



W • S E N S E  
INTEGRATED CABLELESS SOLUTIONS



SAPIENZA  
UNIVERSITÀ DI ROMA

# Programmazione di sistemi multicore **CUDA #2**

Fabrizio Gattuso

# Avvisi

12 Dicembre **NO LEZIONE PER IT MEETING**

**Il secondo esonero sarà svolto insieme al primo appello**

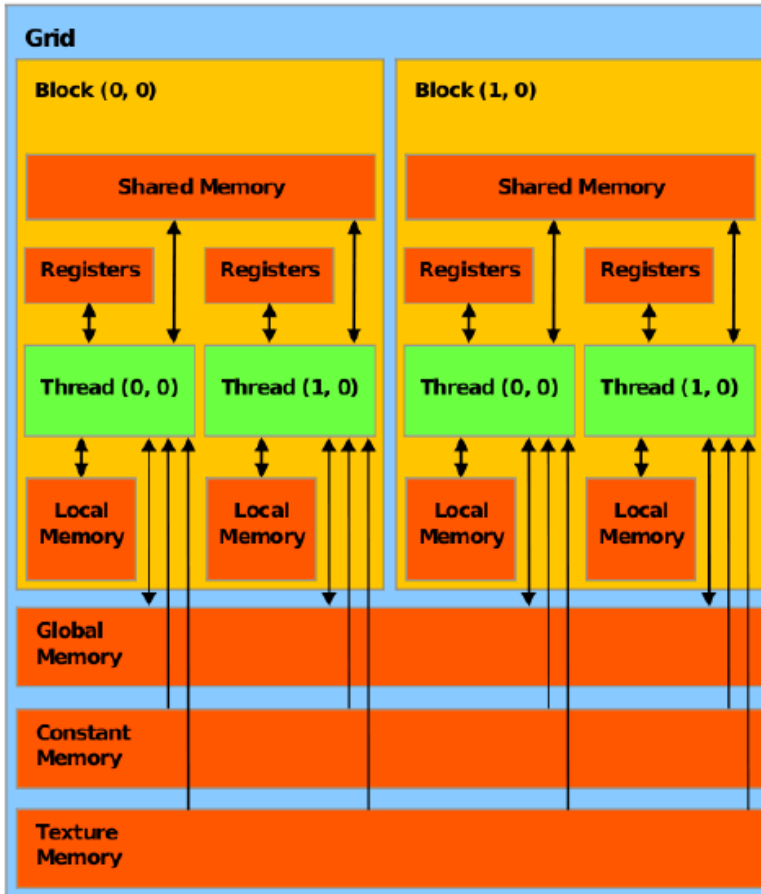
Date esame completo e primo esonero:

- **14 Gennaio** ore 16-18, Aula 2L - Via del Castro Laurenziano 7A
- **4 Febbraio** ore 16-18, Aula 2L - Via del Castro Laurenziano 7A

È uscito il codice OPIS **SIU975JZ** per la **valutazione del corso**

È importante che cominciate ad **inviare la mail con la richiesta materiale per il progetto**

# Gestione della memoria



Variable	Memoria	Scope
<code>int var</code>	register	thread
<code>int array[10]</code>	local	thread
<code>__shared__</code>	shared	block
<code>__device__</code>	global	grid
<code>__constant__</code>	constant	grid

# Lavorare con i thread

I thread all'interno dello stesso blocco sono *utili* perché possono comunicare tra loro.

Attraverso la direttiva

```
void __syncthreads ()
```

Tutti i thread saranno sincronizzati tra loro e solamente quando tutti avranno raggiunto la barriera il codice andrà avanti.

# Lavorare con i thread (2)

Questo è utile anche per prevenire tre tipi di hazards noti:

- RAW (Read after write)
- WAR (Write after read)
- WAW (Write after write)

**Fare sempre attenzione alle dipendenze temporali dei dati.**

# Lavorare con i thread (3)

Come già dovrete sapere la memoria utilizzabili per condividere dati tra due thread è quella:

1. **Global**
2. **Shared**

**\_\_shared\_\_ per varabili su area di memoria condivisa**

**Non usare la global se non strettamente necessario, in quanto, è molto più lenta rispetto alla shared. Al Massimo utilizzare la memoria constant.**

# DimGrid e DimBlock

Come abbiamo visto se si avvia un kernel nel seguente modo:

**myKernel**<<<N, M>>> verrà lanciato con N blocchi da M threads.

Se si vogliono specificare più nel dettaglio tutte le dimensioni, tra cui la griglia, è necessario usare l'indirizzamento esteso.

dim3 grid( 512, 1, 1) o dim3 grid( 512)

Dim3 block( 1024, 1024,

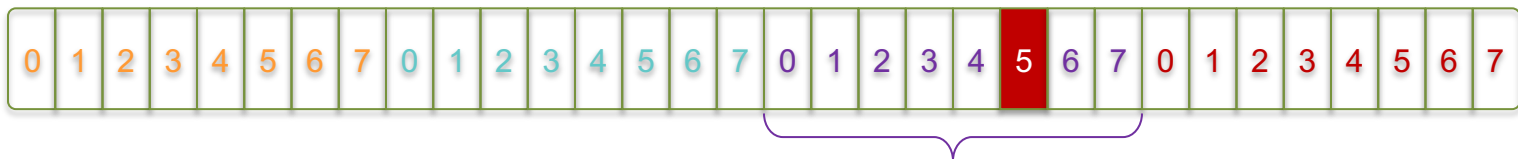
# Indexing

Abbiamo già visto come si può ricavare l'id univoco per un'architettura 1D:

$$\text{index} = \text{threadIdx.x} + \text{blockIdx.x} * \text{blockDim.x}$$

$M = 8$

$\text{threadIdx.x} = 5$





# Indexing (2)

**Index2D** = blockIdx.x \* blockDim.x \* blockDim.y +  
threadIdx.y \* blockDim.x + threadIdx.x

**Index3D** = blockIdx.x \* blockDim.x \* blockDim.y \*  
blockDim.z + threadIdx.z \* blockDim.y \* blockDim.x +  
threadIdx.y \* blockDim.x + threadIdx.x

**E questo solo per le griglie in 1D...**

# Gestione dei device (GPU)

Esistono delle comode funzioni per la gestione del device.

**cudaError\_t** cudaGetDeviceCount ( **int\*** count )

*Ritorna il numero di device disponibili.*

**cudaError\_t** cudaGetDevice ( **int\*** device )

*Ritorna il device in uso.*

**cudaError\_t** cudaSetDevice ( **int** device )

*Imposta il device da utilizzare.*

**cudaError\_t** cudaGetDeviceProperties ( **cudaDeviceProp\*** prop, **int** device )

*Ritorna la CC del device selezionato.*

# Gestione degli errori

Gestire gli errori è sempre una pratica essenziale per ogni sviluppatore. Nel calcolo parallelo è essenziale.

**cudaError\_t è l'enumeration utilizzata per rappresentare gli errori.**

cudaSuccess = 0

cudaErrorInvalidValue = 1

cudaErrorMemoryAllocation = 2

cudaErrorInvalidConfiguration = 9

cudaErrorDuplicateVariableName = 43



# Gestione degli errori (2)

**cudaError\_t** cudaGetLastError ( **void** )

*Ritorna l'ultimo errore lanciato.*

**char\*** cudaGetErrorString( **cudaError\_t** error )

*Ritorna una stringa sull'errore passato in input.*

# Calcolo dello speedup

L'incremento delle prestazioni di un programma parallelo rispetto a uno sequenziale è detto **speedup**.

Si può calcolare come: **tempo sequenziale / tempo parallelo**

$$\text{Speedup}(N) = \frac{1}{(1-P) + \frac{P}{N}}$$

Serial part of job = 1 (100%) - Parallel part

Parallel part is divided up by N workers

# Come scegliere i parametri da passare al kernel

1. Guardare sempre I limiti massimi dell'hw (**warp size, SM, max blocks, max threads per block**)
2. I thread all'interno del blocco dovrebbero essere multipli dello WARP size
3. Partire da numeri grandi di threads e blocks, calcolare lo speed up e provare a diminuire i parametri cercando la soluzione ideale



# Compute Capability (CC)

	FERMI GF100	FERMI GF104	KEPLER GK104	KEPLER GK110	KEPLER GK210
<b>Compute Capability</b>	2.0	2.1	3.0	3.5	3.7
<b>Threads / Warp</b>	32				
<b>Max Threads / Thread Block</b>	1024				
<b>Max Warps / Multiprocessor</b>	48		64		
<b>Max Threads / Multiprocessor</b>	1536		2048		
<b>Max Thread Blocks / Multiprocessor</b>	8		16		
<b>32-bit Registers / Multiprocessor</b>	32768		65536		131072
<b>Max Registers / Thread Block</b>	32768		65536		65536
<b>Max Registers / Thread</b>	63			255	
<b>Max Shared Memory / Multiprocessor</b>	48K				112K
<b>Max Shared Memory / Thread Block</b>	48K				
<b>Max X Grid Dimension</b>	2 <sup>16</sup> -1		2 <sup>32</sup> -1		
<b>Hyper-Q</b>	No			Yes	
<b>Dynamic Parallelism</b>	No			Yes	



# Piccolo esercizio

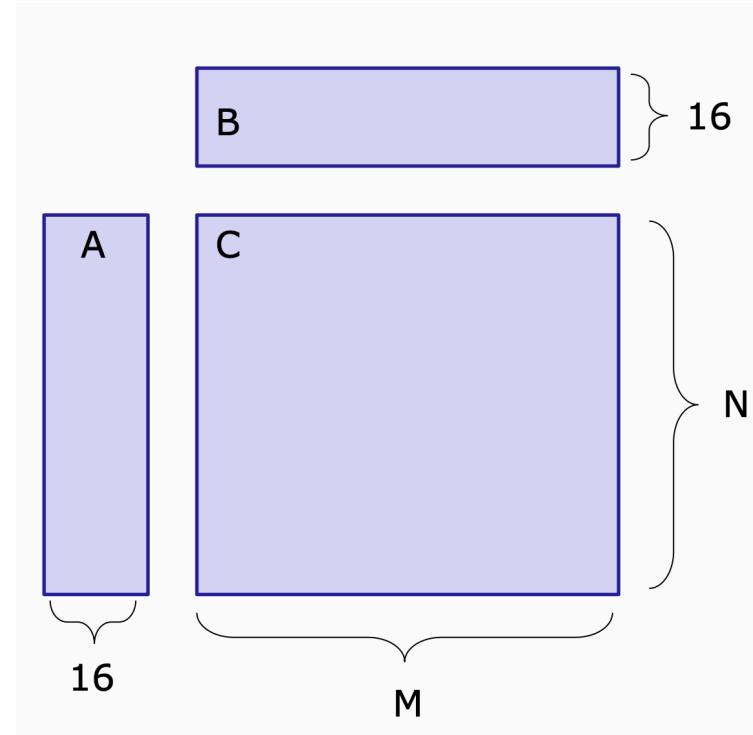
Calcolo parallelo del prodotto di due matrici

$$C = AB$$

**A** una matrice  $N \times 16$

**B** una matrice  $16 \times M$

**M** e **N** multipli di 16



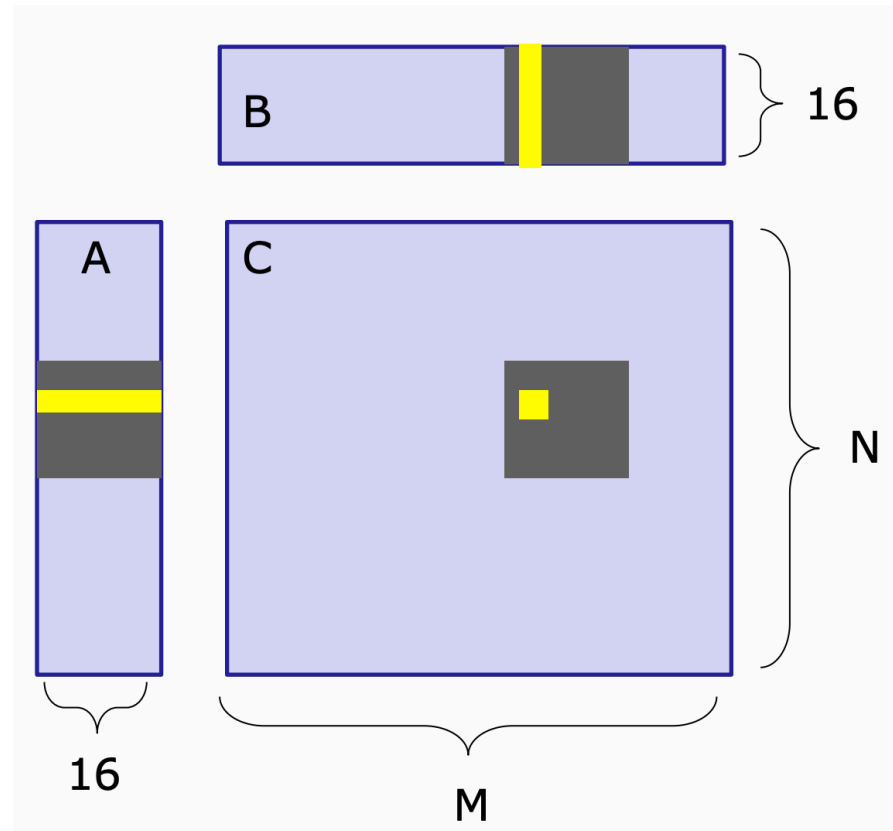


# Piccolo esercizio (2)

La matrice viene divisa in tile da  $16 \times 16$

Ogni tile viene passata a un blocco di thread

Ogni thread calcola un elemento della matrice risultante



# Piccolo esercizio (3)

```
__global__ void simpleMultiply(float *a, float* b, float  
*c, int N) {  
    int row = blockIdx.y * blockDim.y + threadIdx.y;  
    int col = blockIdx.x * blockDim.x + threadIdx.x;  
    float sum = 0.0f;  
    for (int i = 0; i < TILE_DIM; i++) {  
        sum += a[row*TILE_DIM+i] * b[i*N+col];  
    }  
    c[row*N+col] = sum;  
}
```

