

ECE408 /CS483/CSE408 Spring 2020

Applied Parallel Programming

## Lecture 2: Introduction to CUDA C and Data Parallel Programming

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

1

1

## Objective

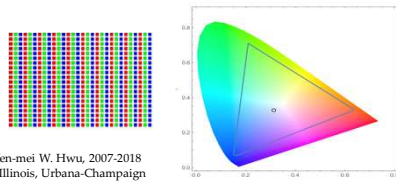
- To learn the basic concept of data parallel computing
- To learn the basic features of the CUDA C programming interface

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

2

2

*A Data Parallel Computation Example:  
Conversion of a color image to grey-  
scale image*

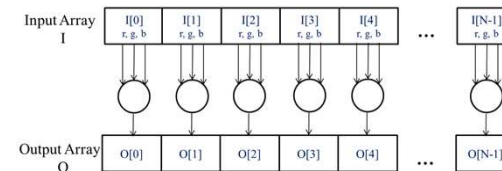


© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

3

3

*The pixels can be calculated  
independently of each other*



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

4

4

## CUDA/OpenCL – Execution Model

- Integrated host+device app C program
  - Serial or modestly parallel parts in **host** C code
  - Highly parallel parts in **device** SPMD kernel C code

Serial Code (host)

Parallel Kernel (device)  
KernelA<<< nBlk, nTid >>>(args);

Serial Code (host)

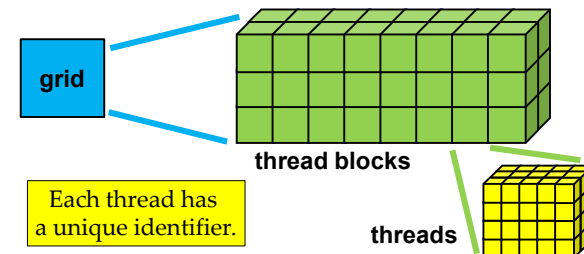
Parallel Kernel (device)  
KernelB<<< nBlk, nTid >>>(args);

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

5

## Logical Execution Model for CUDA

- Each CUDA kernel
  - is executed by a **grid**,
  - a 3D array of **thread blocks**, which are
  - 3D arrays of **threads**.

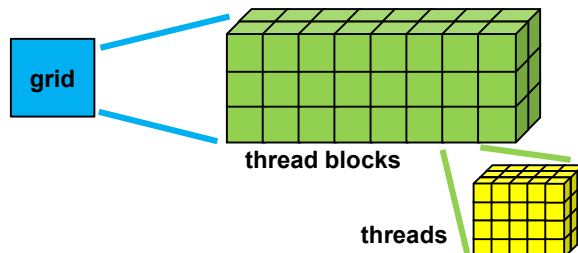


© 2020 by S. S. Lumetta

6

## Single Program, Multiple Data

- Each thread
  - executes the **same program**
  - on **distinct data inputs**,
  - a single-program, multiple-data (**SPMD**) model



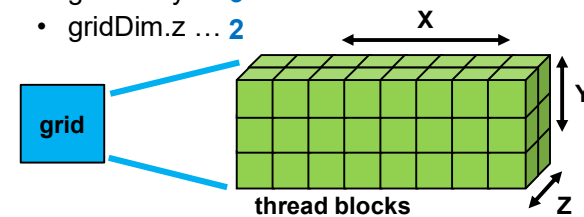
© 2020 by S. S. Lumetta

7

7

## gridDim Gives Number of Blocks

- Number of blocks in each dimension is
  - gridDim.x ... **8**
  - gridDim.y ... **3**
  - gridDim.z ... **2**



For 2D (and 1D grids), simply use grid dimension 1 for Z (and Y).

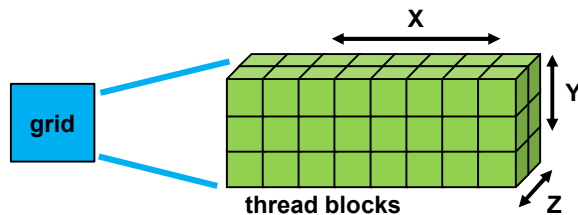
© 2020 by S. S. Lumetta

8

8

## blockIdx is Unique for Each Block

- Each block has a unique index tuple
  - blockIdx.x (from 0 to (gridDim.x - 1))
  - blockIdx.y (from 0 to (gridDim.y - 1))
  - blockIdx.z (from 0 to (gridDim.z - 1))



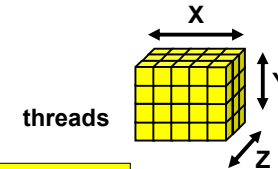
© 2020 by S. S. Lumetta

9

9

## blockDim: # of Threads per Block

- Number of blocks in each dimension is
  - blockDim.x ... 5
  - blockDim.y ... 4
  - blockDim.z ... 3



For 2D (and 1D blocks), simply use block dimension 1 for Z (and Y).

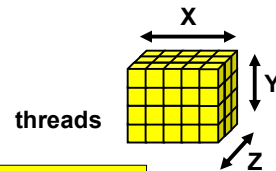
© 2020 by S. S. Lumetta

10

10

## threadIdx Unique for Each Thread

- Each thread has a unique index tuple
  - threadIdx.x (from 0 to (blockDim.x - 1))
  - threadIdx.y (from 0 to (blockDim.y - 1))
  - threadIdx.z (from 0 to (blockDim.z - 1))



threadIdx tuple is unique to each thread  
WITHIN A BLOCK.

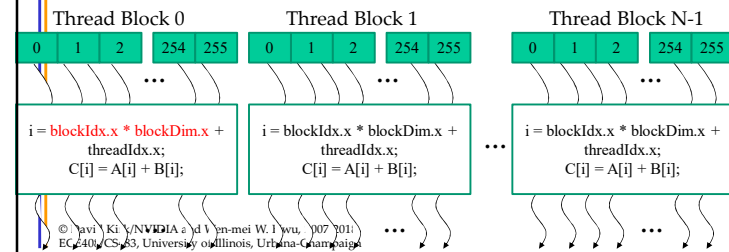
© 2020 by S. S. Lumetta

11

11

## Thread Blocks: Scalable Cooperation

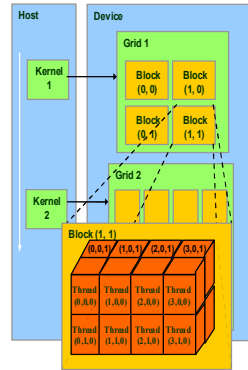
- Threads within a block cooperate via **shared memory, atomic operations** and **barrier synchronization** (to be covered later)
- Threads in different blocks cooperate less.



12

## blockIdx and threadIdx

- Thread block and thread organization
- simplifies memory addressing
- when processing multidimensional data
  - Image processing
  - Vectors, matrices, tensors
  - Solving PDEs on volumes
  - ...

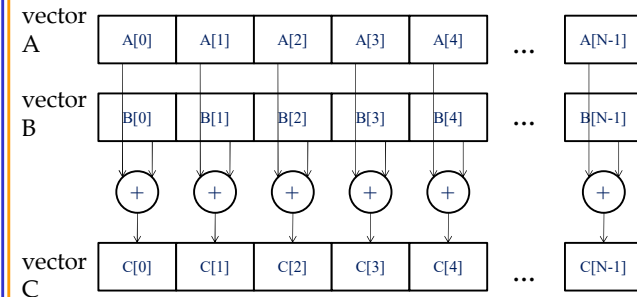


© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

13

13

## Vector Addition – Conceptual View



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

14

14

## Vector Addition – Traditional C Code

```
// Compute vector sum C = A+B
void vecAdd(float* A, float* B, float* C, int n)
{
    for (i = 0, i < n, i++)
        C[i] = A[i] + B[i];
}

int main()
{
    // Memory allocation for A_h, B_h, and C_h
    // I/O to read A_h and B_h, N elements
    ...
    vecAdd(A_h, B_h, C_h, N);
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

15

15

## Heterogeneous Computing vecAdd Host Code

```
#include <cuda.h>
void vecAdd(float* A, float* B, float* C, int n)
{
    int size = n * sizeof(float);
    float *A_d, *B_d, *C_d;

    ...

    1. // Allocate device memory for A, B, and C
       // copy A and B to device memory

    2. // Kernel launch code - to have the device
       // to perform the actual vector addition

    3. // copy C from the device memory
       // Free device vectors
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

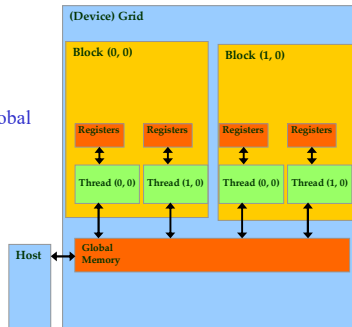
16

16

## Partial Overview of CUDA Memories

- Device code can:
  - R/W per-thread **registers**
  - R/W per-grid **global memory**
- Host code can
  - Transfer data to/from per grid **global memory**

We will cover more later.



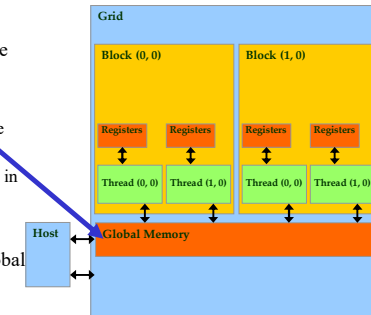
© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

17

17

## CUDA Device Memory Management API functions

- **cudaMalloc()**
  - Allocates object in the device **global memory**
  - Two parameters
    - **Address of a pointer** to the allocated object
    - **Size of the allocated object** in terms of bytes
- **cudaFree()**
  - Frees object from device global memory
  - **Pointer** to freed object



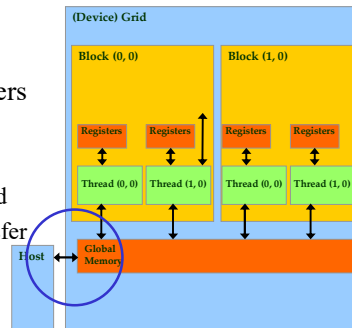
© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

18

18

## Host-Device Data Transfer API functions

- **cudaMemcpy()**
  - memory data transfer
  - Requires four parameters
    - Pointer to destination
    - Pointer to source
    - Number of bytes copied
    - Type/Direction of transfer



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

19

19

```
void vecAdd(float* A, float* B, float* C, int n)
{
    int size = n * sizeof(float);
    float *A_d, *B_d, *C_d;

    1. // Transfer A and B to device memory
    // (error-checking omitted)
    cudaMalloc((void **) &A_d, size);
    cudaMemcpy(A_d, A, size, cudaMemcpyHostToDevice);
    cudaMalloc((void **) &B_d, size);
    cudaMemcpy(B_d, B, size, cudaMemcpyHostToDevice);

    // Allocate device memory for
    cudaMalloc((void **) &C_d, size);

    2. // Kernel invocation code - to be shown later
    ...

    3. // Transfer C from device to host
    cudaMemcpy(C, C_d, size, cudaMemcpyDeviceToHost);
    // Free device memory for A, B, C
    cudaFree(A_d); cudaFree(B_d); cudaFree(C_d);
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

20

20

## Example: Vector Addition Kernel

Device Code

```
// Compute vector sum C = A+B
// Each thread performs one pair-wise addition
__global__
void vecAddKernel(float* A_d, float* B_d, float* C_d, int n)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    if(i < n) C_d[i] = A_d[i] + B_d[i];
}

int vecAdd(float* A, float* B, float* C, int n)
{
    // A_d, B_d, C_d allocations and copies omitted
    // Run ceil(n/256) blocks of 256 threads each
    vecAddKernel<<<ceil(n/256.0), 256>>>>(A_d, B_d, C_d, n);
}
```

21

## Example: Vector Addition Kernel

Host Code

```
// Compute vector sum C = A+B
// Each thread performs one pair-wise addition
__global__
void vecAddKernel(float* A_d, float* B_d, float* C_d, int n)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    if(i < n) C_d[i] = A_d[i] + B_d[i];
}

int vecAdd(float* A, float* B, float* C, int n)
{
    // A_d, B_d, C_d allocations and copies omitted
    // Run ceil(n/256) blocks of 256 threads each
    vecAddKernel<<<ceil(n/256.0), 256>>>>(A_d, B_d, C_d, n);
}
```

22

22

## More on Kernel Launch

Equivalent Host Code

```
int vecAdd(float* A, float* B, float* C, int n)
{
    // A_d, B_d, C_d allocations and copies omitted
    // Run ceil(n/256) blocks of 256 threads each
    dim3 DimGrid(n/256, 1, 1);
    if (0 != (n % 256)) { DimGrid.x++; }
    dim3 DimBlock(256, 1, 1);

    vecAddKernel<<<DimGrid, DimBlock>>>>(A_d, B_d, C_d, n);
}
```

- Any call to a kernel function is asynchronous from CUDA 1.0 on, explicit synch needed for blocking

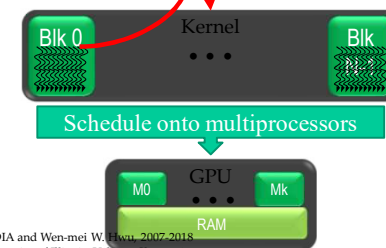
23

23

## Kernel execution in a nutshell

```
__host__
void vecAdd()
{
    dim3 DimGrid(ceil(n/256.0), 1, 1);
    dim3 DimBlock(256, 1, 1);
    vecAddKernel<<<DimGrid, DimBlock>>>>(A_d, B_d, C_d, n);
}

__global__
void vecAddKernel(float* A_d, float* B_d, float* C_d, int n)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    if(i < n) C_d[i] = A_d[i] + B_d[i];
}
```



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

24

24

## More on CUDA Function Declarations

	Executed on the:	Only callable from the:
<code>__device__ float DeviceFunc()</code>	device	device
<code>__global__ void KernelFunc()</code>	device	host
<code>__host__ float HostFunc()</code>	host	host

- `__global__` defines a kernel function
- Each “`__`” consists of two underscore characters
- A kernel function must return `void`
- `__device__` and `__host__` can be used together

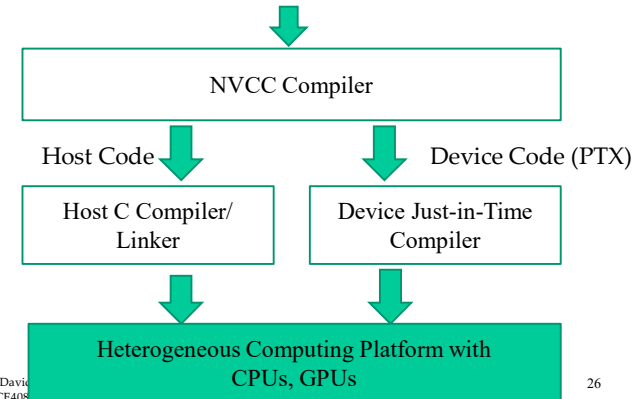
© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

25

25

## Compiling A CUDA Program

Integrated C programs with CUDA extensions



© David  
ECE408

26

26

# QUESTIONS?

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2018  
ECE408/CS483, University of Illinois, Urbana-Champaign

27

27