

PAP

Parámetros en los computadores paralelos

- Tipo y número de procesadores
 - ↳ Procesador de grano fino (operaciones simples)
 - ↳ " " grueso (operaciones complejas)
- Presencia y ausencia de mecanismos de control
- Funcionamiento síncrono y asíncrono
- Formas de comunicación entre procesos

Paralelismo interno:

Queda oculto a la arquitectura del computador

Paralelismo explícito:

Aquel que queda visible al usuario

* Clasificación de Hwang-Briggs

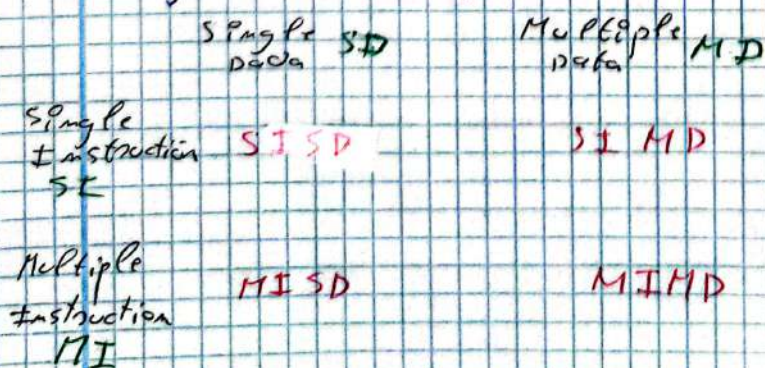
Establecen una primera aproximación a las clases de computadores paralelos fijando 3 configuraciones básicas:

- Computadores pipeline
- " Matriciales
- Sistemas multiprocesador

De esta forma se pueden distinguir 2 tipos de paralelismo:

- temporal (segmentación)
 - Aquellos que se les solapan varias operaciones en un instante
- Espacial (Existencia de procesos replicados)
 - Síncrono
 - Asíncrono

* Clasificación Flynn



SISD

- Modelo Von Neuman (secuencial)
- Instrucciones:
De memoria a procesador
- Datos:
Entre memoria y procesador

SIMD

- Una única unidad de control (vectorial)
- Misma instrucción se ejecuta simultáneamente en todas las unidades de procesamiento

MISD

Utilizado por la NASA para comprobar los errores al lanzar un cohete

MIMD

Cada procesador ejecuta un programa diferente

• CPU

- Dominadas por memorias cache
- La latencia a estas memorias es importante
- Ejecutan instrucciones complejas

• Pocos núcleos

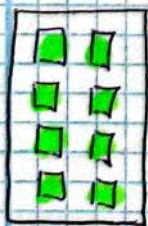
• GPU

- Poca importancia de la latencia
- Carga de trabajo sobre pixel shader y vertex shader es relativamente baja

• Muchos núcleos

CPU (- Cores + cache + Control)

GPU (+ Cores - Cache - Control)



Streaming procesador //

Scalar procesador

1 hilo = 1 scalar procesador

1 WARP = 32 hilos

- L +
- A • Global memory
- T • Local memory
- N • Constant memory
- C • Shared memory (sobre la que se puede leer y escribir)
- I • Registro

*
$$\text{Ancho de banda} = \frac{\text{Frecuencia memoria} \cdot \text{Ancho bus de memoria}}{8}$$

- Frecuencia memoria = 10 GB/s
- bus memoria = 256 bits

$$AB = \frac{10 \cdot 256}{8} = 320 \text{ GB/s}$$

- Frecuencia memoria = 1600 MHz = 1'6 GHz
- bus memoria = 512 bits

$$AB = \frac{16 \cdot 2 \cdot 512}{8} = 2048 \text{ GB/s}$$

• Cauda gráfico, proceso de renderización

- A la GPU le llega una lista de primitivas (triángulos), que tienen la info asociada a peso, color, ...
- Pasa al procesador de vértices y aplica la geometría (puede hacer pequeños cambios)
- Hace el rasterizado (pasar imágenes 3D a 2D) en 2 fases
 - Culling (Marcar triángulos que puede eliminar porque están fuera de la cámara)
 - Culling (eliminar los triángulos marcados previamente)
- Transforma los triángulos en tramas de píxeles
- Se pasa al procesador de píxeles
- Se prepara la imagen para sacarla por pantalla (blending compone la imagen)

¿Que es un renderizado?

- Proceso de transformación de una imagen de 3 dimensiones a 2

¿Que es el rasterizado?

Conversión de triángulos a fragmentos

Anti-Aliasing

Reduccion de los dientes de sierra de las imágenes

La GPU utiliza Single Program Multiple Data stream

(un programa se ejecuta en múltiples hilos que acceden a diferentes flujos de datos)

1 bloque = 512 hilos

(Máximo 8 bloques por SM)

2 o más Bloques = 768 hilos

ej.

• Si lanzo una matriz $32 \times 32 = 1024$ hilos necesito

▮

Necesitare 2 bloques SM (1 de 512 y otro de 512)

¿Cuántas matrices de 256 hilos puedo ejecutar en 1 SM?

$$\hookrightarrow \frac{768}{256} = 3$$

ej.

$16 \times 16 = 256$ hilos, si tengo $\frac{768}{256} = 3$ matrices por SM



ej.

$8 \times 8 = 64$ hilos en un SM puedo ejecutar

Matriz Bloque

Ocupación $\frac{1}{3}$ (lo máximo es 768)

$$16 \cdot 16 = 256 \text{ hilos} \cdot 8 \text{ bloques (no entra)}$$

↓

$$\frac{768}{256} = 3 \text{ Bloques Máximo puedo lanzar}$$

$$\frac{312 \cdot 512}{256} = 1024 \text{ Bloques de hilos de 256}$$

$$\text{Cada SM puede 3 Bloques} \Rightarrow \frac{1024}{3} = 342 \text{ SM necesito}$$

ej:

Matriz $40 \cdot 40$ (a máxima ocupación)

¿Cuántos SM necesito utilizar?

$$40 \cdot 40 = 1600 \text{ hilos}$$

Usaré una matriz $16 \cdot 16 = 256 \text{ hilos}$

$$\frac{768}{256} = 3 \text{ bloques por SM}$$

↓

$$\frac{1600}{256} = 6.25 \text{ bloques de hilos necesito}$$

↓

$$\left\lceil \frac{6.25}{3} \right\rceil = 3 \text{ SM necesarios}$$

Memoria

- Global 768 MB - 1024 MB
- Local 16 KB
- Compartida 16 KB
- Registros 8 KB
- Constante 64 KB
- Texturas 8 KB

Compilación

- Pasamos el fichero .cu al compilador
- Separa lo que va a ejecutar la CPU y la GPU

CPU { - Lo de la CPU pasa al compilador y devuelve un fichero .cpp

GPU { - Crea un fichero en ensamblador para la parte GPU .s
- Este indica que procesos son paralelos
- se lo pasa al OCC que unifica este fichero con las características de nuestra tarjeta
- Devuelve un fichero .SASS que es lo que ejecutará sobre la GPU

- Para acceder a un hilo en concreto de un bloque
 $N^{\circ} \text{ bloque} \cdot N^{\circ} \text{ hilos por bloque} + \text{threadIdx.x}$

$$\text{blockIdx.x} \cdot \text{blockDim.x} + \text{threadIdx.x}$$

- bloques 2 Dim, hilos 3 dim

grid(10, 21, 1) ✓
block(1, 2, 3)

grid(10, 21, 2) ✗
block(1, 2, 3)

- 16 KB por SM de memoria compartida
- Matriz de 1024 hilos y consume 32 KB de memoria compartida
¿cuántos SM necesito?

Opciones

- 3 Bloques de 256 en 1 SM y 1 Bloque de 256 en otro bloque
- 1 \times 3
- 2 \times 2

$$\frac{32 \text{ KB}}{4 \text{ Bloques}} = 8 \text{ KB por Bloque (cada SM tiene 16 KB)}$$

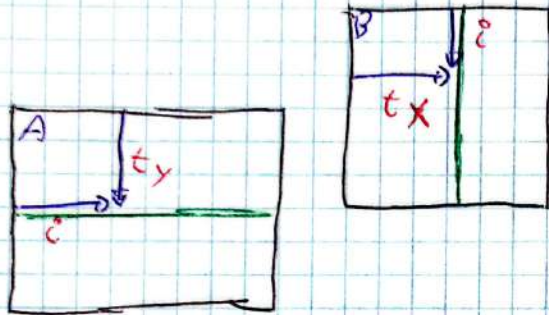


Solo puedo hacer 2 y 2


```

int valor = 0;
for (int i = 0; i < width; i++) {
    int primeiraMatriz = a[threadIdx.y * width + i];
    int segundaMatriz = b[i * width + threadIdx.x];
    valor += primeiraMatriz * segundaMatriz;
}
c[threadIdx.y * width + threadIdx.x] = valor;

```



```

int Row = BlockIdx.y * TileWidth + threadIdx.y;
int Col = BlockIdx.x * TileWidth + threadIdx.x;

int valor = 0;
for (int i = 0; i < width; i++) {
    valor += a[Row * width + i] * b[i * width + Col];
}
c[Row * width + Col] = valor;

```

```

dim3 dimGrid (width / tileWidth, width / tileWidth);
dim3 dimBlock (tileWidth, tileWidth);

<<< dimGrid, dimBlock >>> ( ... );

```