



GPU Teaching Kit  
Accelerated Computing



# CUDA Parallelism Model

Kernel-Based SPMD Parallel Programming

## Objective

- To learn the basic concepts involved in a simple CUDA kernel function
  - Declaration
  - Built-in variables
  - Thread index to data index mapping

# Example: Vector Addition Kernel

## Device Code

```
// Compute vector sum C = A + B
// Each thread performs one pair-wise addition
```

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

3

NVIDIA

ILLINOIS

# Example: Vector Addition Kernel Launch (Host Code)

## Host Code

```
void vecAdd(float* h_A, float* h_B, float* h_C, int n)
{
    // d_A, d_B, d_C allocations and copies omitted
    // Run ceil(n/256.0) blocks of 256 threads each
    vecAddKernel<<<ceil(n/256.0),256>>>(d_A, d_B, d_C, n);
}
```

The ceiling function makes sure that there are enough threads to cover all elements.

4

NVIDIA

ILLINOIS

# More on Kernel Launch (Host Code)

## Host Code

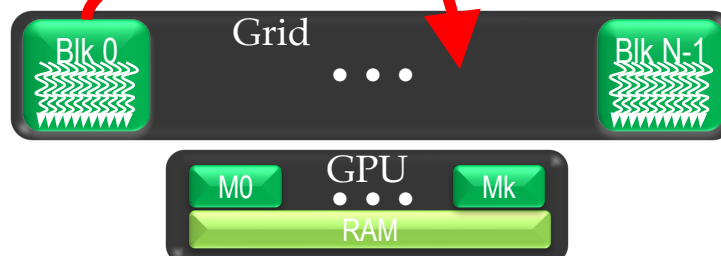
```
void vecAdd(float* h_A, float* h_B, float* h_C, int n)
{
    dim3 DimGrid((n-1)/256 + 1, 1, 1);
    dim3 DimBlock(256, 1, 1);
    vecAddKernel<<<DimGrid,DimBlock>>>(d_A, d_B, d_C, n);
}
```

This is an equivalent way to express the ceiling function.

## Kernel execution in a nutshell

```
__host__
void vecAdd(...)
{
    dim3 DimGrid(ceil(n/256.0),1,1);
    dim3 DimBlock(256,1,1);
    vecAddKernel<<<DimGrid,DimBlock>>>(d_A,d_B
    ,d_C,n);
}

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



## More on CUDA Function Declarations

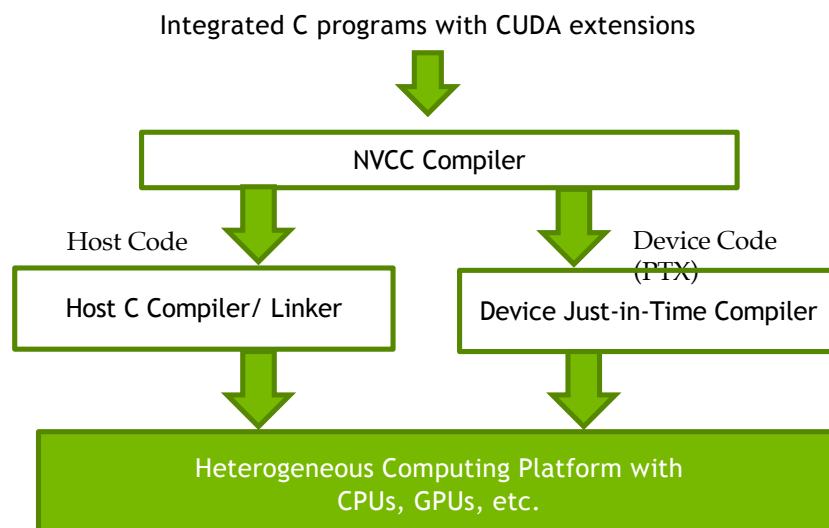
	Executed on the:	Only callable from the:
<code>__device__ float DeviceFunc ()</code>	device	device
<code>__global__ void KernelFunc ()</code>	device	host
<code>__host__ float HostFunc ()</code>	host	host

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

7



## Compiling A CUDA Program



8





GPU Teaching Kit  
Accelerated Computing



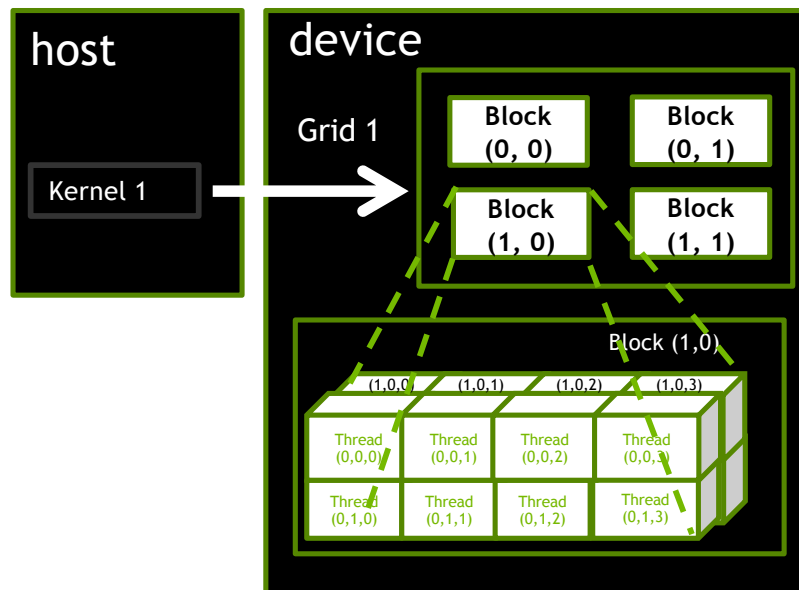
## Lecture 3.2 – CUDA Parallelism Model

Multidimensional Kernel Configuration

### Objective

- To understand multidimensional Grids
  - Multi-dimensional block and thread indices
  - Mapping block/thread indices to data indices

# A Multi-Dimensional Grid Example

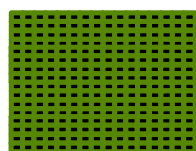


11

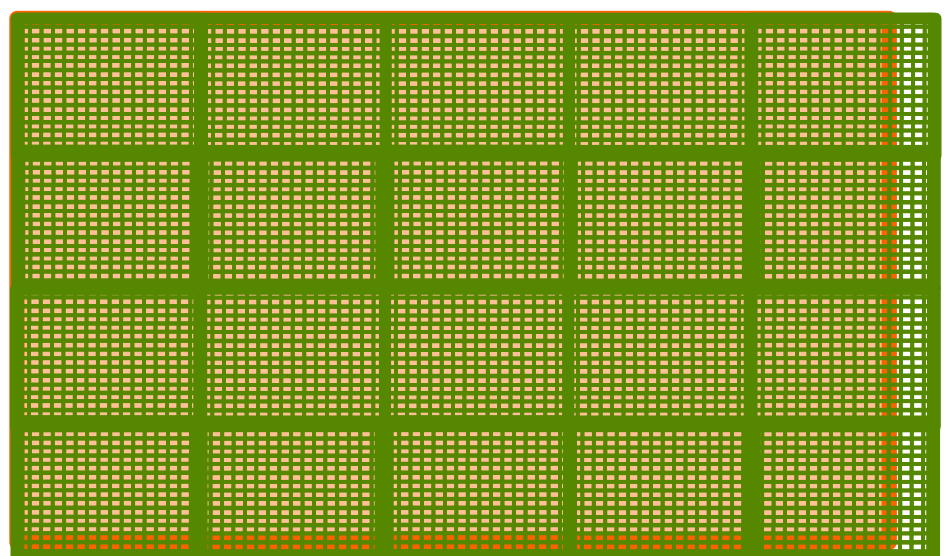
NVIDIA

ILLINOIS

# Processing a Picture with a 2D Grid



16×16 blocks



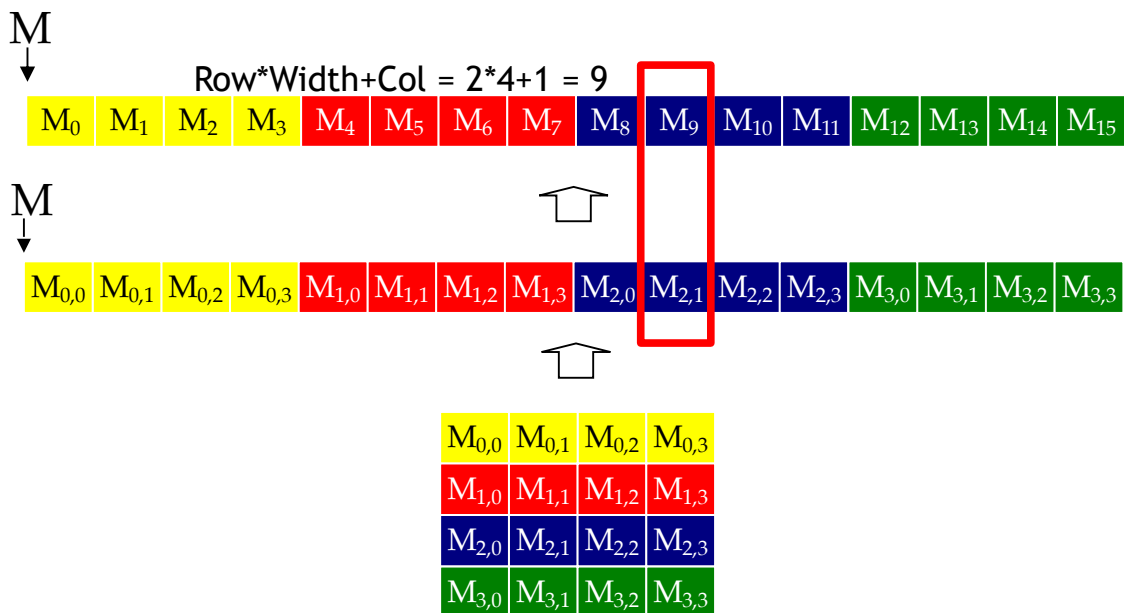
62×76 picture

12

NVIDIA

ILLINOIS

# Row-Major Layout in C/C++



## Source Code of a PictureKernel

```
__global__ void PictureKernel(float* d_Pin, float* d_Pout,
                             int height, int width)
{
    // Calculate the row # of the d_Pin and d_Pout element
    int Row = blockIdx.y*blockDim.y + threadIdx.y;

    // Calculate the column # of the d_Pin and d_Pout element
    int Col = blockIdx.x*blockDim.x + threadIdx.x;

    // each thread computes one element of d_Pout if in range
    if ((Row < height) && (Col < width)) {
        d_Pout[Row*width+Col] = 2.0*d_Pin[Row*width+Col];
    }
}
```

Scale every pixel value by 2.0

# Host Code for Launching PictureKernel

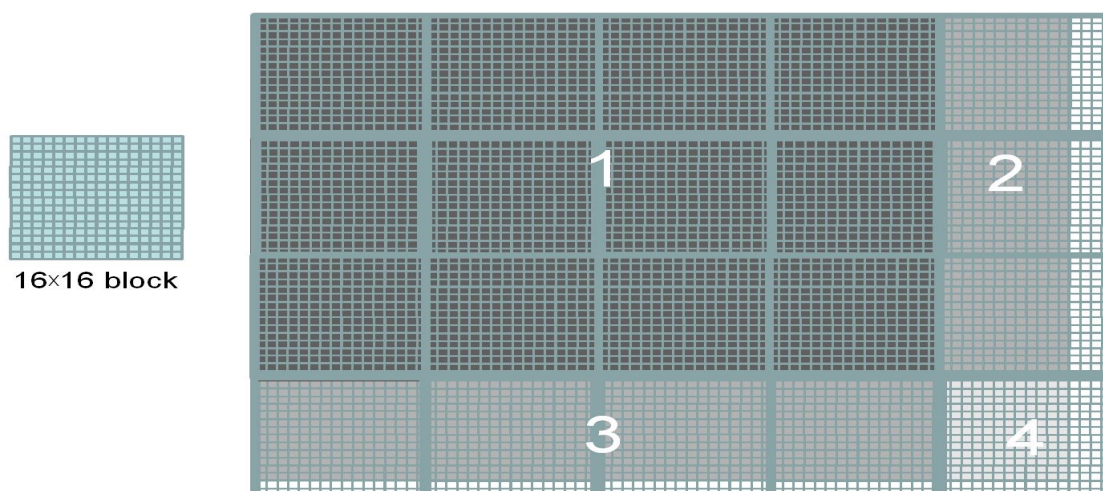
```
// assume that the picture is m × n,  
// m pixels in y dimension and n pixels in x dimension  
// input d_Pin has been allocated on and copied to device  
// output d_Pout has been allocated on device  
...  
dim3 DimGrid((n-1)/16 + 1, (m-1)/16+1, 1);  
dim3 DimBlock(16, 16, 1);  
PictureKernel<<<DimGrid,DimBlock>>>(d_Pin, d_Pout, m, n);  
...
```

15

NVIDIA

ILLINOIS

## Covering a 62×76 Picture with 16×16 Blocks



Not all threads in a Block will follow the same control flow path.

16

NVIDIA

ILLINOIS