

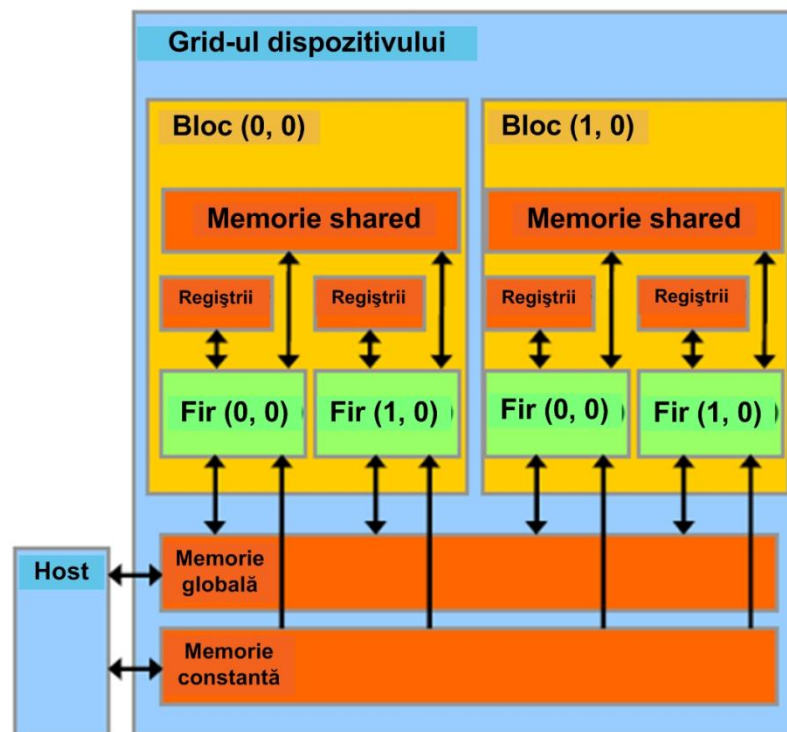
Introducere în CUDA

Continuare

Cosmin – Ioan Niță

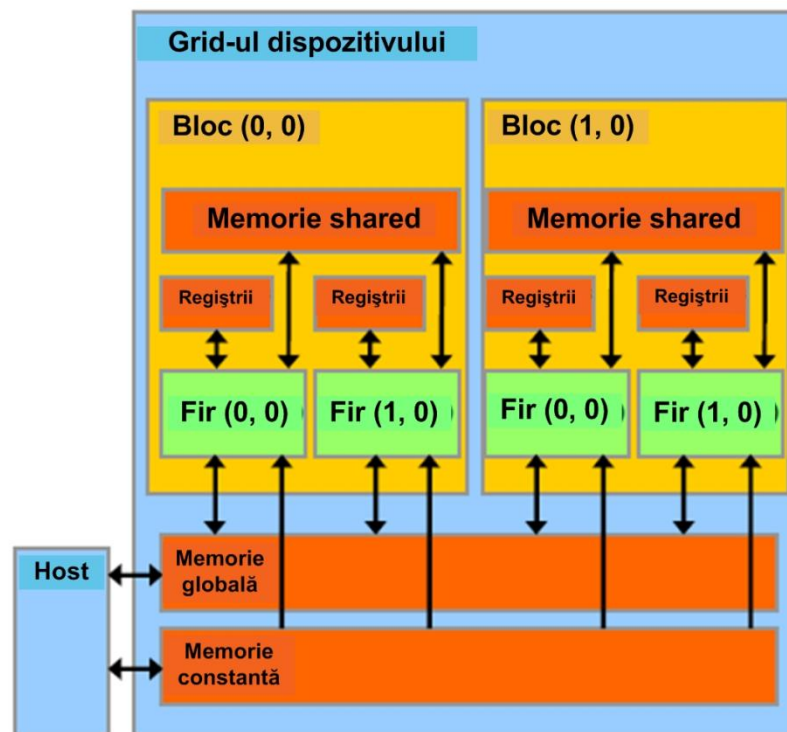
Tipuri de memorii CUDA

- Memorie globală
- Memorie partajată (shared memory)
- Regiștrii
- Caracteristici:
 - Vizibilitate
 - Perioadă de viață
 - Latență



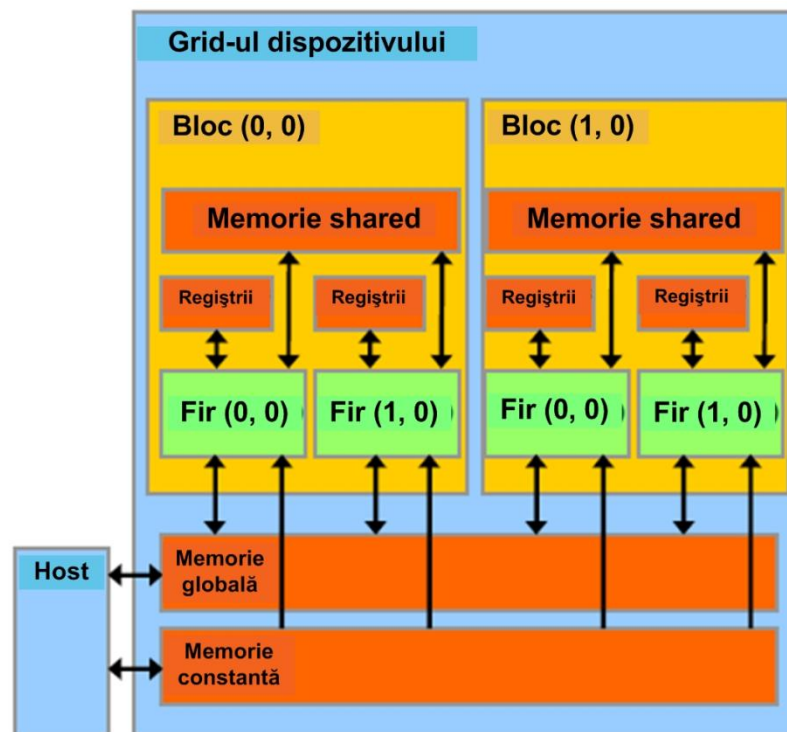
Tipuri de memorii CUDA

- Memorie globală
 - Se alocă și eliberează doar de pe CPU (*cudaMalloc*, *cudaFree*)
 - Poate fi accesată atât de CPU cât și de GPU (*cudaMemset*, *cudaMemcpy*).
 - La nivel de GPU este vizibilă oricărui thread.
 - Accesul este foarte lent.



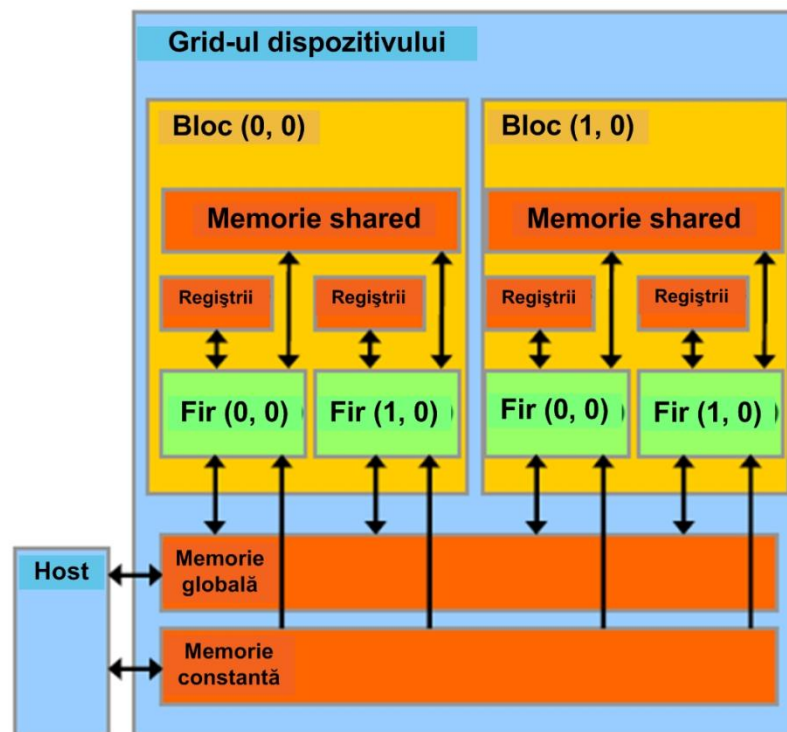
Tipuri de memorii CUDA

- Memorie partajată (shared memory)
 - Vizibilă doar de pe GPU la nivelul unui bloc.
 - Se alocă în interiorul unui kernel. Ex:
`__shared__ float f[128];`
 - Accesul este foarte rapid.
 - Limitată, mult mai mică decât memoria globală



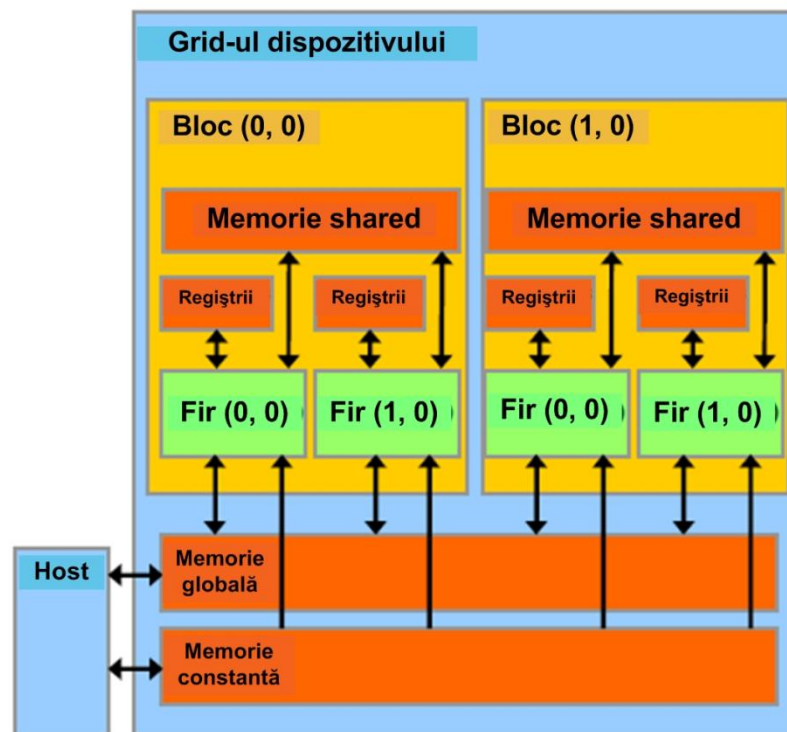
Tipuri de memorii CUDA

- Regiștrii
 - Vizibili doar la nivel de thread.
 - Variabilele declarate local într-un thread vor fi automat plasate în regiștrii.
 - Nu pot fi declarați ca un tablou.
 - Accesul este foarte rapid.



Tipuri de memorii CUDA

- Durata de viață
 - Regiștrii: kernel
 - Mem. Partajată: kernel
 - Mem. Globală: Aplicație



Tipuri de memorii CUDA

- Utilizare

- Regiștrii:

```
__global__ void kernel ()  
{  
    //Registrii  
    float f1, f2, f3;  
}
```

Tipuri de memorii CUDA

- Utilizare
 - Memoria partajata:

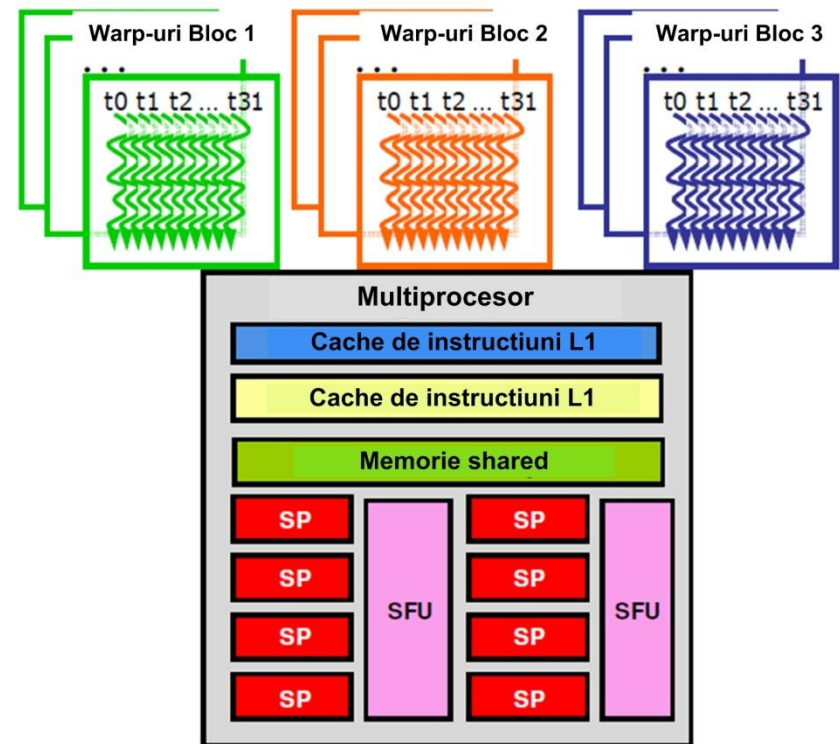
```
__global__ void kernel ()  
{  
    //Shared memory  
    __shared__ float tmp[128];  
}
```


Factori ce limitează paralelismul

- Dacă se depășește numărul maxim de regiștrii aceștia vor fi alocați în memoria globală.
- Dacă se depășește cantitatea de memorie shared se reduce paralelismul
- DeviceQuery
- Opțiunea `--ptxas-options=-v`

Optimizarea accesului la memoria globală

- Accesul la memoria globală – factor ce limitează viteza de execuție.
- Mascarea latenței
- Când un warp așteaptă date de la mem. Globală, GPU-ul porneste execuția altui warp.
- Citirea din mem. Globală se poate face în paralel pentru jumătate de warp dacă se accesează adrese de memorie consecutive.



SFU – special function unit = unitate de calcul a operațiilor speciale

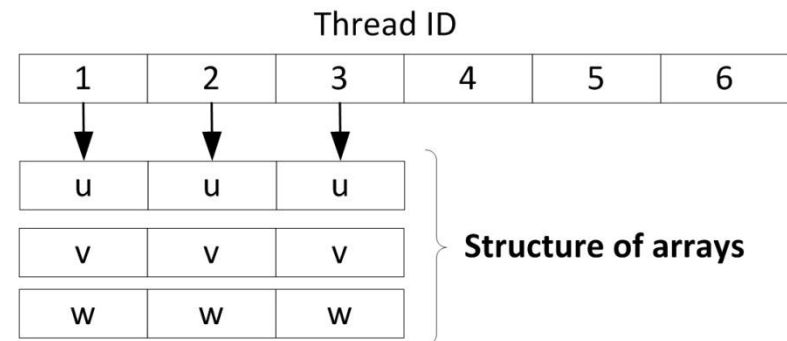
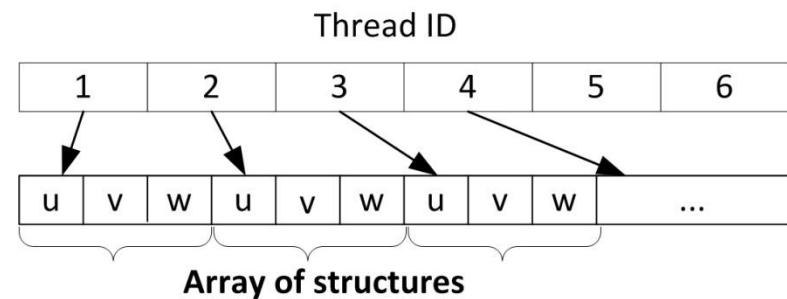
Optimizarea accesului la memoria globală

- Exemplu:
 - Transformarea unei mulțimi de puncte 3D.
 - Fiecare thread va executa o transformare.

$$p_i \leftarrow T \cdot p_i$$

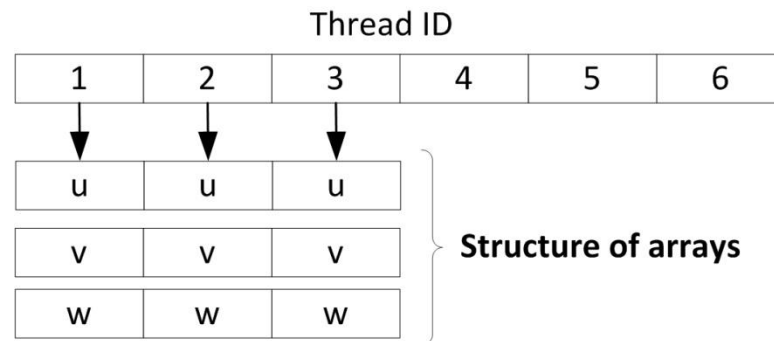
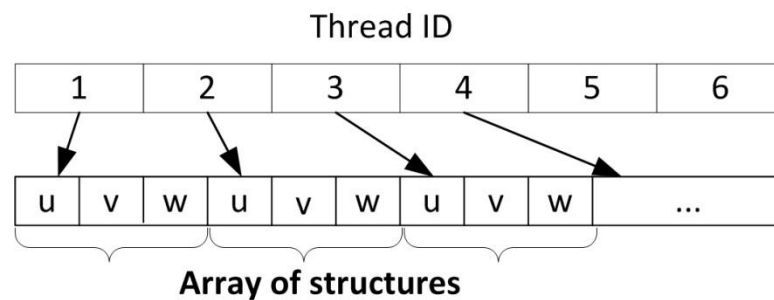
$$p_i = [u_i, v, w_i]^T$$

$$T = \begin{bmatrix} t_{xx} & t_{xy} & t_{xz} \\ t_{yx} & t_{yy} & t_{yz} \\ t_{zx} & t_{zy} & t_{zz} \end{bmatrix}$$



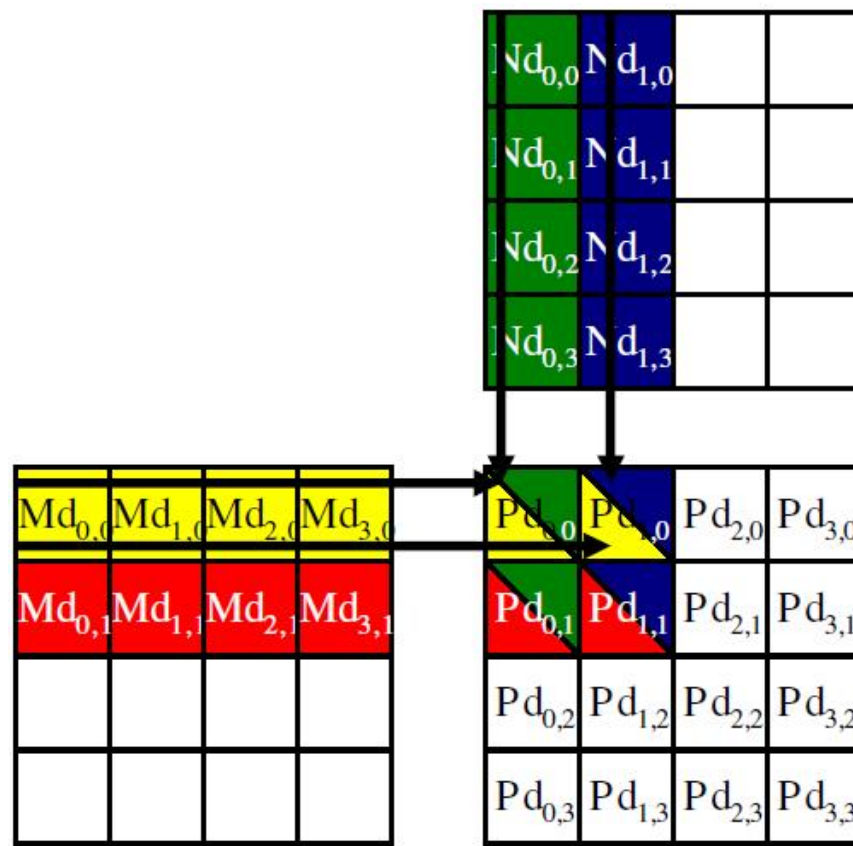
Optimizarea accesului la memoria globală

- Exemplu:
 - Înmulțirea a două matrici.



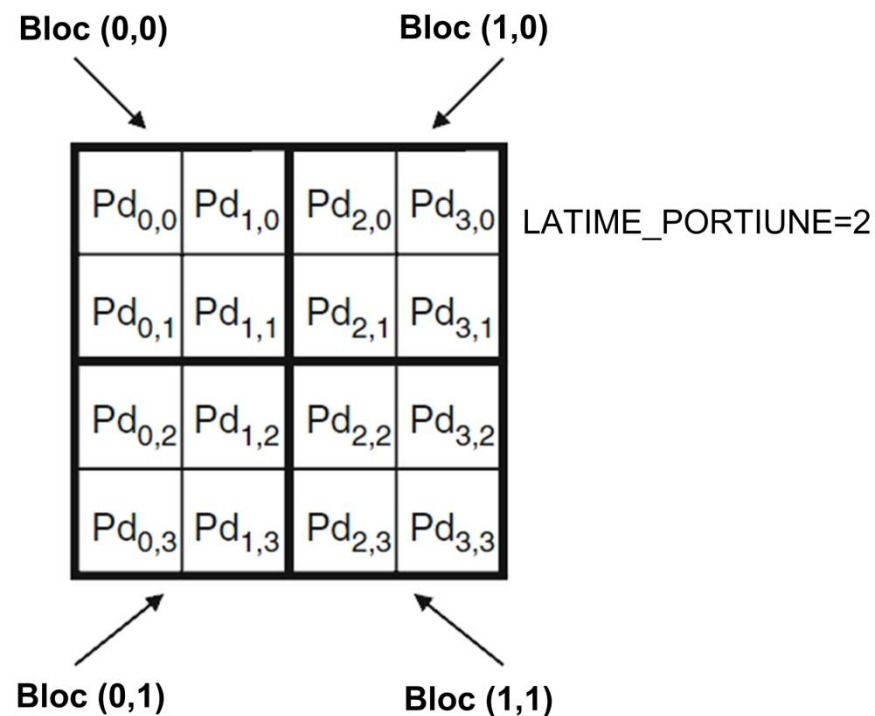
Optimizarea unei aplicații CUDA folosind memoria partajata

- Exemplu pentru înmulțirea a două matrici
 - Accesări redundante ale memoriei globale.

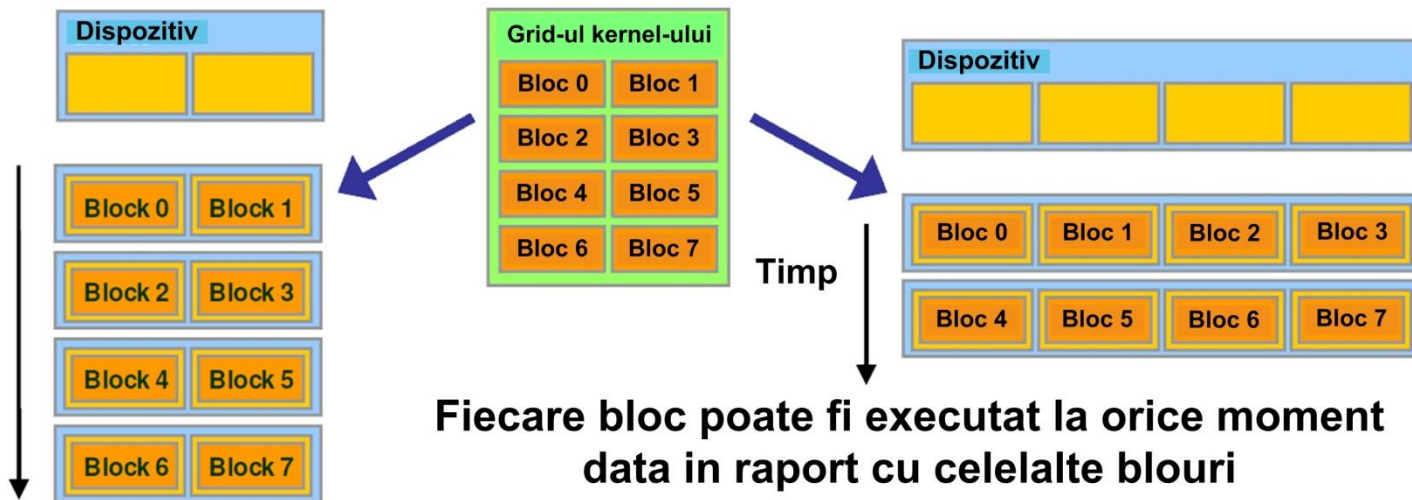


Optimizarea unei aplicații CUDA folosind memoria partajata

- Exemplu pentru înmulțirea a două matrici
 - Accesări redundante ale memoriei globale.
 - Produse parțiale
 - Fiecare thread va copia valori din mem. Globală în memoria shared.



Optimizarea unei aplicații CUDA folosind memoria partajata



Optimizarea unei aplicații CUDA folosind memoria partajata

- Bariera de sincronizare a thread-urilor dintr-un bloc.
 - `__syncthreads();`

Optimizarea unei aplicații CUDA folosind memoria partajata

- Bariera de sincronizare a thread-urilor dintr-un bloc.
 - `__syncthreads();`
 - Deadlock

Optimizarea unei aplicații CUDA folosind memoria partajata

