



DEPARTAMENTO
DE COMPUTACION

Facultad de Ciencias Exactas y Naturales - UBA

Trabajo Práctico II

Organización del Computador II
Primer Cuatrimestre de 2019

Integrante	LU	Correo electrónico
Rodrigo Laconte	193/18	rola1475@gmail.com
Julia Rabinowicz	48/18	julirabinowicz@gmail.com
Amalia Sorondo	281/18	sorondo.amalia@gmail.com



Facultad de Ciencias Exactas y Naturales
Universidad de Buenos Aires

Ciudad Universitaria - (Pabellón I/Planta Baja)

Intendente Güiraldes 2160 - C1428EGA

Ciudad Autónoma de Buenos Aires - Rep. Argentina

Tel/Fax: (54 11) 4576-3359

<http://www.fcen.uba.ar>

Resumen

En el presente trabajo implementamos tres filtros para imágenes. En base a estas implementaciones, realizamos experimentos con el fin de explorar algunos de los factores que influyen en la cantidad de ciclos de clock insumidos por una función en ASM. Analizamos de esta forma cómo afectan en el costo temporal procedimientos como llamados a funciones, saltos condicionales o distintas formas de acceder a memoria. Junto con el análisis de los resultados intentamos sacar conclusiones contrastándolas con nuestras hipótesis iniciales.

Índice

1. Introducción	3
2. Desarrollo	4
2.1. Filtro de Nivel	4
2.2. Filtro Bordes	4
2.3. Filtro Rombos	5
3. Resultados	6
3.1. Nivel	6
3.1.1. Experimento 1: saltos condicionales	6
3.1.2. Experimento 2: llamado a funciones	8
3.2. Bordes	10
3.2.1. Experimento 1: ASM vs C	10
3.2.2. Experimento 2: parámetros	12
3.3. Rombos	14
3.3.1. Experimento 1: Máscaras (no) alineadas	14
3.3.2. Experimento 2: Accesos a memoria dentro y fuera del ciclo	16
4. Conclusión	18

1. Introducción

En la actualidad entre todas las funcionalidades de las computadoras, una a destacar es la edición y procesamiento de archivos multimedia, tales como audio, video o imágenes. En este contexto, los algoritmos suelen ser repetitivos ya que generalmente se quiere aplicar el mismo procedimiento a los datos de cada archivo. En la década del '60, surgió la idea de poder procesar, en una sola instrucción, varios de estos datos¹. Con este propósito, y con la ventaja que brindó la aparición de los registros de mayor capacidad de almacenamiento en los procesadores, comenzaron a aparecer las instrucciones SIMD (Single Instruction Multiple Data).

A lo largo de este informe vamos a enfocarnos en el procesamiento de imágenes. Su utilidad e importancia consisten en el reconocimiento de patrones en imágenes por ejemplo, tomar de una foto una patente, o una cara, pero también tiene un lado más enfocado al aspecto comercial, como lo es la edición de imágenes.

En este trabajo vamos a analizar las implementaciones de tres filtros en particular: nivel, bordes y rombos. Los mismos se encuentran explicados en la sección de desarrollo. Con este objetivo, en primer lugar implementamos los filtros en lenguaje ASM, utilizando instrucciones SIMD. Esto nos sirvió como base de la experimentación que realizamos de los mismos.

Todos los experimentos estudiados en este informe comparan distintas implementaciones de los filtros mencionados, midiendo la eficiencia de cada una en función de la cantidad de ciclos de reloj. El propósito de la experimentación es determinar qué factores influyen en la performance del código de ASM, por lo que nos centramos en aspectos tales como accesos a memoria, cantidad de saltos, contraste con código de C, entre otras cuestiones.

Para cada experimento planteado, proponemos hipótesis previas sobre lo que esperamos que ocurra, es decir cómo creemos que se modificará la eficiencia de cada implementación, basada en nuestras ideas y conocimientos previos. Luego, efectuamos cada uno de ellos y analizamos los resultados, comparándolos con los esperados e intentando explicar, en los casos en los que nuestras predicciones difieren de lo obtenido, las razones por las cuales esto sucedió.

¹<https://en.wikipedia.org/wiki/SIMD>

2. Desarrollo

2.1. Filtro de Nivel

Nuestra implementación del filtro de nivel es la más simple de las tres en cuanto al código. Consiste en un ciclo en el cual en cada iteración son procesados cuatro píxeles de la imagen, levantándolos en un registro XMM. Esta cantidad de píxeles ocupa el registro de 128 bits entero ya que cada uno está compuesto por cuatro componentes de un byte: azul (b), verde (g), rojo (r) y transparencia (a).

```
XMM0 = [ PIXEL 3 | PIXEL 2 | PIXEL 1 | PIXEL 0 ]
XMM0 = [A3|R3|G3|B3|A2|R2|G2|B2|A1|R1|G1|B1|A0|R0|G0|B0]
```

Previo al ciclo, nos guardamos en otros registros XMM dos máscaras que luego utilizaremos. De esta manera evitamos acceder a memoria en cada iteración. Una de las máscaras se encuentra relacionada directamente con el parámetro índice que recibimos, entre cero y siete. En la sección `.rodata` guardamos las máscaras en orden, escribiendo primero la necesaria para el índice igual a cero, luego la del uno, y así hasta el siete. De esta manera, accederemos a las máscaras direccionando con `[mask0 + índiceParámetro]`.

```
section .rodata
align 16
mask0: times 16 db 0x01
.
.
.
mask7: times 16 db 0x80
```

La segunda máscara la usamos para setear la transparencia de los píxeles a 0xFF al final del proceso. También nos armamos un contador para recorrer la imagen fuente utilizando los parámetros que la describen (ancho y alto).

Luego, el cuerpo del ciclo primero se encarga de levantar 4 píxeles en un registro de 128 bits para correrle un PAND contra la máscara. Como resultado, obtendríamos un registro R1 con el bit que nos interesa seteado en 1 (en cada paquete) solo si esta en 1 en los componentes de los píxeles que levantamos. Paso siguiente, comparamos por igualdad los píxeles con la máscara de R1, así tendremos cada paquete (uno por cada componente de cada píxel) en 0xFF si está prendido el bit en cuestión. Previo a mover los píxeles a la imagen de destino, y por último, debemos setear los componentes de transparencia en 0xFF por si se hayan modificado accidentalmente. Para esto, realizamos un POR entre la máscara de la transparencia y los píxeles modificados.

2.2. Filtro Bordes

El filtro de bordes consiste en aplicar operaciones matriciales sobre los píxeles, es decir que modifica cada píxel en función de sus píxeles vecinos y las matrices operadores Gx y Gy. El filtro opera sobre imágenes en escala de grises, cuyos píxeles cuentan con una única componente y por lo tanto ocupan un solo byte. Luego, en un registro XMM se pueden almacenar 16 píxeles. Sin embargo, el filtro requiere realizar operaciones cuyo resultado podría irse de la representación en byte. Por este motivo decidimos procesar de a 8 píxeles en paralelo para poder desempaquetarlos y así evitar pérdida de información al realizar las cuentas correspondientes. Levantamos entonces datos ubicándolos en la parte baja de los registros XMM gracias a la instrucción `movq`.

Para poder operar con los píxeles vecinos de forma paralela, en cada iteración del ciclo levantamos datos de la siguiente forma:

Sea el cuadro 1 un fragmento de la imagen a la cual le estamos aplicando el filtro y supongamos que nos encontramos en una iteración arbitraria del ciclo y queremos procesar los píxeles "B1,C1,D1,E1,F1,G1,H1,I1". En cada iteración vamos incrementando rdi, el puntero a la fuente, en 8 bytes. Desreferenciamos posiciones relativas a rdi para acceder a los vecinos de los píxeles a procesar. Utilizamos inicialmente 8 registros XMM:

0xA0	0xB0	0xC0	0xD0	0xE0	0xF0	0xG0	0xH0	0xI0	0xJ0
0xA1	0xB1	0xC1	0xD1	0xE1	0xF1	0xG1	0xH1	0xI1	0xJ1
0xA2	0xB2	0xC2	0xD2	0xE2	0xF2	0xG2	0xH2	0xI2	0xJ2

Cuadro 1: Fragmento de imagen

```

XMM1 = [0xA0,0xB0,0xC0,0xD0,0xE0,0xF0,0xG0,0xH0]
XMM2 = [0xB0,0xC0,0xD0,0xE0,0xF0,0xG0,0xH0,0xI0]
XMM3 = [0xC0,0xD0,0xE0,0xF0,0xG0,0xH0,0xI0,0xJ0]
XMM4 = [0xA1,0xB1,0xC1,0xD1,0xE1,0xF1,0xG1,0xH1]
XMM6 = [0xC1,0xD1,0xE1,0xF1,0xG1,0xH1,0xI1,0xJ1]
XMM7 = [0xA2,0xB2,0xC2,0xD2,0xE2,0xF2,0xG2,0xH2]
XMM8 = [0xB2,0xC2,0xD2,0xE2,0xF2,0xG2,0xH2,0xI2]
XMM9 = [0xC2,0xD2,0xE2,0xF2,0xG2,0xH2,0xI2,0xJ2]

```

Es decir, que los registros contienen la siguiente información según las cuentas a realizar:

```

↖  ↑  ↗ = xmm1 xmm2 xmm3
←    → = xmm4      xmm6
↙  ↓  ↘ = xmm7 xmm8 xmm9

```

Luego desempaquetamos: extendemos la representación de los píxeles de un byte a una word. Como las cuentas con las matrices Gx y Gy involucran las dos a los vecinos "de las esquinas" salvamos la información duplicándola en los registros XMM10, XMM11, XMM12, XMM13. Luego realizamos las operaciones correspondientes a Gx y Gy, que consisten en mantener dos acumuladores totalGx (XMM14) y totalGy (XMM15) a los que les sumamos el valor del píxel vecino por el valor de la posición asociada al vecino en la matriz Gy o Gx. Estos acumuladores operan de forma paralela: cada word representa la acumulación de uno de los píxeles procesados. Tomamos el valor absoluto de los acumuladores y realizamos la suma saturada empaquetada. Finalmente empaquetamos: el resultado final de cada píxel pasa a ocupar un byte y será el nuevo valor del píxel que es copiado en la posición apuntada por rsi, que es el puntero al destino.

Decidimos considerar la última iteración del ciclo por separado por cuestiones prácticas, evitando de esta forma accesos a memoria no correspondiente a la imagen. Finalmente, los píxeles que forman los bordes de la imagen son puestos en 255 recorriendo primero la primer fila de la imagen, luego la última columna, seguida por la última fila y por último la primer columna.

2.3. Filtro Rombos

Al aplicar el filtro "Rombos" a una imagen, cada uno de sus píxeles se ve modificado en función de sus coordenadas en la imagen. Cada píxel está compuesto en este caso por 4 componentes que ocupan un byte cada una, luego un píxel ocupa 4 bytes. Como los resultados de las operaciones efectuadas con las componentes pertenecen al rango [-128,127], no es necesario desempaquetar los datos ya que pueden representarse en un byte sin correr riesgo de overflow. Podemos afirmar esto en base a las cuentas: primero le restamos a size/2, es decir 32, un número entre 0 y 63. Obtenemos entonces un número entre -31 y 32, y nos quedamos con el valor absoluto: un número entre 0 y 32. Finalmente, sumamos dos valores pertenecientes al rango [0,32], obteniendo un número entre 0 y 64 al que le restamos 32.

Dependiendo de este último resultado asignamos a x el valor 0 o un número entre 0 y 64. Es decir luego de la primera cuenta realizada nuestros datos ocupan un byte como máximo. Levantamos entonces de 4 píxeles, siendo ésta la mayor cantidad de píxeles que pueden ser procesados en simultáneo, ocupando un registro XMM completo.

Como es necesario tener en cuenta las coordenadas de cada píxel, recorremos la imagen gracias a dos ciclos anidados, llevando un contador que indica en cuál fila nos encontramos y otro que indica en cuál columna. Nos posicionamos entonces en una fila y vamos recorriendo sus columnas. En cada iteración realizamos lo detallado a continuación. Levantamos cuatro píxeles contiguos, los cuales comparten la misma fila ya que el ancho de las filas son múltiplo de 8 y por ende también múltiplo de 4. De esta forma, tenemos en un registro XMM las cuatro componentes de cada uno de los cuatro píxeles, a través del cual vamos a aplicarle las operaciones a cada componente, es decir a cada byte del registro.

Primero calculamos el x que debemos sumarle a cada componente de cada píxel, para luego guardar estos x en un registro XMM $[x_3, x_2, x_1, x_0]$ ² donde cada uno ocupa el byte menos significativo de la double word que le corresponde. Como $size$ es igual a 64, para representar $size/2$ y $size/16$ levantamos de memoria las máscaras conformadas por $[32, 32, 32, 32]$ y $[4, 4, 4, 4]$. Por otro lado, siendo i la fila actual y j las 4 columnas actuales, tenemos en un registro $[j+3, j+2, j+1, j]$ y en otro $[i, i, i, i]$. Para obtener los restos de ambos registros módulo 64, basta con quedarnos con los 6 bits más significativos y poner en 0 el resto ya que 2^6 es igual a 64. Realizamos entonces un shift lógico a izquierda y luego a derecha de los paquetes de doublewords. Le restamos al registro $[32, 32, 32, 32]$ los restos y nos quedamos con los valores absolutos de cada doubleword. Con cuentas empaquetadas similares calculamos $[i+(j+3)-32, i+(j+2)-32, i+(j+1)-32, i+j-32]$. Gracias a la instrucción `pcmpgtd` nos creamos una máscara que tiene 1s en las doublewords donde $i+(j+k)-32$ es mayor a $size/16$. Utilizando esa misma máscara negada guardamos los $i+(j+k)-32$ menores a $size/16$. Luego comparando con una máscara de ceros, separamos positivos de negativos para restar el valor absoluto de éstos últimos y sumar el valor absoluto de los primeros a cada componente de los píxeles. Antes de realizar esto último utilizamos la instrucción `pshufb` para pasar de $[0,0,0,x_3,0,0,0,x_2,0,0,0,x_1,0,0,0,x_0]$ a $[x_3,x_3,x_3,x_3,x_2,x_2,x_2,x_2,x_1,x_1,x_1,x_1,x_0,x_0,x_0,x_0]$, ya que le queremos sumar x a cada componente de cada píxel. Movemos el resultado al destino y avanzamos los punteros a destino y a fuente y sumamos a la columna actual 4. Para actualizar el registro conteniendo $[j+3, j+2, j+1, j]$ le sumamos la máscara $[4, 4, 4, 4]$. Cuando llegamos al final de la columna, salimos del ciclo interno de columna: incrementamos la fila y reseteamos las columnas a $[3, 2, 1, 0]$.

3. Resultados

Para todos los experimentos, analizamos el costo temporal de un filtro dado un parámetro, haciendo el promedio de la cantidad de ciclos de clock de 10000 corridas. Esto lo realizamos 5 veces, devolviendo el promedio obtenido. Esto lo corremos para la implementación modificada para el fin de cada experimento, y para la implementación original, comparando ambos resultados. Todos los experimentos fueron corridos en la computadora con las siguientes características: versión del sistema operativo Ubuntu 18.04.3 LTS, con 12GB RAM y procesador Intel core i5-8265U CPU @1.60GHz x8.

3.1. Nivel

Buscamos determinar mediante la experimentación con diferentes implementaciones en ASM del filtro de nivel los costos en términos temporales de distintos procedimientos o recursos puestos en juego al implementar funciones en este lenguaje.

3.1.1. Experimento 1: saltos condicionales

Como primer paso nos cuestionamos si los saltos condicionales implican un costo considerable en lo que es el tiempo de ejecución de la función, medido en ciclos de clock.

²Representamos los registros de forma que a la derecha se encuentran los bits menos significativos

Hipótesis

Esperamos que el tiempo de ejecución de la implementación del filtro con el ciclo desarrollado sea menor a la implementación original. Esto sería producto de ahorrarnos el costo particular de realizar un salto. El procesador, para optimizar la cantidad de ciclos de clock utilizados para realizar el fetch, decode, execute de una instrucción, realiza lo que se denomina como pipeline. El pipelining consiste en superponer en el tiempo la ejecución de varias instrucciones a la vez. De esta forma los distintos bloques del procesador trabajan en paralelo y de forma simultánea con distintas instrucciones. Un salto condicional representa una discontinuidad en el flujo de ejecución, por lo que es considerado un obstáculo en esta optimización que se caracteriza por buscar instrucciones en secuencia. El salto puede hacer que todas, o muchas de las instrucciones que se encontraban preprocesadas deban descartarse.

Experimentos

Para analizar el costo que generan los saltos en el código desarrollamos el ciclo una vez, reduciendo por la mitad la cantidad de iteraciones del ciclo, por lo que se producen la mitad de saltos condicionales. No pudimos seguir desarrollándolo ya que el ancho de la imagen es múltiplo de 8, entonces si levantábamos más de 8 píxeles por iteración podíamos tener accesos inválidos a memoria. Para realizar una comparación "justa", buscamos tener la misma cantidad de instrucciones dentro ciclo en ambas implementaciones, luego si hay diferencias en los costos temporales van a deberse únicamente al costo de la instrucción de salto. Para lograr esto incluimos una instrucción innecesaria en la implementación con el ciclo desarrollado.

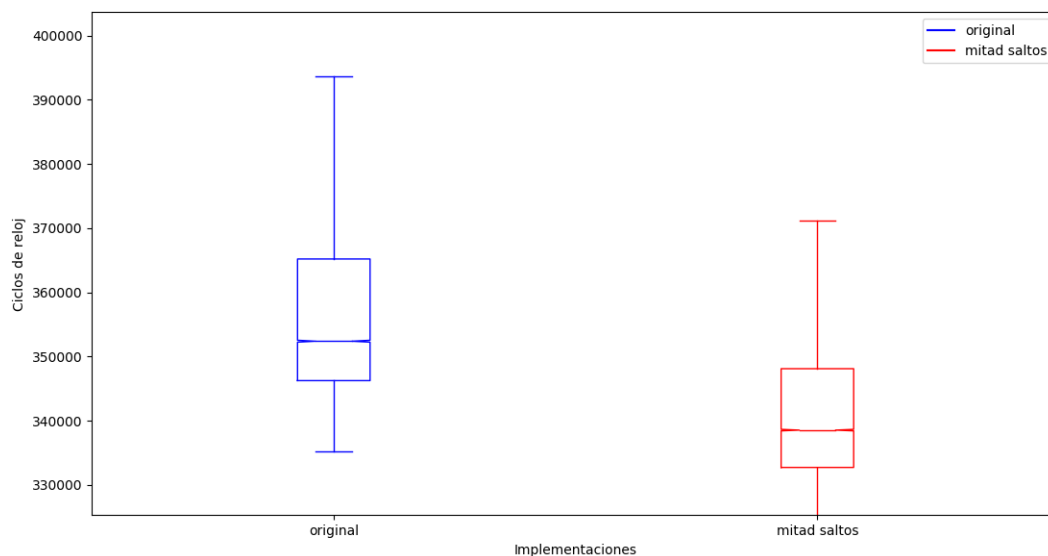


Figura 1: Experimento saltos condicionales

Análisis resultados y conclusiones

Podemos observar en el gráfico que la implementación original lleva en promedio más ciclos de clock que la nueva implementación que cuenta con la mitad de saltos condicionales. Ambas implementaciones tienen la misma cantidad de instrucciones en sus ciclos, luego la diferencia en la cantidad de ciclos de clock promedio que llevan se produce por el número de saltos condicionales totales que realizan. Nuestra hipótesis sobre el costo considerable de las instrucciones de saltos condicionales parece verificarse. Este

costo puede justificarse, como ya explicamos, por la obstaculización que implica un salto condicional en el proceso de pipelining. La varianza y la dispersión entre los ciclos de clock insumidos por las diferentes corridas son similares entre las dos implementaciones, es decir en ambas se producen las mismas diferencias de performance entre una corrida y otra de la misma implementación.

3.1.2. Experimento 2: llamado a funciones

Vimos que los saltos implican un costo en cuanto al tiempo de ejecución de una función, ahora bien, ¿qué sucede con los llamados a funciones?

Hipótesis

Los llamados a funciones implican pushear la dirección de retorno a la pila, luego realizar un salto a la dirección a la cual hace referencia la instrucción "call". Al finalizar la función llamada, se popea la dirección de retorno y se efectúa un nuevo jump a ésta. La desventaja de hacer llamados a funciones entonces, además de pagar el costo asociado a realizar un salto que estudiamos anteriormente, es que al trabajar con la pila se está trabajando con la memoria. Por estas razones esperamos que sea más ineficiente en cuanto al tiempo.

Experimentos

Para estudiar esto hicimos tres nuevas implementaciones del filtro de nivel:

Modificación 1: el ciclo hace un llamado a una función que efectúa las operaciones que hace el ciclo en la implementación original.

Modificación 2: el ciclo llama a tres funciones distintas que efectúan cada una una operación simple realizada por ciclo en la implementación original.

Modificación 3: el ciclo llama a tres funciones distintas, donde éstas realizan a su vez dos llamados adicionales cada una, y la última función que es llamada es la que se encarga de realizar una operación del ciclo.

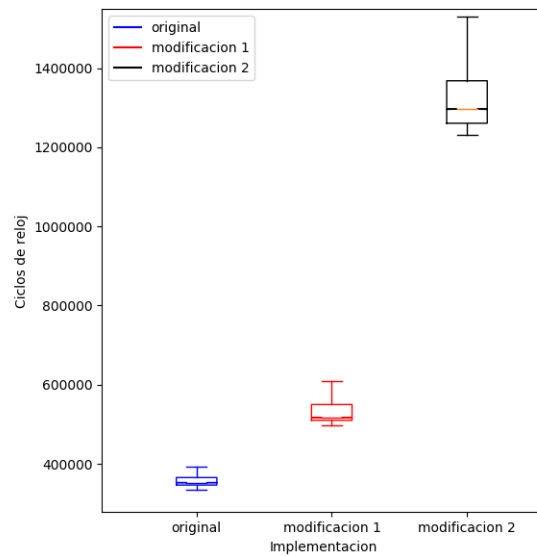


Figura 2: Experimento llamado a funciones

En la figura 2 no incluimos la representación de la modificación 3 ya que la cantidad de ciclos de clocks insumidos por ésta era tanto mayor a la de las otras modificaciones que se iba de la escala.

En la siguiente tabla comparamos el promedio de ejecutar 3 veces 10000 corridas de cada modificación y la implementación original del filtro:

	Original	Modificación 1	Modificación 2	Modificación 3
Cantidad de ciclos de clock promedio	$3,5 \cdot 10^5$	$5,7 \cdot 10^5$	$13 \cdot 10^5$	$38 \cdot 10^5$

Cuadro 2: Promedio de la cantidad de ciclos de clock efectuados por cada implementación

Análisis resultados y conclusiones

Podemos ver en el cuadro cómo crece notablemente la cantidad de ciclos de clock al realizar más llamados a funciones dentro del ciclo. Podemos concluir que llamar a funciones representa un gran costo, lo que coincide con lo que supusimos inicialmente. Sin embargo, esto era predecible ya que anteriormente verificamos que efectuar saltos es costoso. Luego, como llamar a una función implica un jump, se deduce que llamar a funciones perjudica el desempeño temporal de una función. Decidimos entonces realizar un nuevo experimento para poner en evidencia de forma más explícita el gasto temporal en la ejecución de una función producido por el llamado a otras funciones.

Buscamos contrastar el costo de ejecutar un jump con el de realizar un llamado a función. Como los llamados modifican la pila, esperamos que estos insuman una mayor cantidad de ciclos de clock. Para experimentar esta cuestión, alteramos la modificación 1 para que en lugar de realizar un llamado a una función que realiza las operaciones del ciclo efectúe un salto a la etiqueta donde se encuentran las operaciones y luego otro salto para volver al ciclo. Alcanzamos los siguientes resultados:

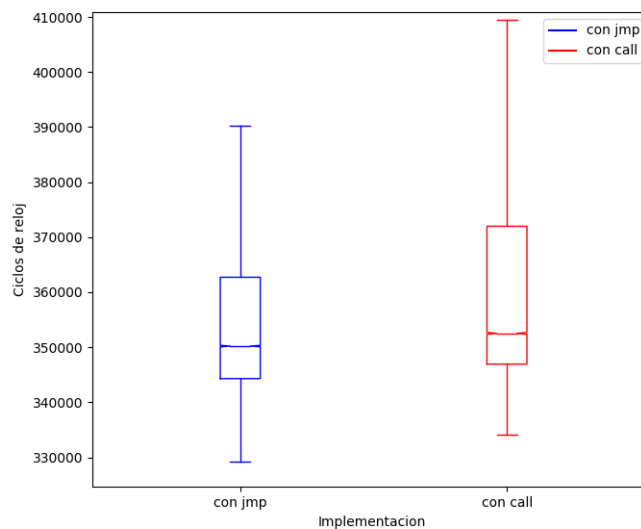


Figura 3: Experimento llamado a función vs saltos

Constatamos que en concordancia con nuestra conjetura, los llamados a función implican un mayor costo temporal que el implicado por los saltos. Esto se explica por la interacción con la pila, y por ende con la memoria, que conlleva el `call` a una función al pushear primero la dirección de retorno y luego poppearla para regresar. Además los llamados a función implican un salto para ir a la función llamada y otro para volver. No obstante, la diferencia no es muy significativa, por lo que podemos inferir entonces que el costo de la modificación de la pila es despreciable en comparación con los demás costos involucrados en la ejecución.

3.2. Bordes

3.2.1. Experimento 1: ASM vs C

El primer encuentro con ASM no resultó tan reconfortante, por lo menos para nosotros. Códigos poco legibles, instrucción-operando-operando, parecía que estábamos martillando la computadora. Nada parecido a los lenguajes que veníamos viendo. Pero si uno se fija, al estar codeando en un nivel tan bajo, se cuenta con una interacción más cercana con los componentes del procesador y su arquitectura, ya que las instrucciones del lenguaje son provistas por los productores de estos. Por lo tanto, nos preguntamos cuál de los dos lenguajes resulta en una implementación más rápida.

Hipótesis

Como tenemos entendido, los códigos en lenguaje C deben pasar por un proceso de compilación, para transformarse en lenguaje ensamblador. Aunque no contamos con completos conocimientos sobre cómo funciona cada compilador, esta compilación es probable que genere instrucciones innecesarias, o algunas que podríamos implementar de manera más eficiente en ASM directamente nosotros. Por esta razón, suponemos que la implementación del filtro de Bordes en ASM se efectúe en menos ciclos de clock que la de C. Además, como C no funciona con procesamiento vectorizado como SIMD, esperamos que ese sea otro factor por el cual el filtro implementado en ASM finalice su ejecución en menos ciclos de reloj.

Experimentos

Comprobaremos si nuestras suposiciones son correctas tomando una imagen y realizando varias corridas con ambas implementaciones.

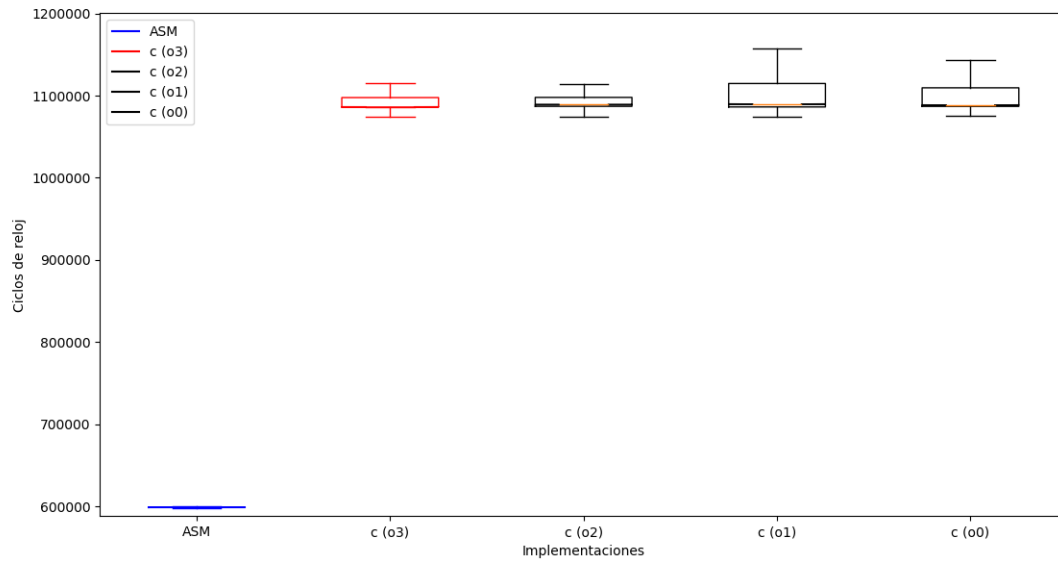


Figura 4: Comparación entre la implementación de ASM y la de C utilizando diferentes flags de optimización

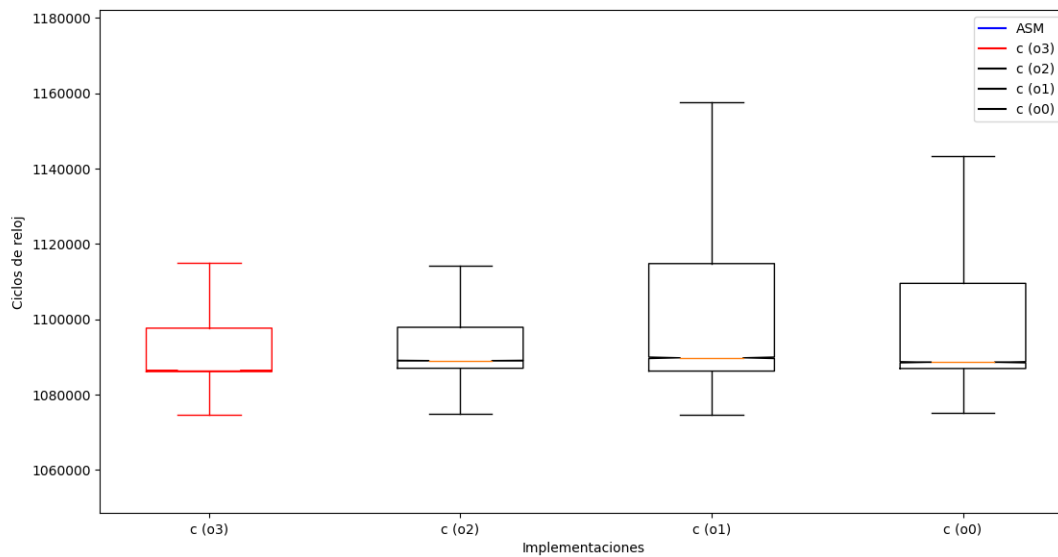


Figura 5: Diferencias entre la implementación de C con los diferentes flags de optimización

Análisis de resultados y conclusiones

Como podemos observar en el gráfico de la figura 4, hay una diferencia abismal entre la implementación del filtro en ASM y la de C, incluso al utilizar la mayor optimización para este último. Efectivamente utilizar SIMD presenta una gran ventaja al momento de procesar datos de forma repetitiva.

Por otro lado, como muestra la figura 5, hay ciertos cambios en los tiempos de ejecución que resultan de utilizar diferentes flags de optimización. Esto era de esperar ya que cada uno funciona de maneras distintas y se encargan de optimizar ciertos criterios de eficiencia de un programa³. El flag -O0 es el flag por default y optimiza el tiempo de ejecución sobre el de compilación, al revés de lo que hace el flag -O1. Los flags -O2 y -O3 optimizan aún más que el -O1 el tiempo de ejecución pero con la desventaja de que puede tardar más la compilación. En nuestro caso, contrastando con el código de ASM, ninguno de los flags mencionados optimizan lo suficiente como para contrarrestar el costo de procesar de a un píxel contra procesar de a 8.

3.2.2. Experimento 2: parámetros

El filtro "Bordes" modifica cada píxel de la imagen en función de sus píxeles cercanos, es decir que el resultado de cada píxel depende de los píxeles vecinos correspondientes a la imagen pasada por parámetro. Podemos preguntarnos entonces: ¿el desempeño de la función depende de la imagen a la cual se le aplica el filtro? El parámetro tiene dos características relevantes que podemos hacer variar para intentar responder a esta pregunta: su tamaño y las componentes de sus píxeles.

Hipótesis

El filtro realiza la misma cantidad de operaciones sin importar el valor de los píxeles. Por este motivo suponemos que aunque el valor por el cual el píxel es modificado depende de la imagen pasada por parámetro, el desempeño de la función no varía. Esperamos también que el tiempo que tarda la función al procesar una imagen es proporcional al tamaño de la imagen, ya que el tamaño del parámetro se traduce en la cantidad de iteraciones del ciclo.

Experimentos

Para contrastar nuestras hipótesis con resultados experimentales realizamos dos experimentos. El primero consiste en comparar la cantidad de ciclos de clock de la función con una foto blanca como parámetro y con una foto con filtro "ruido". Estos dos parámetros buscan representar dos opuestos; la foto blanca no presenta bordes y la otra solo tiene bordes. Comparamos estos dos extremos con la imagen "puente" que nos brindó la cátedra. Esta última imagen sería un ejemplo de imagen promedio, es decir, no perteneciente a ninguno de los extremos. Observamos los resultados:

³<https://www.rapidtables.com/code/linux/gcc/gcc-o.html>

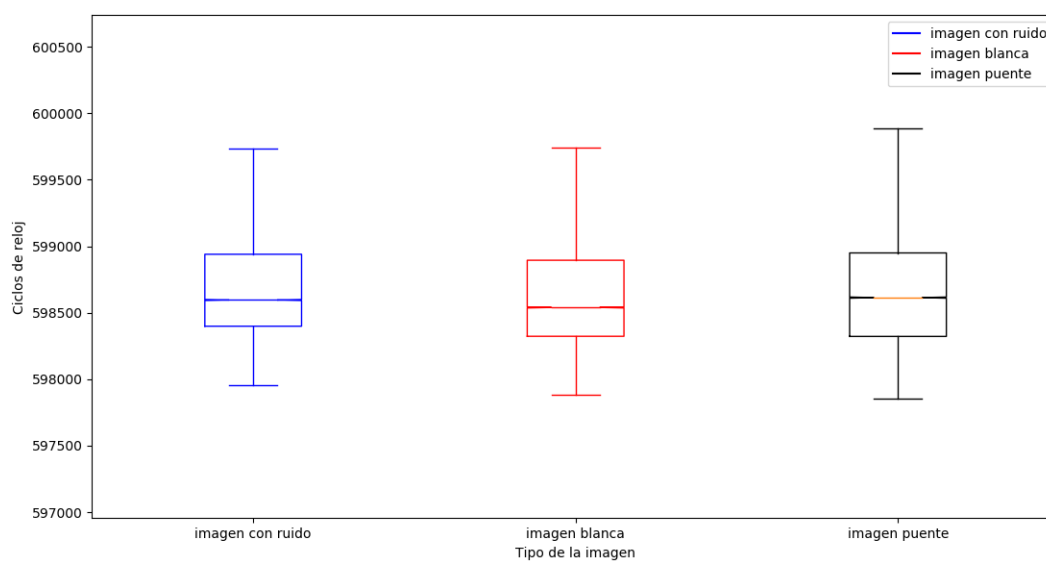


Figura 6: Experimento contenido de la imagen

El segundo experimento consiste en correr el filtro con la misma foto en dos tamaños diferentes: una de 64x37 (2368 píxeles) y otra de 128x75 (9600 píxeles).

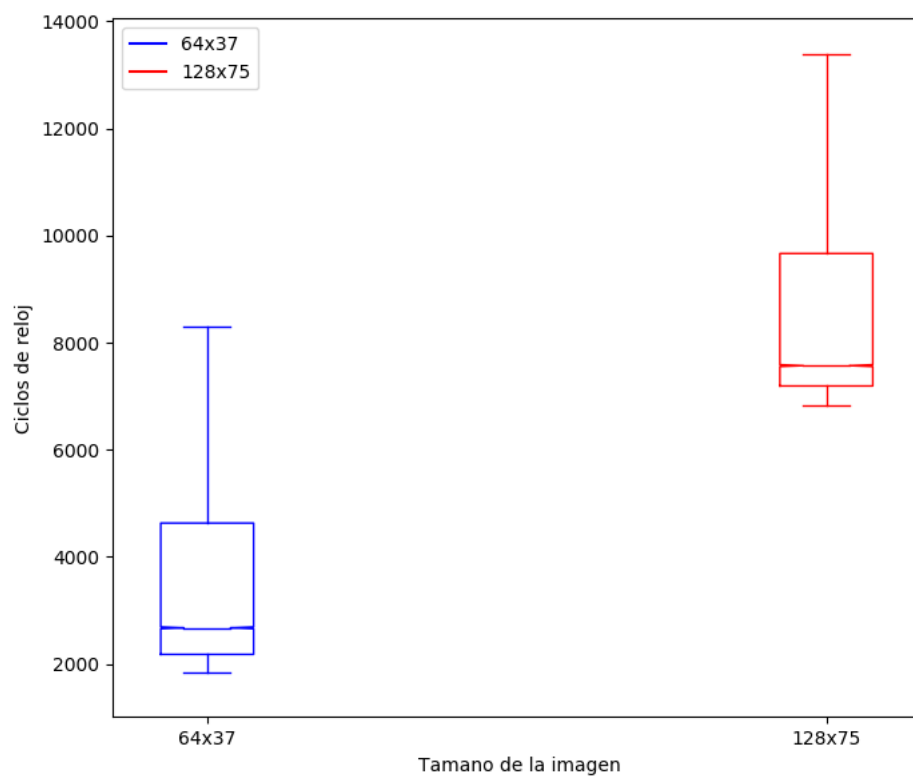


Figura 7: Experimento tamaño de la imagen

Análisis de resultados y conclusiones

Observamos gracias a la figura 6 que el contenido de la imagen no repercute en el tiempo de ejecución al aplicarle el filtro. Esto coincide a lo supuesto y se debe a que el código del filtro no está ramificado en función del contenido de los píxeles; no importa el parámetro pasado, las instrucciones ejecutadas son las mismas.

Respecto a la segunda parte del experimento, una foto posee aproximadamente el cuádruple de píxeles que la otra. Analizando los resultados exhibidos en el gráfico de la figura 7, vemos que el promedio de la cantidad de ciclos de clock insumidos por la corrida con la imagen grande no es cuatro veces la cantidad insumida por la corrida con la imagen chica, sino que es más bien tres veces ésta. Descartamos entonces nuestra hipótesis de proporcionalidad lineal entre el tamaño del parámetro y la performance temporal, pero verificamos la intuición de un aumento en el tiempo de ejecución al aumentar el tamaño de la imagen procesada. Podemos atribuir esta diferencia entre los resultados esperados y los obtenidos a los costos "fijos" al ejecutar el filtro, es decir aquellos que no dependen de los parámetros por ejemplo algunos accesos a memoria.

3.3. Rombos

Para el filtro "Rombos" decidimos llevar a cabo experimentos centrados en la cuestión del acceso a memoria. Buscamos responder a las siguientes preguntas; ¿cuán costoso resulta el acceso a memoria dentro del ciclo principal de una función? ¿Es significativa la diferencia en cuanto a la performance temporal del acceso a memoria de datos alineados en comparación con la del acceso a datos desalineados?

3.3.1. Experimento 1: Máscaras (no) alineadas

Cuando nos encontramos con la implementación del filtro de rombos, vemos que usa una gran cantidad de máscaras, al menos en comparación con los otros dos filtros. Éstas están predefinidas y escritas en memoria en la sección `.rodata` con sus respectivas etiquetas, para luego ser levantadas de memoria durante el procesamiento de los píxeles. Para este evento INTEL nos provee de dos instrucciones (entre varias), una que tiene como precondition que los datos a levantar de memoria estén alineados a dieciséis bytes mientras que la otra no, `MOVDQA` y `MOVDQU` respectivamente. Por esta razón se nos ocurrió experimentar con estas cuestiones: ¿Qué instrucción es más rápida? ¿Siempre se puede alinear la memoria? ¿Qué tan difícil es alinearla?

Hipótesis

Al ver que `MOVDQA` tiene una precondition, suponemos que tiene un trabajo de arquitectura que hace que el pasaje de la memoria al procesador sea más rápido, mientras que si uno levanta 16 bytes de memoria sin estar alineada, creemos que puede llegar a significar más de un "uso" del bus. Esto se debe a que el bus siempre levanta memoria de forma alineada, por lo que para levantar datos que no se encuentran alineados va a tener que buscarlos más de una vez. En cuanto a cuándo podemos alinear la memoria y cómo, suponemos que simplemente tendremos que dejar unos bytes con "basura" de tal manera que podamos guardar las máscaras de manera alineada.

Experimentos

Para comparar estas dos instrucciones vamos a tener dos implementaciones diferentes: una con las máscaras escritas de forma alineada que se aproveche de esto y utilice `MOVDQA`, y otra sin aclaraciones ni suposiciones sobre la alineación. A continuación observaremos los resultados obtenidos.

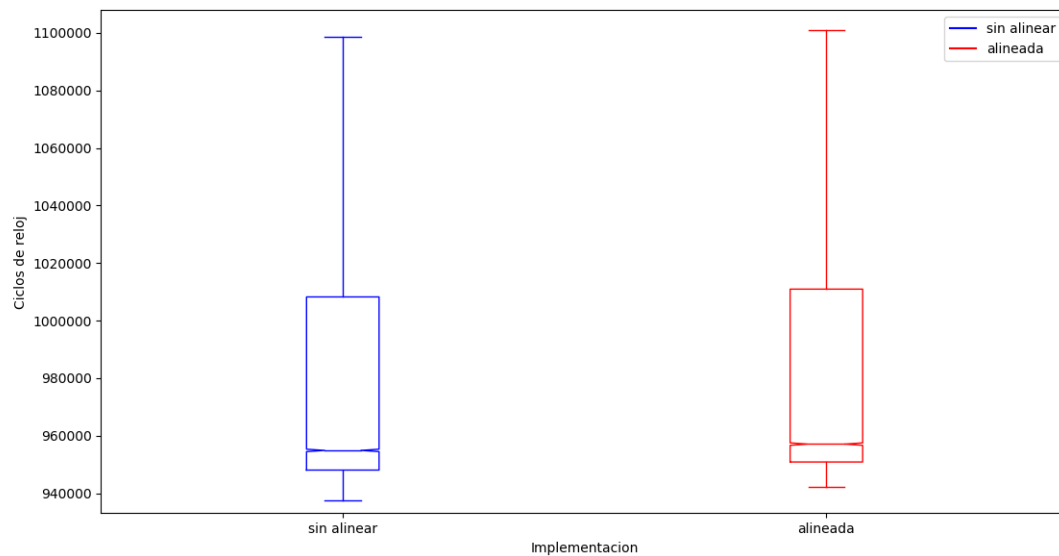


Figura 8: Experimento sobre la alineación de las máscaras

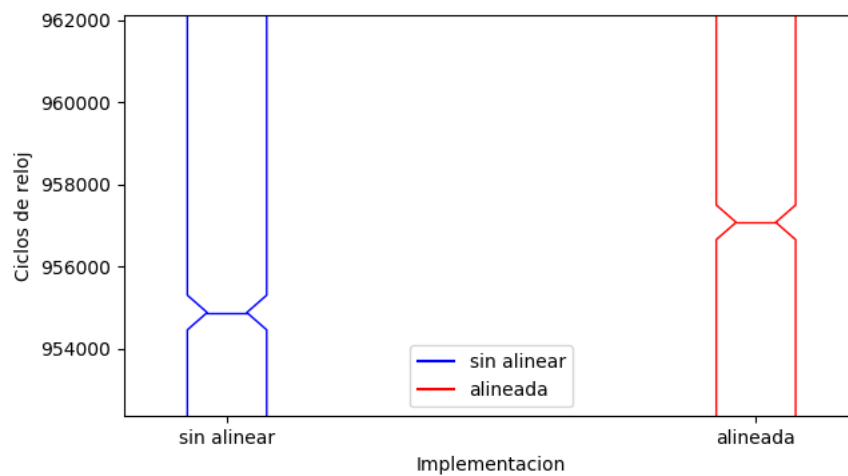


Figura 9: Medias de los resultados del experimento sobre la alineación de las máscaras

Análisis de resultados y conclusiones

Podemos observar en el gráfico que los resultados no fueron los esperados: de la figura 8 podemos concluir las dos implementaciones tardaron en general cantidades similares de ciclos de reloj, aunque si miramos con más atención, en la figura 9 vemos que la media de la implementación en la que la memoria se deja sin alinear es ligeramente menor que la que sí alinea la memoria. Sin embargo, esta diferencia no resulta significativa a fines de este análisis del tiempo de ejecución ya que siempre debemos tener en cuenta que el procesador podría estar realizando otras operaciones mientras se ejecuta nuestro programa.

Una razón por la cual pudo haber ocurrido que la alineación de la memoria no haya sido de gran influencia en la cantidad de ciclos de clock es que podría haber estado alineada, en cuyo caso no cambiaría la instrucción *"align 16"* al principio del código. Otra razón es que únicamente movemos datos de memoria a un registro al principio de la función, antes de iniciar el ciclo, y como traer cosas de memoria ya representa un costo grande en sí, la alineación no puede ser un cambio relevante.

3.3.2. Experimento 2: Accesos a memoria dentro y fuera del ciclo

Queremos analizar ahora el costo del acceso a memoria; ¿levantar datos en cada iteración representa un obstáculo significativo para la optimización temporal de una función?

Hipótesis

Suponemos que la cantidad de ciclos de clock de una función que accede a memoria dentro de su ciclo principal es mayor a una que realiza los accesos fuera del ciclo. Esto se debe no únicamente a que el número de instrucciones dentro del ciclo se ve incrementado por las operaciones de mover de memoria a algún registro, si no también a que el acceso a memoria representa un gran costo temporal, más aún si los datos buscados no se encuentran en la memoria cache.

Además, los accesos a memoria perjudican el buen funcionamiento del pipeline ya que en el caso de que se necesite un operando de memoria en la decodificación de una instrucción, el acceso a memoria para traer este operando interferirá con el fetch de la siguiente instrucción. De esta forma, en cada iteración la optimización buscada con pipelining se ve obstaculizada. Es decir que los accesos a memoria son incompatibles con el funcionamiento paralelo de los distintos bloques del procesador.

Experimentos

En la implementación original del filtro "Rombos" las máscaras son levantadas de memoria y almacenadas en registros XMM antes de iniciar el ciclo. Con el objetivo de analizar las ventajas o desventajas que presentan las distintas formas de levantar las máscaras creamos una nueva implementación del filtro, en la cual se levantan las máscaras en cada iteración del ciclo. Queremos comparar las dos implementaciones de forma tal que los cambios en eficacia temporal sean producto únicamente de la diferencia de acceso a las máscaras. Por este motivo comparamos la implementación que levanta las máscaras en el ciclo con la implementación del filtro original con 6 instrucciones adicionales en el ciclo. Estas instrucciones no modifican el funcionamiento del ciclo.

Comparamos entonces la performance de las dos implementaciones obteniendo los siguientes resultados:

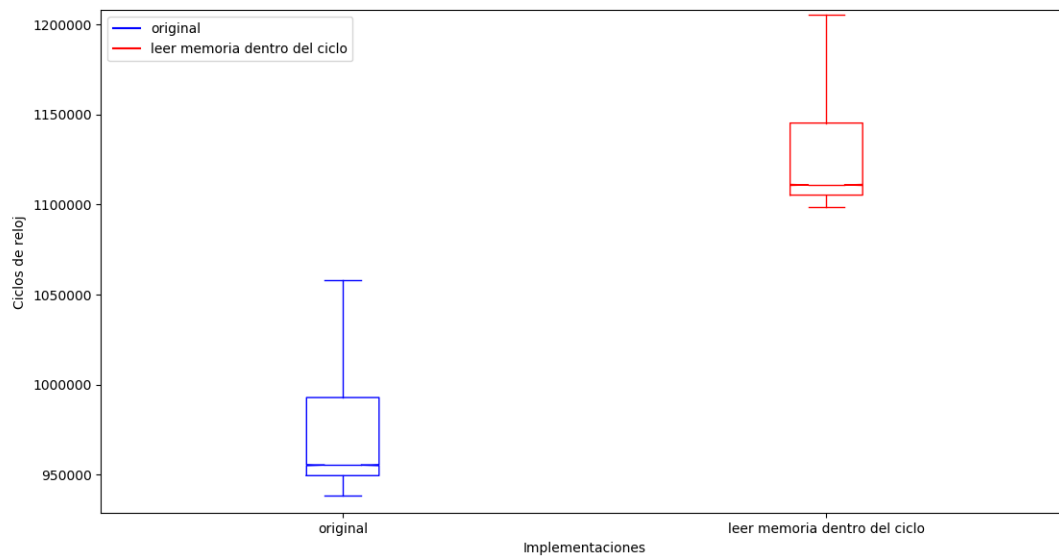


Figura 10: Resultado de levantar memoria en cada iteración del ciclo

Análisis de resultados y conclusiones

Como muestra la Figura 10, la implementación en la que se traen datos de memoria en cada una de las iteraciones del ciclo siempre está por arriba de la que trae de memoria una única vez antes de comenzar el ciclo, en lo que refiere a cantidad de ciclos de reloj. Es decir, no solo se confirmó nuestra hipótesis inicial de que acceder una mayor cantidad de veces a memoria implica una mayor cantidad de ciclos de reloj, si no que además podemos observar que este incremento en el tiempo que tarda en ejecutarse es realmente alto en la mayoría de los casos. Esto se genera debido a que los accesos a memoria representan uno de lo que más tiempo de manda en cuanto a características de un programa (por eso la memoria cache es un tema de gran relevancia, ya que permite reducir este costo). En conclusión, por más que los datos puedan llegar a encontrarse en la memoria cache, igualmente acceder a los datos que se encuentran en memoria (ya sea en la principal o en la cache) significa un gran costo para el tiempo de ejecución.

4. Conclusión

A lo largo del trabajo estudiamos las implementaciones de distintos filtros para imágenes que trabajan con SIMD, es decir que procesan datos de forma simultánea. Analizamos cómo y en qué medida diferentes factores influyen en el tiempo de ejecución de los códigos.

Luego de realizar las experimentaciones presentadas en este informe, podemos concluir que hay ciertas características de los programas que afectan su desempeño temporal en mayor medida, como lo son el hecho de que haya sido implementado en C o en ASM, la cantidad de accesos a memoria o la cantidad de saltos.

Por un lado, analizando el costo temporal de los saltos condicionales y comparándolo con el costo de los llamados a funciones pudimos darnos cuenta de que efectivamente cambiando la implementación a una en la que se realizaban la mitad de los saltos se reducía considerablemente la cantidad de ciclos de clock. Esto no nos sorprendió ya que suponíamos que era un factor que generaría cambios en la cantidad de ciclos de reloj que lleva la ejecución al influir en gran parte en el proceso de pipelining.

También pudimos constatar que en algunos filtros implementados el desempeño temporal no depende de la foto en sí y que existe una fuerte relación entre el tamaño de la imagen y el costo temporal de aplicar el filtro, sin embargo esta relación no es lineal debido al costo fijo implicado por la ejecución del filtro.

Comparando los experimentos realizados pudimos observar que no solo cuantos más accesos a memoria, mayor es el tiempo de ejecución si no que también es, dentro de las características analizadas, una de las que más influyen en la performance del programa. Esto también fue coherente con lo que creíamos ya que la optimización de los accesos a memoria es un tema de interés: es uno de los incentivos que llevaron a incluir la memoria cache en las computadoras.

Otro hecho que generó grandes cambios en la cantidad de ciclos de reloj que toma el programa en ejecutar fue el cambio de la implementación de C a ASM. Vimos que la implementación en C tardó mucho más que la de ASM. Esto en parte se debe a la forma en la que se compilan los códigos pero, en mayor medida, se debe a que las instrucciones de SIMD nos permiten procesar varios datos simultáneamente, mientras que en C modificamos de a un dato por iteración.