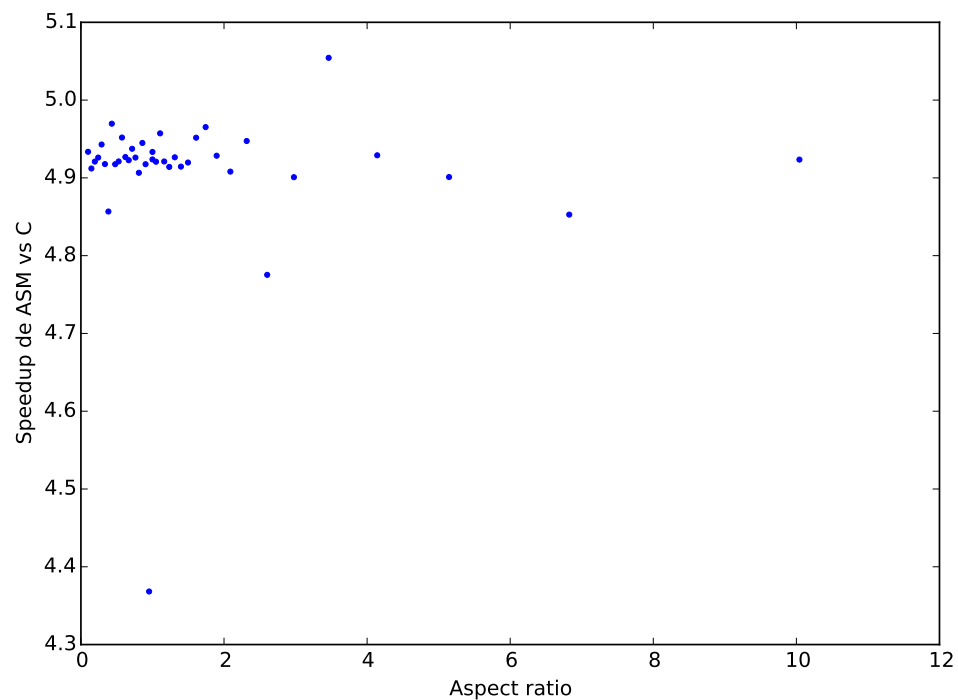
**Figura 7:** Tiempo de procesamiento por píxel para distintos aspect ratio con filtro diff.



**Figura 8:** *Speedup* de procesamiento por píxel para distintos aspect ratio con filtro diff.

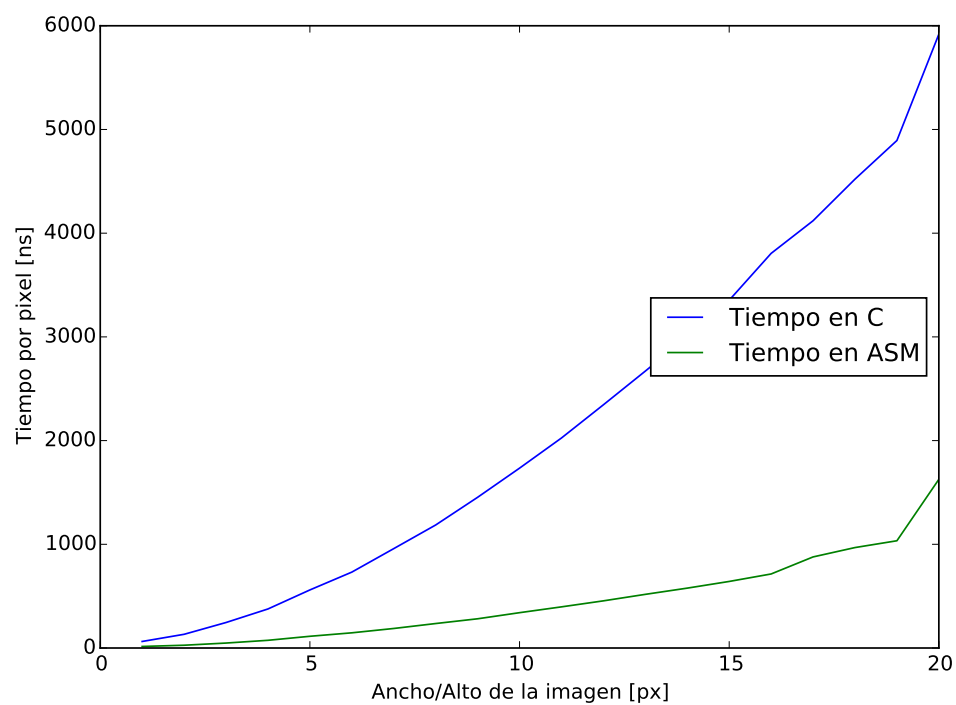


**Figura 9:** Tiempo de procesado por píxel para distintos aspect ratio con filtro blur.

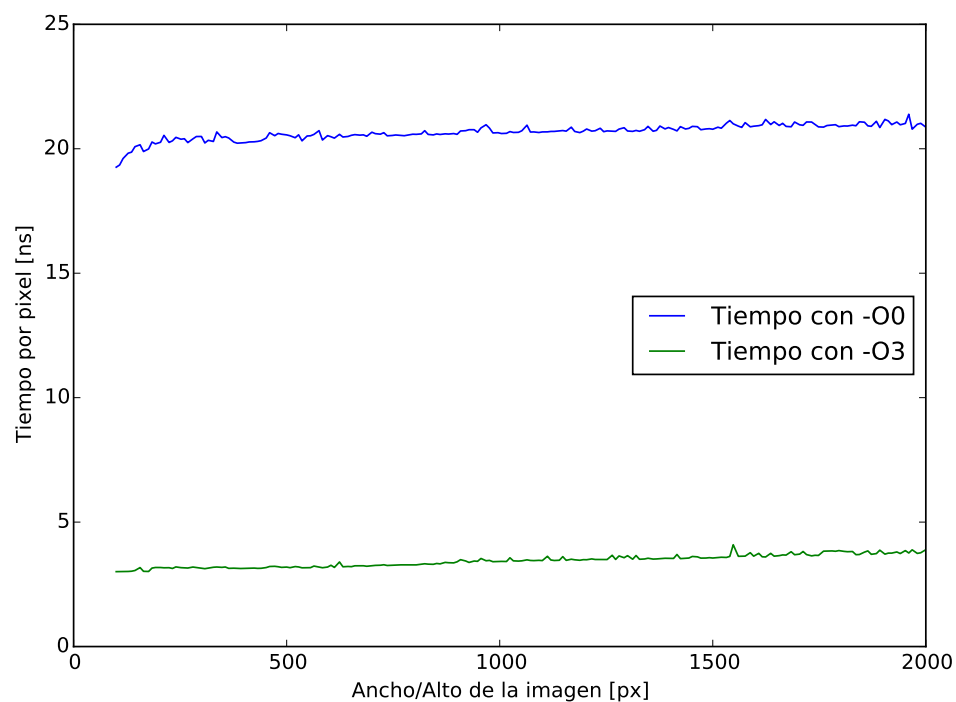


**Figura 10:** *Speedup* de procesado por píxel para distintos aspect ratio con filtro blur.

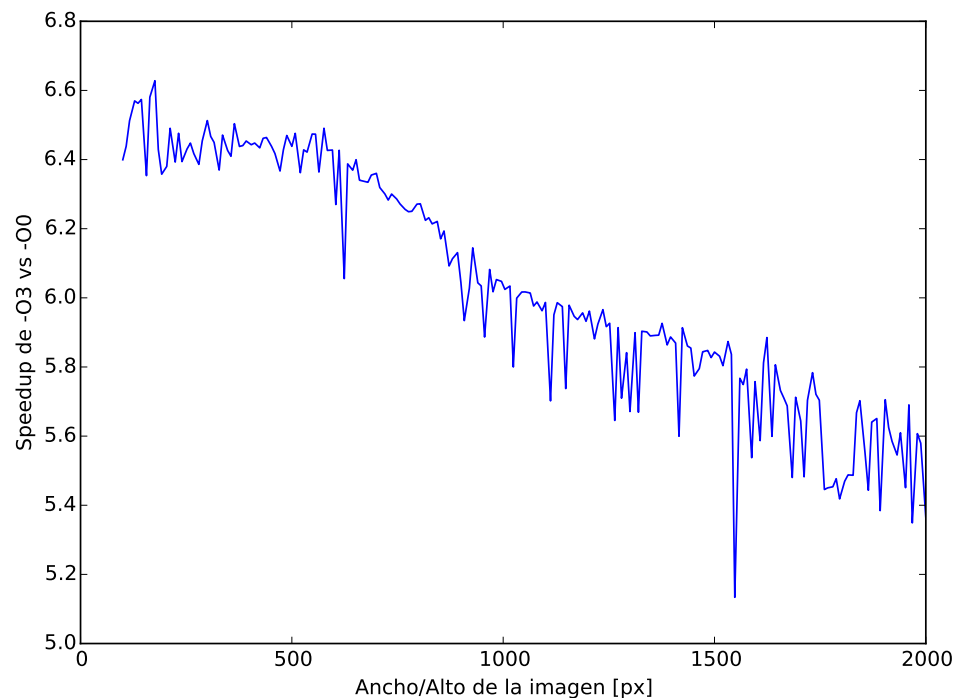




**Figura 11:** Tiempo de procesado por píxel para distintos radios con filtro blur.



**Figura 12:** Tiempo de procesado por píxel para distintos tamaños.



**Figura 13:** *Speedup* de procesamiento por píxel para distintos tamaños.

- Las implementaciones hechas en *Assembler* son efectivamente más rápidas que las implementaciones hechas en *C*. Esto era lo esperado porque las implementaciones en *Assembler* hacen uso de las instrucciones SSE y se procesan más de un píxel simultáneamente.
- En los algoritmos en los que no hubo conversión a Float de los datos, la velocidad de resolución fue mayor porque al no ser necesario usar floats, cada píxel ocupaba menos bytes y por ende podíamos procesar mayor cantidad simultáneamente.
- El SpeedUp de las funciones no cambió demasiado entre las funciones. Y todas estas operaciones se pueden obtener directamente con instrucciones SIMD, mientras que en C hubo que programarlas, realentizando el algoritmo.
- Ambos códigos en ASM son mejores que la optimización más agresiva de gcc (-O3). Restaría comparar con otros compiladores.
- A partir de la inspección del código generado en ASM de gcc pudimos observar que, en -O3, a pesar de utilizar muchos trucos, nunca usa las instrucciones SIMD, es decir, no autovectoriza.
- La mejora por utilizar -O3 en el filtro de diff es de 5x, sólo con escoger un flag del compilador!

## 4. Conclusión

Básicamente, este trabajo práctico nos ha hecho pensar como implementar ciertas funciones de una manera distinta a la que estamos acostumbrados. Además, nos ha hecho conocer en profundidad el modelo *SIMD*.

La programación vectorial es muy útil para optimizar funciones. Creemos que nuestras implementaciones muestran claramente la ganancia en performance que hay al programar así. Sin embargo, hay un costo. Los algoritmos son más complejos debido a que el número de instrucciones aumenta significativamente, se debe pensar cuidadosamente como se cargan los datos de memoria y como se guardan para evitar errores de *Violación de segmento*. Muchas veces los datos deben ser transformados y reacomodados para poder usarlos con las operaciones *SIMD*. También hay que analizar bien cuando finalizar los ciclos y si se debe tratar aparte y como se deben tratar los últimos datos, esos que no pudieron ser tratados en el ciclo debido a que no eran suficientes como para poder cargarlos en un registro `xmmx`. La otra desventaja que tiene este modo de operación es que nos ata a una arquitectura específica haciendo nuestro programa no portable. Nos ha pasado tener que repensar ciertos partes del código por no contar, por ejemplo, con la extensión **SSE3**.

Para cerrar, nos resultó interesante ver la aplicación de la programación vectorial en un tema concreto como el tratamiento de imágenes.