

Práctico 3 - GPGPU

Santiago Calvo

Estudiante de Ingeniería en Computación

5.055.578-2

Montevideo, Uruguay

santiago.calvo.vello@fing.edu.uy

Resumen—Informe de los resultados de la realización del práctico 3 de la materia Computación de Propósito General en Unidades de Procesamiento Gráfico.

I. INTRODUCCIÓN

Con el objetivo de realizar un estudio de la performance de distintos algoritmos utilizando CUDA, se realizan diversas implementaciones con accesos a memorias distintos.

II. INSTRUCCIONES DE EJECUCIÓN Y COMPILACIÓN

La ejecución de los algoritmos se realizó sobre cluster.uy utilizando los archivos bash para cada ejercicio contenido en el entregable.

III. EJERCICIO 1

Para la generación del cuadrado se realizó la escritura utilizando el índice según bloque, su dimensión y el índice de thread. De esta manera se escribe en la dirección necesaria.

Luego para el llamado se dividió el espacio en bloques de 32×32 teniendo bloques de tamaño 1024. Para la grilla se divide en el $\frac{\text{número_de_puntos}}{32}$ en cada eje (x,y) cubriendo todas las celdas.

Luego se utilizó bloques y grillas de igual distribución para la generación de la matriz con función *seno* de los puntos.

IV. EJERCICIO 2

Para la generación del cubo se realizó la escritura utilizando el índice según bloque, su dimensión y el índice de thread. De esta manera se escribe en la dirección necesaria.

Luego para el llamado se dividió el espacio en bloques de $8 \times 8 \times 8$ teniendo bloques de tamaño 512. Posteriormente se divide la grilla en $\frac{\text{número_de_puntos}}{8}$ en cada eje (x,y,z) cubriendo todas las celdas.

Luego se utilizó bloques y grillas de igual distribución para la generación del cubo con función *tangente* de los puntos.

V. EJERCICIO 3

Este ejercicio trataba sobre el desarrollo de una función con tres órdenes de acceso distintas.

- 1. $S(y, x) = \sum_{z=-\pi}^{\pi} f(x, y, z)$
- 2. $S(z, x) = \sum_{y=-\pi}^{\pi} f(x, y, z)$
- 3. $S(z, y) = \sum_{x=-\pi}^{\pi} f(x, y, z)$

Estas presentan una diferencia importante en los tiempos de ejecución como se puede ver en la figura V. La 3er forma de acceso tiene un tiempo de ejecución de 125ms en contraposición a los 16ms de las otras funciones.

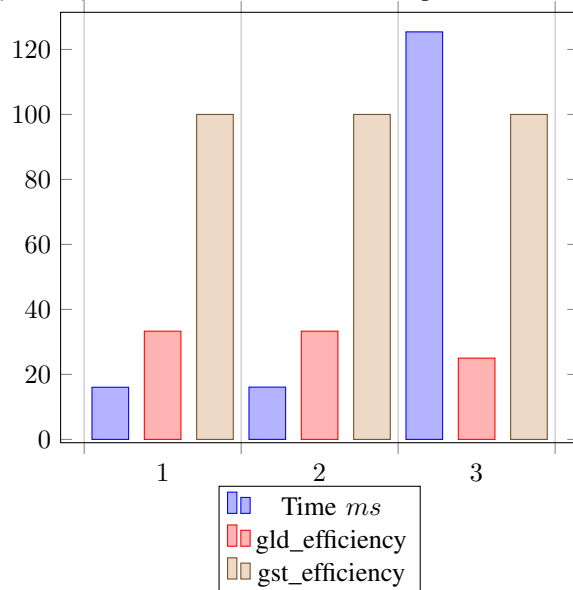
Estos se debe deber a los accesos globales a memoria. Como se puede ver, al utilizar *nvprof* se obtienen los valores de eficiencia de estos accesos nombrados:

- *gld_efficiency*: Relación de lecturas globales a memoria con lecturas requeridas
- *gst_efficiency*: Relación de escrituras a memoria con escrituras a memoria requeridas

En la figura V se puede ver como las lecturas tienen una particularmente baja en el tercer método lo cual se manifiesta como menor tiempo de ejecución dado que los datos deben ser obtenidos de memoria global cada vez.

Al acceder a direcciones contiguas se aprovecha el principio de localidad espacial permitiendo aprovechar las lecturas, necesitando luego menos accesos globales. En el caso del tercer algoritmo no se está aprovechando a su máximo los accesos dado que las direcciones accedidas están "espagadas" por la manera en que se organiza el cubo del cual se leen los datos.

Los accesos por el eje Z e Y no permiten aprovechar el uso de la memoria dado su separación. Tienen distancia de *fila* y *fila x fila* entre los dos elementos respectivamente.



VI. EJERCICIO 4

En este ejercicio se realizó un algoritmo de tipo stencil. Este algoritmo realiza lecturas de múltiples elementos dentro de un mismo thread para realizar un cálculo final. Por esto se verán los accesos a memoria realizados por este.

Utilizando *nvprof* se visualizó la cantidad de transacciones que se realizaron a memoria global, visualizables en la figura VI como la barra 1.

Esté número es coherente con la cantidad de transacciones esperadas aunque levemente menor.

$$\text{tamaño_fila} = \text{tamaño_columna} = \frac{2 \times \pi}{0,001} = 6280$$

$$\begin{aligned} \text{tamaño_matriz} &= \text{tamaño_fila} \times \text{tamaño_columna} \\ &= 39438400 \end{aligned}$$

$$\text{transacciones_esquinas} = 3 \times 4 = 12$$

$$\begin{aligned} \text{transacciones_edges} &= 2 \times 4 \times \\ &\quad (\text{tamaño_fila} + \text{tamaño_columna} - 4) \\ &= 100448 \end{aligned}$$

$$\begin{aligned} \text{transacciones_centro} &= (\text{tamaño_f} - 2) \times (\text{tamaño_c} - 2) \times 5 \\ &= 197066420 \end{aligned}$$

$$\begin{aligned} \text{transacciones} &= t_{\text{esquinas}} + t_{\text{edges}} + t_{\text{centro}} \\ &= 12 + 100448 + 197066420 \\ &= 197166880 \end{aligned}$$

Sin embargo esta cantidad de accesos podríamos aprovechar el uso de warps aprovechando una mayor cantidad de datos por transacción si sincronizamos los warps a la hora de leer de la siguiente manera:

```
__shared__ double matrix_point[32][32];
unsigned int threadIdx_x = threadIdx.x + 1;
unsigned int threadIdx_y = threadIdx.y + 1;

matrix_point[threadIdx_y - 1][threadIdx_x - 1] =
    matrix[i * num_points + j];

__syncwarp();
```

De esta forma guardamos una porción importante de la lectura en memoria Shared de mayor velocidad.

Al comparar ambas implementaciones (encontradas en *ej4.cu* y *ej4_improved.cu*), vemos una diferencia de transacciones y de acceso importante pues se realiza solo un acceso a memoria en vez de 5 por thread (reduciendo 5 veces las transacciones necesarias).

Nota, este ejercicio fue corrido en una maquina personal dado falta de acceso a ClusterUY. Las especificaciones son:

- CPU: Intel i7 3.6 GHz
- GPU: GTX 1060 3GB GDDR5
- RAM: 16 GB 2133 MHz

En las siguientes figuras se tiene las dos implementaciones, 1 sin optimizaciones y 2 con optimizaciones.

