

SIMD

Nadia Heredia

Organización del Computador II, 1er cuat. de 2014

- SSE (Streaming SIMD Extensions) es un set de instrucciones que implementa el **modelo de cómputo SIMD** (una misma instrucción para varios datos a la vez).
- Existen 16 **registros de 128 bits** (XMM0...XMM15).
- Varios posibles **tipos de datos: enteros, flotantes de precisión simple** (32 bits) y **flotantes de precisión doble** (64 bit).
- **Formas de operar:**
 - $P = \text{Packed}$ (**E**mpaquetado / **P**aralelo).
 - $S = \text{Scalar}$ (**E**scalar).
- Sufijos para las instrucciones:
 - SS: scalar single-precision (float)
 - SD: scalar double-precision (double)
 - PS: packed single-precision (float)
 - PD: packed double-precision (double)
- Las instrucciones que comienzan con P sólo son para operar sobre enteros (ejemplo PADDQ). Las otras, sin P , son para operar sobre flotantes (ejemplo ADDPS).

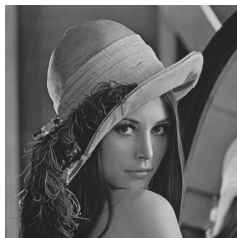
Volviendo al ejemplo del aumento de brillo, teníamos el siguiente código en **C**:

Incremento de brillo

```
for (int i = 0; i < N; i++)  
    for (int j = 0; j < M; j++)  
        I[i][j] += b
```

¿Qué pasa si $I[i][j] + b$ es mayor que 255? Por ejemplo, con $b = 100$

Ejemplo de incremento de brillo sin saturación



(a) Original, brillo ($b = 50$), brillo ($b = 100$)

Problemas con la no saturación

El problema es que el que el resultado de la suma no entra en el registro.

- Se trunca el resultado.
- Se pasa de colores más claros (cercanos al 255) a colores más oscuros (cercanos al 0).

Para evitar esto, habría que hacer

Incremento de brillo con saturación

```
for (int i = 0; i < N; i++)  
  for (int j = 0; j < M; j++)  
    I[i][j] = min(255, I[i][j] + b)
```

Entonces, obtenemos

Ejemplo de incremento de brillo con saturación



(b) Original, brillo saturado ($b = 50$), brillo saturado ($b = 100$)

Como esto es algo común en el procesamiento de señales, existen instrucciones específicas para operar de esta manera:

- *PADDUSB/PSUBUSB*: Suma/Resta enteros sin signo con saturación sin signo.
- *PADDSB/PSUBSB*: Suma/Resta enteros con signo con saturación con signo.

En este caso la suma es de a **byte**, también está para **word**.

Supongamos que queremos pasar una imagen a escala de grises. Una forma de hacerlo es mediante la fórmula:

$$f(r, g, b) = \frac{1}{4} \cdot (r + 2g + b)$$

- En esta situación, es necesario manejar los resultados intermedios en un tipo de datos de mayor precisión para no perder información en los cálculos.

La precisión original es de **byte**.

- Lo primero que tenemos que hacer es pasar los datos de **byte** a algo más grande (en este caso nos alcanza con pasar a **word**).
- ¿Cómo aumentamos la precisión?
Utilizando las instrucciones de desempaquetado.

- Si tenemos un registro con números sin signo:

$$\mathbf{xmm0} = a_{15} \mid a_{14} \mid \dots \mid a_1 \mid a_0$$

- Sólo necesitamos agregar cero delante de cada número. Es decir,

$$\mathbf{xmm0} = 0 \mid a_7 \mid \dots \mid 0 \mid a_0$$

$$\mathbf{xmm1} = 0 \mid a_{15} \mid \dots \mid 0 \mid a_8$$

Concretamente:

```
pxor xmm7, xmm7      ; xmm7 = 0 | 0 | ... | 0  
  
movdqu xmm2, xmm1     ; xmm2 = xmm1 = a15 | a14 | ... | a1 | a0  
  
punpcklbw xmm1, xmm7  ; xmm1 = 0 | a7 | ... | 0 | a0  
punpckhbw xmm2, xmm7  ; xmm2 = 0 | a15 | ... | 0 | a8
```

Ahora cada dato es de tipo **word**.

Después de extender los datos, realizamos las operaciones que necesitamos. Finalmente, tenemos que guardar los datos (que en este caso representarían los píxeles de una imagen en escala de grises), por lo que tenemos que volver a convertirlos a **byte**

- ¿Cómo hacemos la conversión?
Empaquetando los datos.

Las instrucciones de empaquetamiento son varias, tienen en cuenta distintos tipos de datos y si los datos tienen signo o no. Por ejemplo, en el caso concreto de **byte** tenemos:

- *packsswb* (saturación con signo)
- *packuswb* (saturación sin signo)

Desempaquetado $\text{punpck} + l/h + b/w/d/q + w/d/q/dq$

Empaquetado $\text{pack} + ss/us + w/d + b/w$

Tener en cuenta

- si queremos usar la parte alta o la parte baja
- los tipo de datos con los que trabajamos
- el tipo de saturación a usar (cuando empaquetamos)

Concretamente, si:

$$\mathbf{xmm1} = 0 \mid a_7 \mid \dots \mid 0 \mid a_0$$

$$\mathbf{xmm2} = 0 \mid a_{15} \mid \dots \mid 0 \mid a_8$$

Entonces:

$$\text{packuswb xmm1, xmm2} \quad ; \text{ xmm1} = a_{15} \mid a_{14} \mid \dots \mid a_1 \mid a_0$$

Ahora cada dato es de tipo **byte**

Resumiendo, la forma de trabajar en estos casos es

- Leer datos a procesar.
- Extender precisión (unpack).
- Realizar las cuentas necesarias.
- Volver a la precisión original (pack).
- Guardar datos procesados.

Comparación

En **SSE** también existen instrucciones de comparación, aunque se comportan un poco diferente a lo que veníamos usando.

Por ejemplo, supongamos que estamos trabajando con **words** y queremos saber cuáles de ellos son menores a cero.

Si tenemos los datos en `xmm1`, por ejemplo:

`xmm1` = 1000 | - 456 | - 15 | 0 | 100 | 234 | - 890 | 1

```
pxor xmm7, xmm7      ; xmm7 = 0 | 0 | ... | 0
pcmpgtw xmm7, xmm1    ; xmm7 > xmm1 ?
```

Obtenemos el resultado de la comparación en `xmm7`:

`xmm7` =
0x0000 | 0xFFFF | 0xFFFF | 0x0000 | 0x0000 | 0x0000 | 0xFFFF | 0x0000

Es decir, compara **word** a **word** y si se cumple la condición setea en **unos** (0xFFFF en este caso) el resultado, o en **ceros** si no.

Como vimos en el ejemplo, la comparación nos devuelve un registro con unos y ceros.

Podemos usar este registro

- Para extender el signo de números signados al pasar de un tipo de datos a otro.

Usando el hecho de que los números negativos se van a extender con unos, y los positivos con ceros.

- $-5 = 1011$ se extiende a $1111\ 1011$
 - $5 = 0101$ se extiende a $0000\ 0101$
- A modo de máscaras, acompañado de instrucciones como: *PAND*, *POR*, etc.

Comparación - Extensión de signo

Para extender el signo del registro del ejemplo anterior

xmm1 = 1000 | - 456 | - 15 | 0 | 100 | 234 | - 890 | 1

Si tenemos el resultado de la comparación en xmm7, hacemos

```
movdqu xmm2, xmm1      ; copio xmm1
punpckhwd xmm1, xmm7    ; xmm1 = 1000 | ... | 0
punpcklwd xmm2, xmm7    ; xmm2 = 100 | ... | 1
```

Comparación - Uso de máscaras

Supongamos ahora que se desea sumar 3 a los números menores a 0.
Al igual que antes, tenemos los números en `xmm1`, y el resultado de la comparación en `xmm7`.

Tenemos el número 3, que queremos sumar, en un registro

$$\text{xmm2} = 3 \mid 3 \mid \dots \mid 3$$

Hacemos

```
pand xmm7, xmm2    ; xmm7 = 0 | 3 | ... | 3 | 0  
paddw xmm1, xmm7    ; xmm1 = 1000 | - 453 | ... | - 887 | 1
```

Ejercicio 1

Calcular la distancia máxima entre puntos.

- Los puntos (x,y) están almacenados como dos números contiguos de punto flotante de precisión simple.
- Los vectores tienen n puntos, donde n es un short sin signo.
- n es múltiplo de 2.
- Se debe procesar de a dos puntos en paralelo.
- La distancia se calcula entre los puntos de v y los de w , elemento a elemento.
- El prototipo de la función es:

```
float maximaDistancia(float* v, float* w, unsigned short n);
```

Distancia

$$\text{dist}((x,y), (x',y')) = \sqrt{(x-x')^2 + (y-y')^2}$$

Obtenemos los parámetros, chequeamos que los vectores no estén vacíos y preparamos los registros para el ciclo

```
global maximaDistancia
```

```
section .text
```

```
maximaDistancia:      ; rdi = v, rsi = w; dx = n
```

```
xor rcx, rcx          ; obtener los parámetros
```

```
mov cx, dx            ; rcx = n
```

```
cmp rcx, 0            ; chequear si los vectores están vacíos
```

```
je fin
```

```
shr rcx, 1            ; rcx = n / 2
```

```
; limpiar xmm0 para usar como temporal
```

```
xorps xmm0, xmm0      ; xmm0 = 0 | 0 | 0 | 0
```

ciclo:

movups xmm1, [rsi]	; xmm1 = $y_2 \mid x_2 \mid y_1 \mid x_1$
movups xmm2, [rdi]	; xmm2 = $y_2' \mid x_2' \mid y_1' \mid x_1'$
subps xmm1, xmm2	; xmm1 = $(y_2 - y_2') \mid (x_2 - x_2') \mid (y_1 - y_1') \mid (x_1 - x_1')$
mulps xmm1, xmm1	; xmm1 = $(y_2 - y_2')^2 \mid (x_2 - x_2')^2 \mid$ $(y_1 - y_1')^2 \mid (x_1 - x_1')^2$
movdqu xmm2, xmm1	
psrldq xmm2, 4	; xmm2 = $0 \mid (y_2 - y_2')^2 \mid (x_2 - x_2')^2 \mid (y_1 - y_1')^2$
addps xmm1, xmm2	; xmm1 = $* \mid ((x_2 - x_2')^2 + (y_2 - y_2')^2) \mid$ $* \mid ((x_1 - x_1')^2 + (y_1 - y_1')^2)$
sqrtps xmm1, xmm1	; xmm1 = $* \mid \sqrt{(x_2 - x_2')^2 + (y_2 - y_2')^2} \mid$ $* \mid \sqrt{(x_1 - x_1')^2 + (y_1 - y_1')^2}$
maxps xmm0, xmm1	; xmm0 = $* \mid \max(\text{dist_impares}) \mid$ $* \mid \max(\text{dist_pares})$
add rsi, 16	; avanzar los punteros
add rdi, 16	
loop ciclo	

Tenemos ahora dos máximos (pares e impares). Calculamos cuál de los dos es el mayor. Y terminamos la función.

xmm0 = * | max(dist_impares) | * | max(dist_pares)

<code>movdqu xmm1, xmm0</code>	<code>; xmm1 = xmm0</code>
<code>psrldq xmm1, 8</code>	<code>; xmm1 = 0 0 * max(dist_impares)</code>
<code>maxps xmm0, xmm1</code>	<code>; xmm0 = * * * max(dist_pares, dist_impares)</code>
<code>ret</code>	<code>; retornar</code>

Para hacer:

- ¿Qué pasa si ahora tenemos que comparar todos los puntos de v contra todos los de w ?
- ¿Y si n no fuera múltiplo de 2?
- ¿Podemos procesar más de 2 puntos en paralelo?

Normalizar un vector.

- La longitud del vector es **n**, donde n es un int.
- n es múltiplo de 4.
- Se debe devolver un vector, de manera que el máximo elemento sea uno, y el mínimo cero.
- El prototipo de la función es:

```
float* normalizarVector(float* v, int n);
```

Debemos buscar el máximo y el mínimo, y aplicar a cada elemento

$$x[i] = \frac{(x[i] - \text{min})}{(\text{max} - \text{min})}$$

Obtenemos los parámetros y pedimos memoria para el vector a devolver.

```
extern malloc                ; rdi = v, esi = n
global normalizarVector
section .text
normalizarVector:
push rbp
mov rbp, rsp                ; armar el stack frame
push r12
push r13

mov r12, rdi                ; obtener los parámetros
xor r13, r13
mov r13d, esi

mov rdi, r13                ; rdi = n
shl rdi, 2                  ; rdi = n * 4
call malloc                 ; pedir memoria
```

Preparamos los registros para el ciclo, y buscamos dentro del vector el máximo y el mínimo.

```
mov rcx, r13          ; rcx = n
shr rcx, 2            ; rcx = n / 4

mov rdi, r12          ; rdi = v
movups xmm0, [rdi]    ; setear valores iniciales para los máximos y mínimos
movups xmm1, [rdi]

ciclo:
  movups xmm2, [rdi]  ; leer los próximos 4 valores
  minps xmm0, xmm2    ; actualizar mínimos
  maxps xmm1, xmm2    ; actualizar máximos

  add rdi, 16         ; avanzar el puntero
  loop ciclo
```

Conseguimos el mínimo de `xmm0` y lo ponemos en cada posición

movdqu xmm2, xmm0	; $xmm2 = xmm0$
psrldq xmm2, 4	; $xmm2 = 0 \mid xmm0_3 \mid xmm0_2 \mid xmm0_1$
minps xmm0, xmm2	; $xmm0 = * \mid \min(xmm0_2, xmm0_3) \mid$ $* \mid \min(xmm0_0, xmm0_1)$
movdqu xmm2, xmm0	
psrldq xmm2, 8	; $xmm2 = 0 \mid 0 \mid * \mid \min(xmm0_2, xmm0_3)$
minps xmm2, xmm0	; $xmm2 = * \mid * \mid * \mid \min(xmm0_0, xmm0_1, xmm0_2, xmm0_3)$
xorps xmm0, xmm0	; $xmm0 = 0 \mid 0 \mid 0 \mid 0$
addss xmm0, xmm2	; $xmm0 = 0 \mid 0 \mid 0 \mid \min$
pslldq xmm0, 4	; $xmm0 = 0 \mid 0 \mid \min \mid 0$
addss xmm0, xmm2	; $xmm0 = 0 \mid 0 \mid \min \mid \min$
pslldq xmm0, 4	; $xmm0 = 0 \mid \min \mid \min \mid 0$
addss xmm0, xmm2	; $xmm0 = 0 \mid \min \mid \min \mid \min$
pslldq xmm0, 4	; $xmm0 = \min \mid \min \mid \min \mid 0$
addss xmm0, xmm2	; $xmm0 = \min \mid \min \mid \min \mid \min$

Realizamos algo similar para el máximo y le restamos el mínimo al final

movdqu xmm2, xmm1	; xmm2 = xmm1
psrldq xmm2, 4	; xmm2 = 0 xmm1 ₃ xmm1 ₂ xmm1 ₁
maxps xmm1, xmm2	; xmm1 = * min(xmm1 ₂ , xmm1 ₃) * max(xmm1 ₀ , xmm1 ₁)
movdqu xmm2, xmm1	
psrldq xmm2, 8	; xmm2 = 0 0 * max(xmm1 ₂ , xmm1 ₃)
maxps xmm2, xmm1	; xmm2 = * * * max(xmm1 ₀ , xmm1 ₁ , xmm1 ₂ , xmm1 ₃)
xorps xmm1, xmm1	; xmm1 = 0 0 0 0
addss xmm1, xmm2	; xmm1 = 0 0 0 max
pslldq xmm1, 4	; xmm1 = 0 0 max 0
addss xmm1, xmm2	; xmm1 = 0 0 max max
pslldq xmm1, 4	; xmm1 = 0 max max 0
addss xmm1, xmm2	; xmm1 = 0 max max max
pslldq xmm1, 4	; xmm1 = max max max 0
addss xmm1, xmm2	; xmm1 = max max max max
subps xmm1, xmm0	; xmm1 = diferencia entre el máximo y el mínimo 4 veces

Llenamos las posiciones del nuevo vector, asignamos el valor de retorno, y restauramos el stackframe

```
mov rcx, r13          ; rcx = n
shr rcx, 2             ; rcx = n / 4
mov rdi, r12           ; rdi = v
mov rsi, rax           ; rsi = vector nuevo
```

ciclo2:

```
movups xmm2, [rdi]    ; xmm2 =  $x_{i+3} \mid x_{i+2} \mid x_{i+1} \mid x_i$ 
subps xmm2, xmm0       ; xmm2 =  $x_{i+3} - \text{min} \mid \dots \mid \dots \mid x_i - \text{min}$ 
divps xmm2, xmm1       ; xmm2 =  $\frac{x_{i+3} - \text{min}}{\text{max} - \text{min}} \mid \dots \mid \dots \mid \frac{x_i - \text{min}}{\text{max} - \text{min}}$ 
movups [rsi], xmm2     ; copiar valores normalizados al nuevo vector
```

```
add rdi, 16            ; avanzar punteros
add rsi, 16
loop ciclo2
```

```
pop r13
pop r12
pop rbp
ret                    ; retornar
```

¿Preguntas?