



UNIVERSIDAD NACIONAL DE SAN AGUSTIN

ESCUELA PROFESIONAL DE
CIENCIA DE LA COMPUTACIÓN
ALGORITMOS PARALELOS

CUDA : Multiplicación de Matrices Modificado

Alumna :

Rosa Yuliana Gabriela

Paccotacya Yanque

Profesor:

Mg. Alvaro Hen-ry Mamani

Aliaga

Índice

| | |
|--|----------|
| 1. Multiplicación de matrices con tiles modificada | 2 |
| 1.1. Multiplicación de matrices con tiles modificada | 2 |
| 1.2. Código | 2 |
| 2. Características del Dispositivo | 3 |
| 3. Análisis de Factores limitantes del dispositivo | 4 |
| 3.1. Registros | 4 |
| 3.2. Memoria Compartida | 4 |
| 3.3. Otros | 4 |
| 4. Experimentos | 5 |
| 4.1. Tabla de resultados (ms) | 5 |
| 4.2. Capturas de Pantalla | 5 |
| 5. Conclusiones | 5 |

1. Multiplicación de matrices con tiles modificada

1.1. Multiplicación de matrices con tiles modificada

En la multiplicación de matrices con tiles simple, al momento de generar los productos se cargaban en memoria compartida, los mismos bloques, generandose asi redundancia. Para esto se modifica el código de manera que podamos aprovechar los bloques ya cargados en memoria, haciendo la multiplicación para dos elementos de la matriz resultante.

1.2. Código

```
1  __global__ void MatrixMulTiledMod(int * d_P, int * d_M, int* d_N,int
    Width)
2  {
3      __shared__ int Mds[TILE_WIDTH][TILE_WIDTH];
4      __shared__ int Nds[TILE_WIDTH][TILE_WIDTH];
5
6      __shared__ int Nds2[TILE_WIDTH][TILE_WIDTH];
7      int bx = blockIdx.x; int by = blockIdx.y;
8      int tx = threadIdx.x; int ty = threadIdx.y;
9      // Identify the row and column of the d_P element to work on
10     int Row = by * TILE_WIDTH + ty;
11     int Col = bx * TILE_WIDTH*2 + tx;
12     int Pvalue =0 , Pvalue2=0;
13     // Loop over the d_M and d_N tiles required to compute d_P element
14     for (int ph = 0; ph < Width/TILE_WIDTH; ++ph)
15     {
16         // Collaborative loading of d_M and d_N tiles into shared memory
17         if ((Row< Width) && (ph*TILE_WIDTH+tx)< Width)
18             Mds[ty][tx] = d_M[Row*Width + ph*TILE_WIDTH + tx];
19         if ((ph*TILE_WIDTH+ty)<Width && Col<Width)
20             Nds[ty][tx] = d_N[(ph*TILE_WIDTH + ty)*Width + Col];
21
22         if (( (ph*TILE_WIDTH+ty)<Width) && (Col+TILE_WIDTH<Width))
23             Nds2[ty][tx] = d_N[(ph*TILE_WIDTH + ty)*Width + Col+TILE_WIDTH];
24         __syncthreads();
25         for (int k = 0; k < TILE_WIDTH; ++k)
26         {
27             Pvalue += Mds[ty][k] * Nds[k][tx];
28             Pvalue2 += Mds[ty][k] * Nds2[k][tx];
29         }
30         __syncthreads();
31     }
32     d_P[Row*Width + Col] = Pvalue;
```

```

33     d_P[Row*Width + Col +TILE_WIDTH] = Pvalue2;
34 }

```

Código 1: Multiplicación de Matrices

2. Características del Dispositivo

Se trabajó en una máquina con tarjeta gráfica NVIDIA GT740-M y Cuda 8.0. Las características se detallan a continuación:

```

CUDA Device Query (Runtime API) version (CUDART static linking)

Detected 1 CUDA Capable device(s)

Device 0: "GeForce GT 740M"
  CUDA Driver Version / Runtime Version      8.0 / 8.0
  CUDA Capability Major/Minor version number: 3.5
  Total amount of global memory:              2004 MBytes (2101542912 bytes)
  ( 2) Multiprocessors, (192) CUDA Cores/MP:  384 CUDA Cores
  GPU Max Clock rate:                        1032 MHz (1.03 GHz)
  Memory Clock rate:                         900 Mhz
  Memory Bus Width:                          64-bit
  L2 Cache Size:                             524288 bytes
  Maximum Texture Dimension Size (x,y,z)      1D=(65536), 2D=(65536, 65536), 3D=(4096, 4096, 4096)
  Maximum Layered 1D Texture Size, (num) layers 1D=(16384), 2048 layers
  Maximum Layered 2D Texture Size, (num) layers 2D=(16384, 16384), 2048 layers
  Total amount of constant memory:             65536 bytes
  Total amount of shared memory per block:     49152 bytes
  Total number of registers available per block: 65536
  Warp size:                                   32
  Maximum number of threads per multiprocessor: 2048
  Maximum number of threads per block:         1024
  Max dimension size of a thread block (x,y,z): (1024, 1024, 64)
  Max dimension size of a grid size    (x,y,z): (2147483647, 65535, 65535)
  Maximum memory pitch:                      2147483647 bytes
  Texture alignment:                          512 bytes
  Concurrent copy and kernel execution:       Yes with 1 copy engine(s)
  Run time limit on kernels:                   Yes
  Integrated GPU sharing Host Memory:          No
  Support host page-locked memory mapping:     Yes
  Alignment requirement for Surfaces:          Yes
  Device has ECC support:                     Disabled
  Device supports Unified Addressing (UVA):     Yes
  Device PCI Domain ID / Bus ID / location ID: 0 / 1 / 0
  Compute Mode:
    < Default (multiple host threads can use ::cudaSetDevice() with device simultaneously) >

deviceQuery, CUDA Driver = CUDART, CUDA Driver Version = 8.0, CUDA Runtime Version = 8.0, NumDevs = 1, Device0 = GeForce GT 740M

```

3. Análisis de Factores limitantes del dispositivo

3.1. Registros

El Kernel de la multiplicación de matrices con tiles requiere almacenar 10 variables de tipo int(4 bytes) en registro, por lo que cada hebra requeriría 36 bytes en registros, cada registro es 4 bytes, necesitando así como mínimo 9 registros(16, por cercanía de potencias de 2, siendo 4096 hebras por bloque como máximo, pero el dispositivo solo permite 1024 hebras por bloque, así que el número de registros no limitará las operaciones.

3.2. Memoria Compartida

En memoria compartida se necesitará la cantidad de bytes que usarán Mds, Nds y Nds2 (el tipo int es de 4 bytes) $tile * tile * 4 * 3 = 12 * tile^2$ bytes.

tile=16 $16^2 * 12 = 3072bytes = 2Kb$

tile=32 $32^2 * 12 = 12288bytes = 12Kb$

tile=64 $64^2 * 12 = 49152bytes = 48Kb$

El dispositivo cuenta con 49152 bytes de memoria compartida por bloque, así la memoria compartida no limitará las operaciones.

3.3. Otros

Las limitaciones del hardware son principalmente:

- 2048 hebras por multiprocesador
- 1024 hebras por bloque
- 32 hebras en el warp

Lo recomendable es que el tamaño del bloque sea múltiplo del warp, para así tener todas las hebras activas. Por lo que usando un blocksize de 32 o 16 maximizaría el número de hebras en el SM.

| BlockSize | Nro Blocks | Nro Hebras en SM |
|-----------|------------|------------------|
| 32x32 | 2 | 2048 |
| 16x16 | 8 | 2048 |

Se realizarán experimentos con estas dos dimensiones para evaluar su ejecución.

4. Experimentos

4.1. Tabla de resultados (ms)

| Matriz | 256x256 | | 1024x1024 | | 4096x4096 | | 8192x8192 | |
|------------------|----------|----------|-----------|--------|-----------|---------|-----------|---------|
| Block/Tile Size | 16x16 | 32x32 | 16x16 | 32x32 | 16x16 | 32x32 | 16x16 | 32x32 |
| Mult. Simple | 2.28 | 0.59 | 150.6 | 145.04 | 8745.83 | 8370.59 | 77907.42 | 69143.5 |
| Mult. Tiled | 0.73 | 2.21 | 51.23 | 43.94 | 2764.95 | 2383.15 | 22154.74 | 19305.5 |
| Mult. Tiled Mod. | 1.076256 | 1.079264 | | | | | | |

4.2. Capturas de Pantalla

```

rose@Satellite-S55-A:~/CS_AlgoritmosParalelos/cuda$ ./mt16mod 256
N: 256 Bloques : 16 Hebras: 16
Tiempo : 1.076256 ms
rose@Satellite-S55-A:~/CS_AlgoritmosParalelos/cuda$ ./mt32mod 256
N: 256 Bloques : 8 Hebras: 32
Tiempo : 1.079264 ms

```

5. Conclusiones

La multiplicación con tiles y memoria compartida es más rápida que la multiplicación simple, disminuyendo la lectura de datos en un radio de 1 a 16 o 32 (según sea el tamaño del Tile).

Tener un Tile múltiplo del warp mejora el tiempo ya que siempre se tiene las hebras activas.

Para el dispositivo GeForce 740M, el mejor tamaño de tile es de 32x32.

Al modificar la multiplicación con tiles se obtiene mejor eficiencia y mejor tiempo.