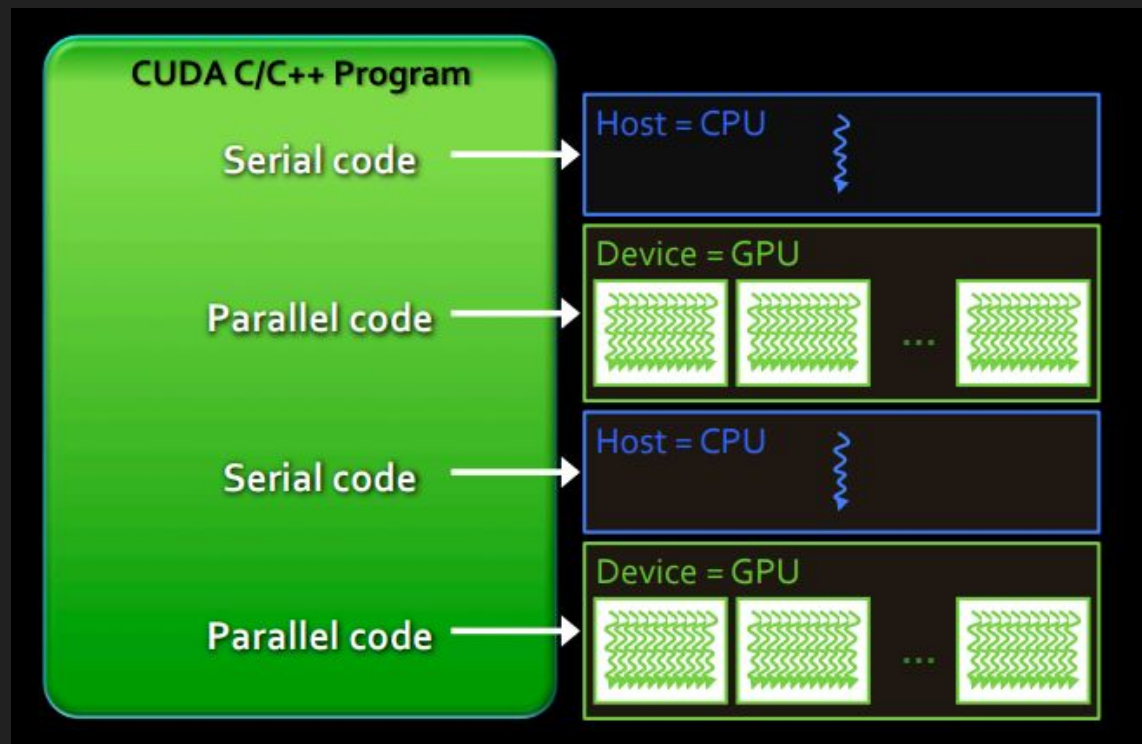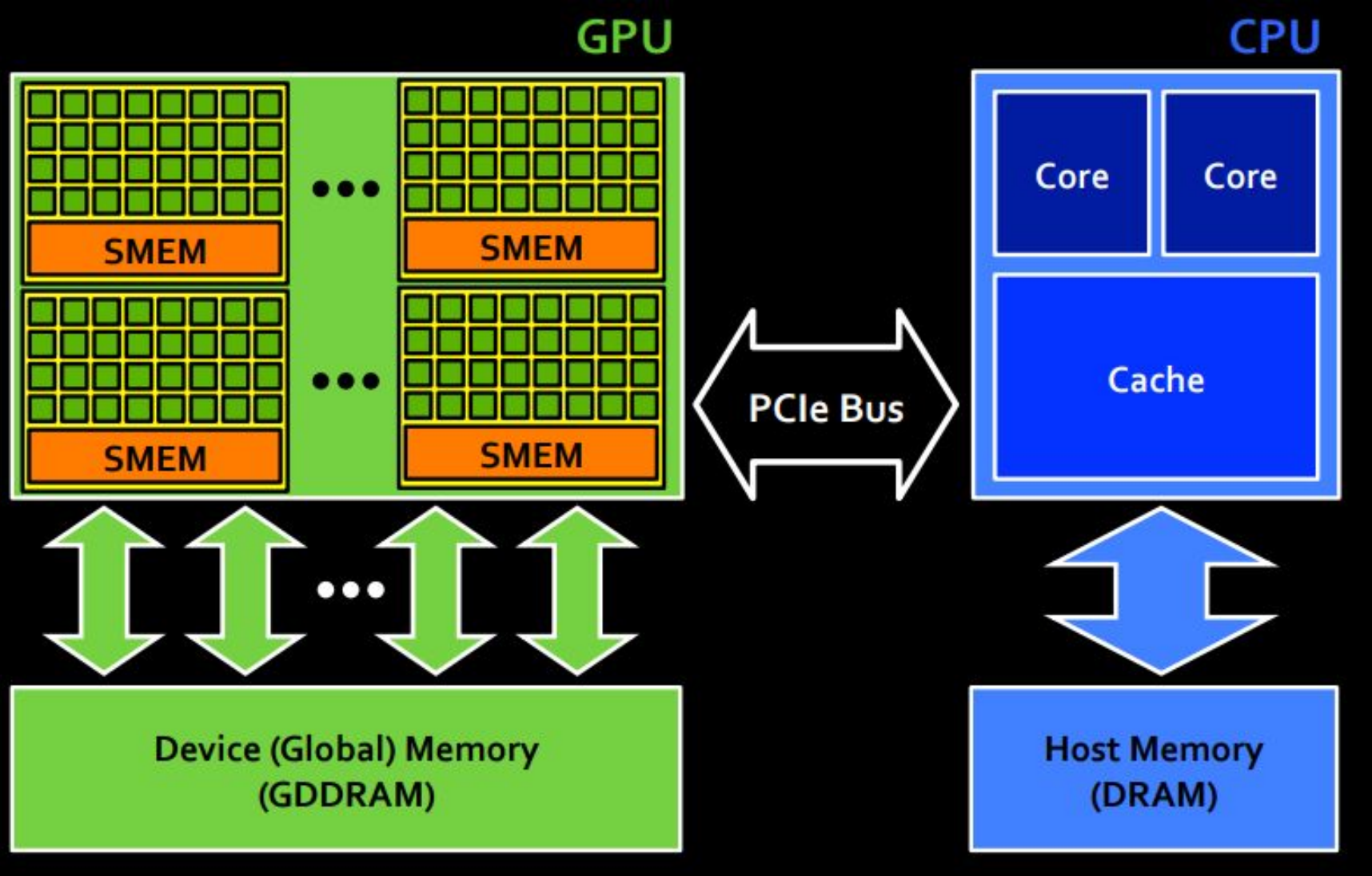# Introduction to CUDA

# CUDA

- Programing system for machines with GPUs
  - Programming Language
  - Compilers
  - Runtime Environments
  - Drivers
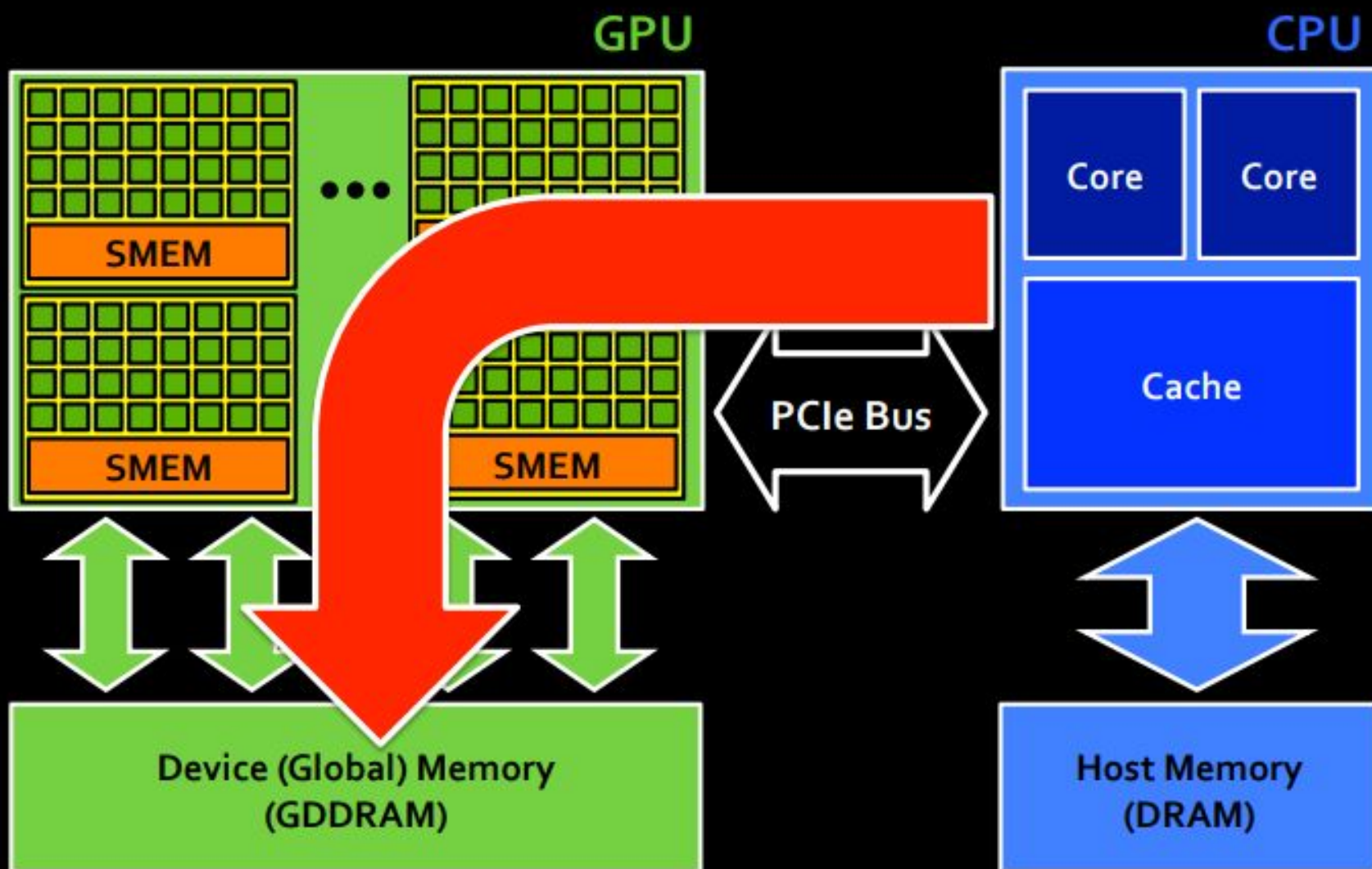  - Hardware

# Behavior of CUDA program

- **Serial** code executes in Host (CPU) thread
- **Parallel** code executes in many concurrent Device (GPU) threads across multiple parallel processing elements
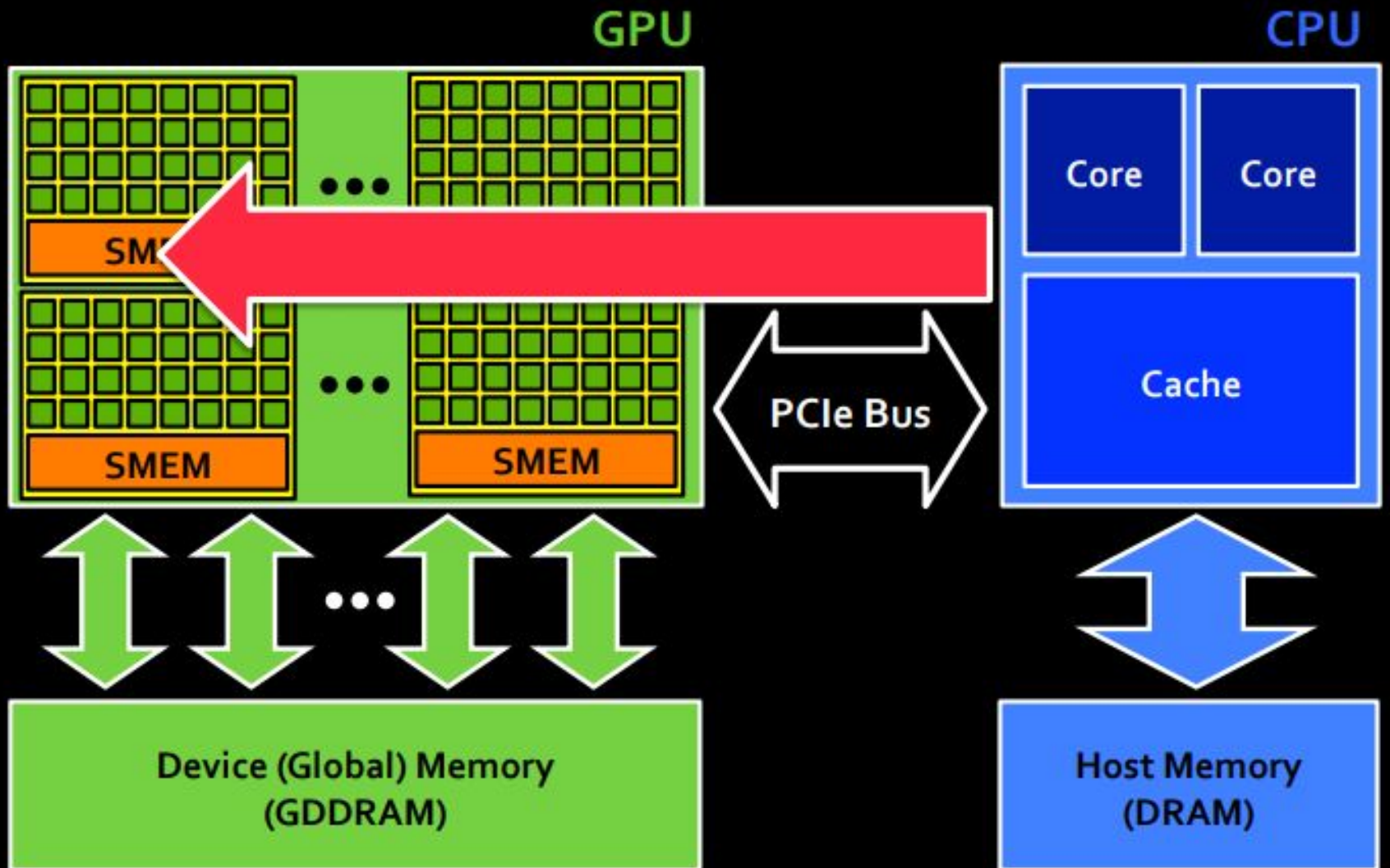
# Execution flow

# Step 2 – launch kernel on GPU

Step 3 – execute kernel on GPU

# CUDA ARCHITECTURE

# CUDA Thread Organization



- **GPUs can handle thousands of concurrent threads**
- **CUDA programming model supports even more**
  — Allows a kernel launch to specify more threads than the GPU can execute concurrently
  — Helps to amortize kernel launch times

# Blocks of threads



- Threads are grouped into blocks

# Grids of blocks



- Threads are grouped into blocks
- Blocks are grouped into a grid
- A kernel is executed as a grid of blocks of threads

# Blocks execute on Streaming Multiprocessors

**Streaming Processor**

**Streaming Multiprocessor**

SMEM

**Thread**

Registers

Global (Device) Memory

**Thread Block**

Per-Block Shared Memory

Global (Device) Memory

Grids of blocks executes across GPU

# CUDA Memory Hierarchy

- **Thread**
  - Registers
  - Local memory

| Thread | Thread | Thread | Thread |
|--------|--------|--------|--------|
| Registers | Registers | Registers | Registers |
| Local | Local | Local | Local |

# CUDA Memory Hierarchy

- **Thread**
  - Registers
  - Local memory

- **Thread Block**
  - Shared memory

# CUDA Memory Hierarchy

- **Thread**
  - Registers
  - Local memory

- **Thread Block**
  - Shared memory

# CUDA Memory Hierarchy

- **Thread**
  - — Registers
  - — Local memory

- **Thread Block**
  - — Shared memory

- **All Thread Blocks**
  - — Global Memory

# Thread and Block ID and Dimensions

- **Threads**
  - 3D IDs, unique within a block
- **Thread Blocks**
  - 2D IDs, unique within a grid
- **Dimensions set at launch**
  - Can be unique for each grid
- **Built-in variables**
  - `threadIdx`, `blockIdx`
  - `blockDim`, `gridDim`
- **Programmers usually select dimensions that simplify the mapping of the application data to CUDA threads**

Device

Grid 1

| Block (0, 0) | Block (1, 0) | Block (2, 0) |
| Block (0, 1) | Block (1, 1) | Block (2, 1) |

Block (1, 1)

| Thread (0, 0) | Thread (1, 0) | Thread (2, 0) | Thread (3, 0) | Thread (4, 0) |
| Thread (0, 1) | Thread (1, 1) | Thread (2, 1) | Thread (3, 1) | Thread (4, 1) |
| Thread (0, 2) | Thread (1, 2) | Thread (2, 2) | Thread (3, 2) | Thread (4, 2) |

# Indexing Arrays With Threads And Blocks

- No longer as simple as just using `threadIdx.x` or `blockIdx.x` as indices

- To index array with 1 thread per entry (using 8 threads/block)

| threadIdx.x | | | | | | | | threadIdx.x | | | | | | | | threadIdx.x | | | | | | | | threadIdx.x | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

blockIdx.x = 0        blockIdx.x = 1        blockIdx.x = 2        blockIdx.x = 3

- If we have `M` threads/block, a unique array index for each entry given by

```
int index = threadIdx.x + blockIdx.x * M;
int index =      x       +      y      * width;
```

# Indexing Arrays: Example

- In this example, the red entry would have an index of 21:

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | | |

M = 8 threads/block

threadIdx.x = 5

blockIdx.x = 2

```
int index = threadIdx.x + blockIdx.x * M;
          =      5       +      2        * 8;
          = 21;
```

# CUDA C

```c
#include <stdio.h>

__global__ void print_kernel() {
    printf("Hello from block %d, thread %d\n", blockIdx.x, threadIdx.x);
}

int main() {
    print_kernel<<<10, 10>>>();
    cudaDeviceSynchronize();
}
```

```
Hello from block 1, thread 0
Hello from block 1, thread 1
Hello from block 1, thread 2
Hello from block 1, thread 3
Hello from block 1, thread 4
Hello from block 1, thread 5
....
Hello from block 8, thread 3
Hello from block 8, thread 4
Hello from block 8, thread 5
Hello from block 8, thread 6
Hello from block 8, thread 7
Hello from block 8, thread 8
Hello from block 8, thread 9
```

# Parallel Programming in CUDA C

- With `add()` running in parallel...let's do vector addition

- Terminology: Each parallel invocation of `add()` referred to as a *block*

- Kernel can refer to its block's index with the variable `blockIdx.x`

- Each block adds a value from `a[]` and `b[]`, storing the result in `c[]`:

```
__global__ void add( int *a, int *b, int *c ) {
    c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];
}
```

- By using `blockIdx.x` to index arrays, each block handles different indices

# Parallel Programming in CUDA C

- We write this code:

```
__global__ void add( int *a, int *b, int *c ) {
    c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];
}
```

- This is what runs in parallel on the device:

Block 0
```
c[0] = a[0] + b[0];
```

Block 1
```
c[1] = a[1] + b[1];
```

Block 2
```
c[2] = a[2] + b[2];
```

Block 3
```
c[3] = a[3] + b[3];
```

# Parallel Addition: `main()`

```
#define N  512
int main( void ) {
    int *a, *b, *c;                 // host copies of a, b, c
    int *dev_a, *dev_b, *dev_c;     // device copies of a, b, c
    int size = N * sizeof( int );   // we need space for 512 integers

    // allocate device copies of a, b, c
    cudaMalloc( (void**)&dev_a, size );
    cudaMalloc( (void**)&dev_b, size );
    cudaMalloc( (void**)&dev_c, size );

    a = (int*)malloc( size );
    b = (int*)malloc( size );
    c = (int*)malloc( size );

    random_ints( a, N );
    random_ints( b, N );
```

# Parallel Addition: `main()` (cont)

```c
    // copy inputs to device
    cudaMemcpy( dev_a, a, size, cudaMemcpyHostToDevice );
    cudaMemcpy( dev_b, b, size, cudaMemcpyHostToDevice );

    // launch add() kernel with N parallel blocks
    add<<< N, 1 >>>( dev_a, dev_b, dev_c );

    // copy device result back to host copy of c
    cudaMemcpy( c, dev_c, size, cudaMemcpyDeviceToHost );

    free( a ); free( b ); free( c );
    cudaFree( dev_a );
    cudaFree( dev_b );
    cudaFree( dev_c );
    return 0;
}
```

# CUDA using Python

- Anaconda/Python 3.6.1/Jupyter notebook
- CUDA Toolkit
- Numba package

# CUDA using Python

# CUDA using Python

```
!apt-get install nvidia-cuda-toolkit
```
```
  fonts-ipafont-mincho fonts-wqy-microhei fonts-wqy-zenhei fonts-indic
  libvdpau-va-gl1 nvidia-vdpau-driver nvidia-legacy-340xx-vdpau-driver
  mesa-utils
Recommended packages:
  libnvcuvid1
The following NEW packages will be installed:
  adwaita-icon-theme at-spi2-core ca-certificates-java cpp-6
  dconf-gsettings-backend dconf-service fontconfig fonts-dejavu-core
  fonts-dejavu-extra g++-6 gcc-6 gcc-6-base glib-networking
```

```
!nvcc --version
```
```
nvcc: NVIDIA (R) Cuda compiler driver
Copyright (c) 2005-2018 NVIDIA Corporation
Built on Tue_Jun_12_23:07:04_CDT_2018
Cuda compilation tools, release 9.2, V9.2.148
```

```
!pip3 install numba
```
```
Collecting numba
  Downloading https://files.pythonhosted.org/packages/42/45/8d5fc45e5f760ac65906ba48dec98e99e7920c96783ac7248c5e31c9464e/numba-0.40.1-cp36-cp36m-manylinux1_x8
    100% |████████████████████████████████| 3.2MB 8.3MB/s
Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-packages (from numba) (1.14.6)
Collecting llvmlite>=0.25.0dev0 (from numba)
  Downloading https://files.pythonhosted.org/packages/34/fb/f9c2e9e0ef2b54c52f0b727cf6af75b68c3d7ddb6d88c8d557b1b16bc1ab/llvmlite-0.25.0-cp36-cp36m-manylinux1
    100% |████████████████████████████████| 16.1MB 2.8MB/s
Installing collected packages: llvmlite, numba
Successfully installed llvmlite-0.25.0 numba-0.40.1
```

# Vector add GPU

```python
from __future__ import print_function
from timeit import default_timer as time
import numpy as np
from numba import cuda

@cuda.jit('(f4[:], f4[:], f4[:])')
def cuda_sum(a, b, c):
    i = cuda.grid(1)
    c[i] = a[i] + b[i]

griddim = 50, 1
blockdim = 32, 1, 1
N = griddim[0] * blockdim[0]
print("N", N)
cuda_sum_configured = cuda_sum.configure(griddim, blockdim)
a = np.array(np.random.random(N), dtype=np.float32)
b = np.array(np.random.random(N), dtype=np.float32)
c = np.empty_like(a)

ts = time()
cuda_sum_configured(a, b, c)
te = time()
print(te - ts)


assert (a + b == c).all()
print(c)
```

http://numba.pydata.org/numba-doc/0.13/CUDAJit.html

# Vector add CPU

```python
from timeit import default_timer as time
import numpy as np
N = 1600
def cpu_sum(a, b, c):
    for i in range(0, N):
        c[i] = a[i] + b[i]


a = np.array(np.random.random(N), dtype=np.float32)
b = np.array(np.random.random(N), dtype=np.float32)
c = np.empty_like(a)

ts = time()
cpu_sum(a, b, c)
te = time()
print(te - ts)

print(c)
```



GPU and CPU