



Formação Inteligência Artificial

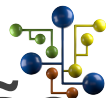


Programação Paralela em GPU



Programação Paralela em CUDA

Parte 2



Programação Paralela em CUDA - Parte 2



Qualificadores



Qualificadores

__global__

Define um Kernel
Executa na **GPU**, sendo
chamado pela **CPU**

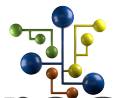
__device__

Pode ser usado para variáveis
Executa na **GPU**, sendo
chamado pela **GPU**

__host__

Para tarefas na CPU
Executa na **CPU**, sendo
chamado pela **CPU**





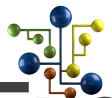
Qualificadores

Os qualificadores também podem ser combinados em uma única instrução:

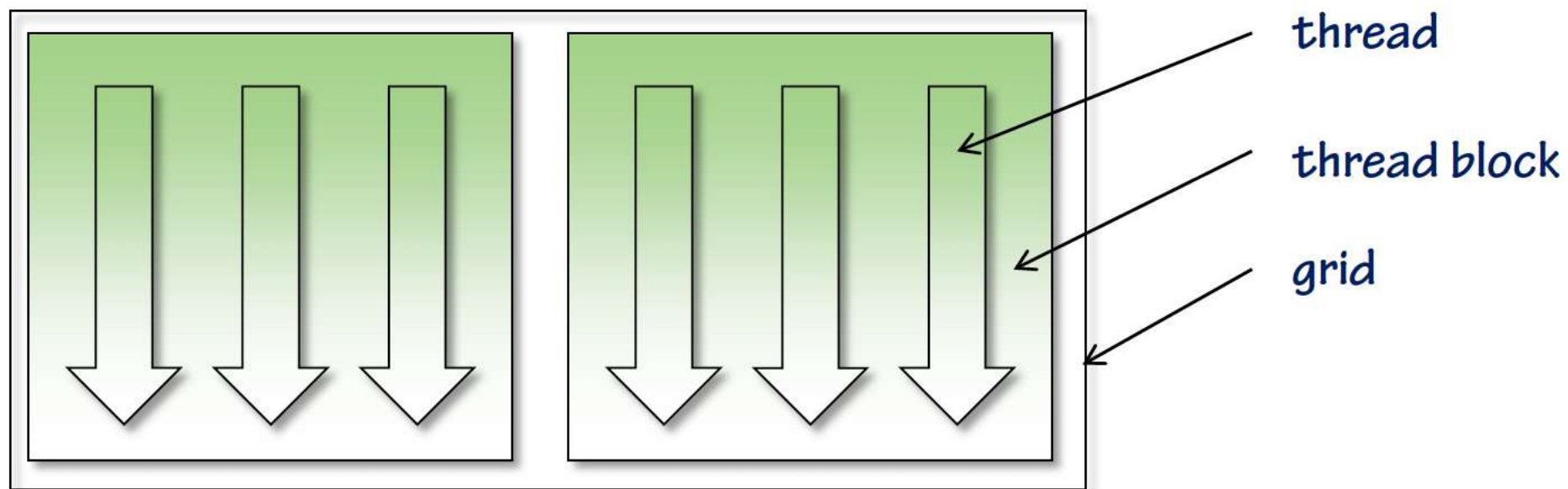
```
__host__ __device__ nome_func()
```



Modelo de Execução

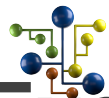


Modelo de Execução

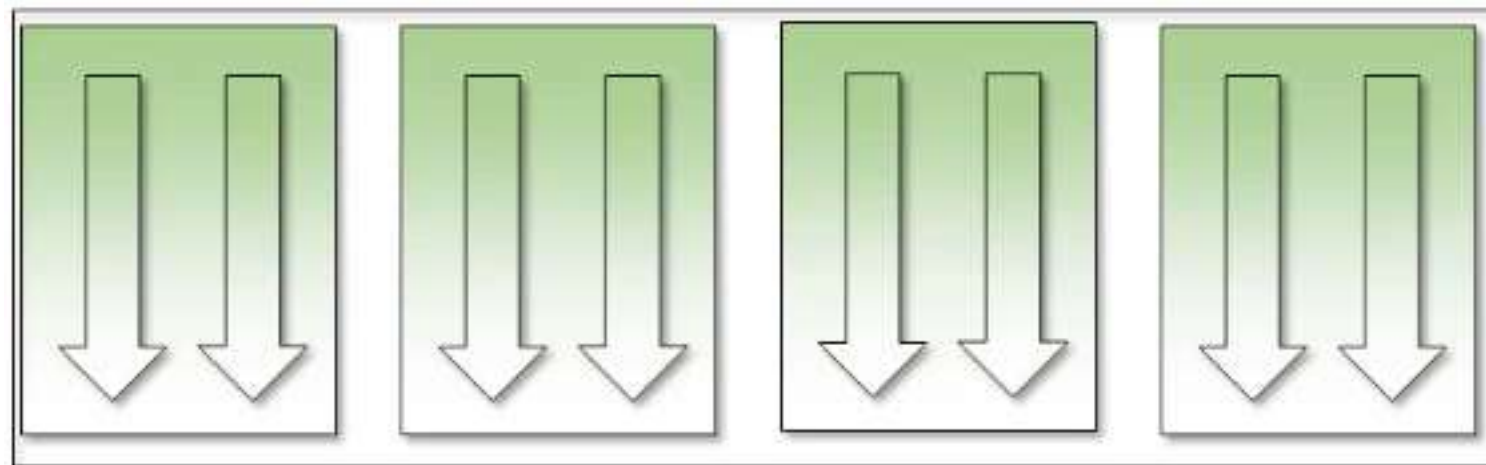


`addArraysGPU <<< 1, count >>> (da, db, dc)`

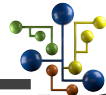




Modelo de Execução

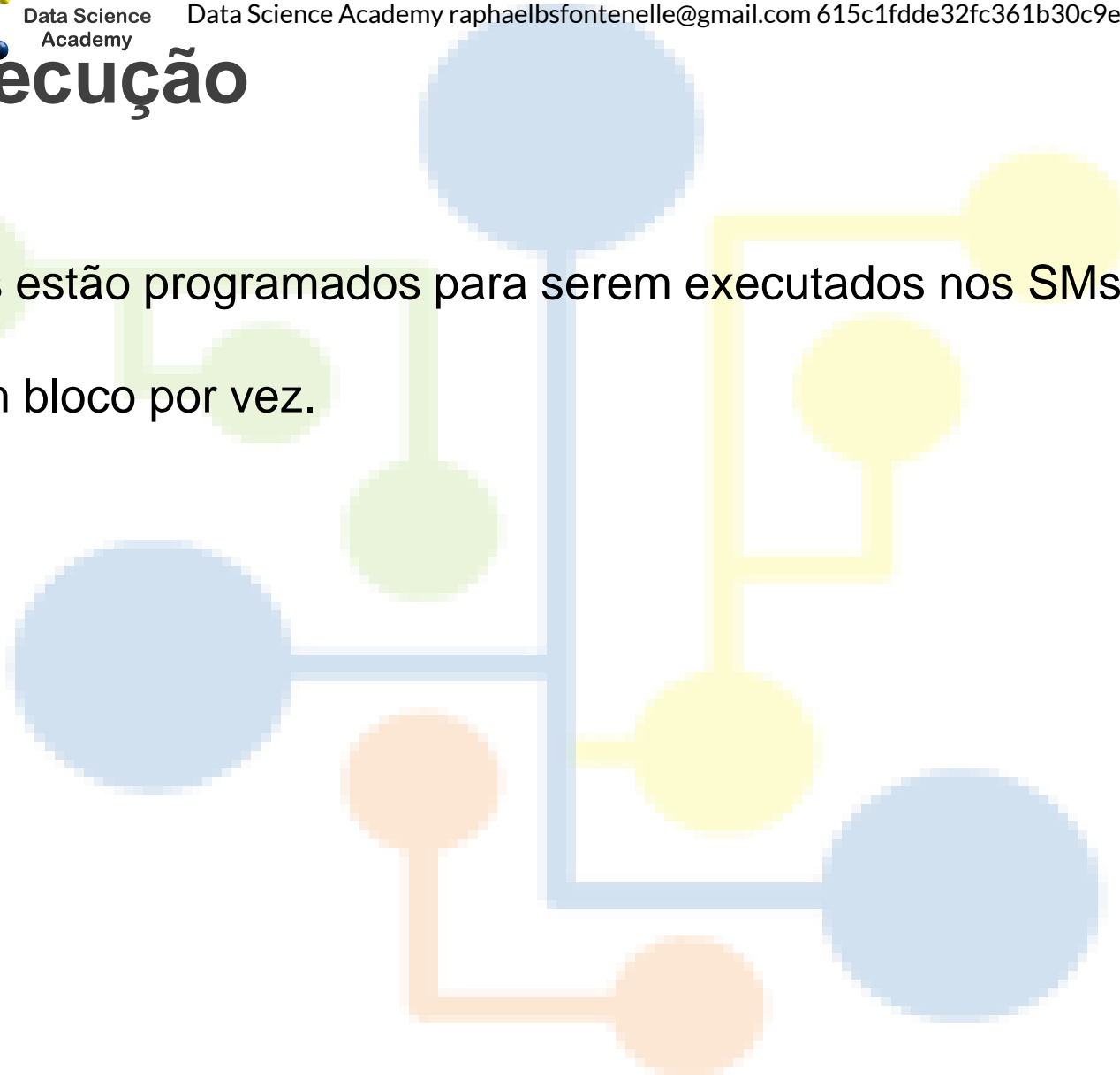


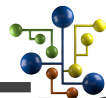
`addArraysGPU <<< 4, 2 >>> (da, db, dc)`



Modelo de Execução

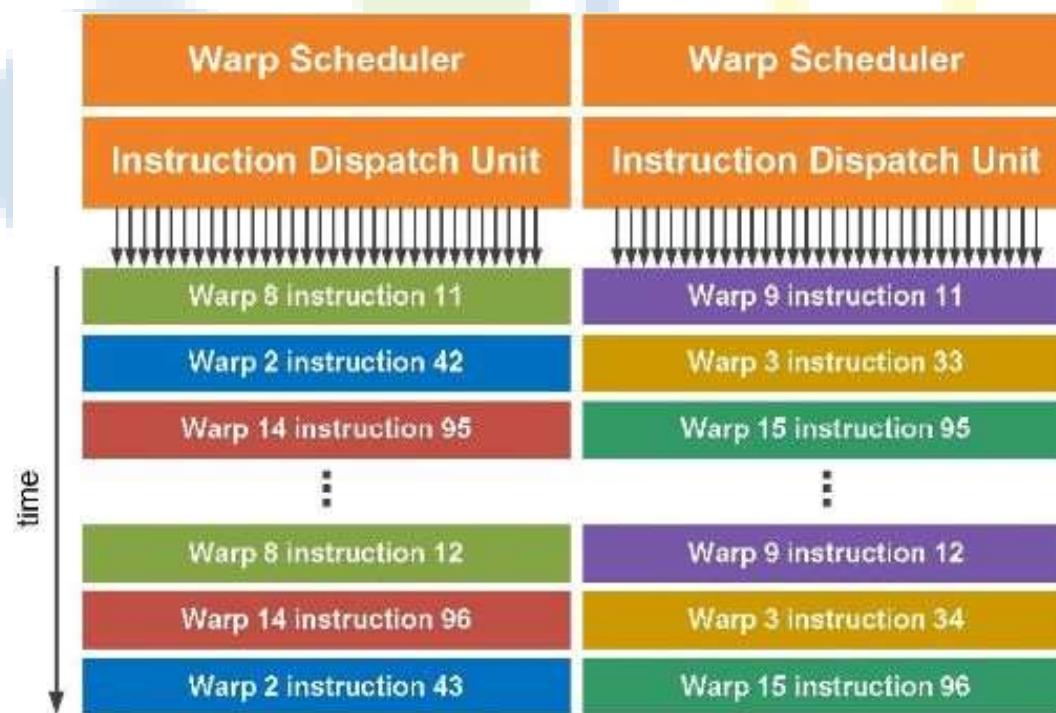
- Os blocos de threads estão programados para serem executados nos SMs disponíveis.
- Cada SM executa um bloco por vez.

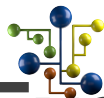




Modelo de Execução

- Cada bloco de thread é dividido em warps (a variável warp size da GPU define quantos warps podem ser usados).
- Os warps são executados de forma paralela (o Warp Watch nos permite visualizar a execução dos warps).





Modelo de Execução

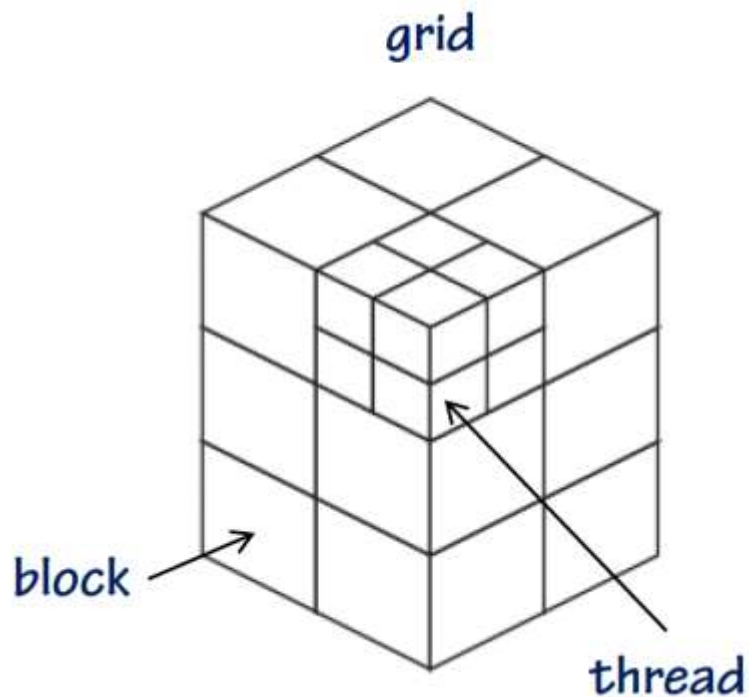
- Se um warp leva muito tempo para sua execução, o scheduler da GPU pode iniciar um novo warp.
- O número de threads por warps depende da capacidade de cálculo (Compute Capability).
- Todas as warps são tratadas em paralelo.



Dimensões de Grids e Blocos



Dimensões de Grids e Blocos



Definimos execução como:

$\langle\langle\langle a, b \rangle\rangle\rangle$

Um grid de a blocos com b threads

$\langle\langle\langle a, b \rangle\rangle\rangle$

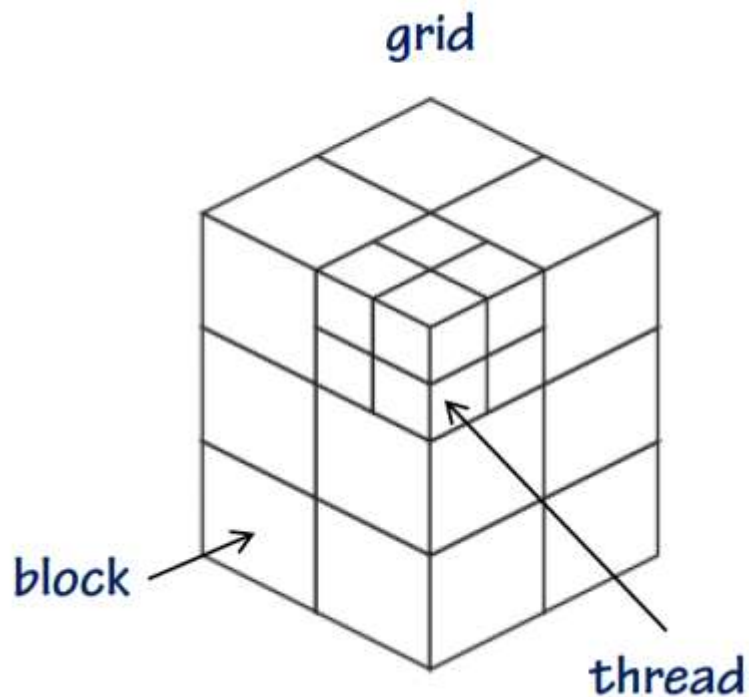
(a, 1, 1) blocos por (b, 1, 1) threads

Estrutura dim3





Dimensões de Grids e Blocos



Definimos execução como:

$\langle\langle\langle a, b \rangle\rangle\rangle$

Um grid de a blocos com b threads

Um grid é uma estrutura 3D

Definimos como $(A \times B \times C)$ blocos de $(x \times y \times z)$ threads





Dimensões de Grids e Blocos

Parâmetros de execução e posição atual:

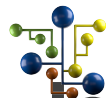
- BlockIdx
 - Onde estamos no grid
- GridDim
 - O tamanho do grid
- ThreadIdx
 - Posição da thread atual no bloco de thread
- blockDim
 - Tamanho do bloco de thread
- Limitações
 - Tamanhos de grid e bloco
 - Número de threads

MAX_BLOCK_DIM_X	512
MAX_BLOCK_DIM_Y	512
MAX_BLOCK_DIM_Z	64
MAX_GRID_DIM_X	65535
MAX_GRID_DIM_Y	65535
MAX_GRID_DIM_Z	1

MAX_THREADS_PER_BLOCK	512
MAX_THREADS_PER_MULTIPROCESSOR	1024



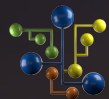
Tratamento de Erros



Tratamento de Erros

- CUDA não gera mensagem de erro
 - Falha silenciosa
- As funções retornam `cudaError_t`
 - Verifique status `cudaSuccess`
 - Para o código de erro use `cudaGetErrorString()`





Data Science
Academy

Data Science Academy raphaelbsfontenelle@gmail.com 615c1fdde32fc361b30c9ec2

Obrigado



Data Science Academy



Data Science Academy