

# GPU Programming Using CUDA

Michael J. Schnieders

Depts. of Biomedical Engineering & Biochemistry

The University of Iowa

&

Gregory G. Howes

Department of Physics and Astronomy

The University of Iowa

Iowa High Performance Computing Summer School

The University of Iowa

Iowa City, Iowa

1-3 June 2015



# Outline

---



- Concepts for GPU Computing
- Programming Model for GPU Computing using CUDA C
- CUDA C Programming
- Advanced CUDA Capabilities

# GPU Computing

---



**Graphics Processing Units** (GPUs) have been developed in response to strong market demand for realtime, high-definition 3D graphics (**video games!**)

**GPUs** are **highly parallel**,  
**multithreaded**,  
**manycore** processors

- Tremendous computational horsepower
- Very high memory bandwidth

**We hope to access this power for scientific computing**

# GPU Programming Languages



- **CUDA** (Compute Unified Device Architecture) is the proprietary programming language for NVIDIA GPUs
- **OpenCL** (Open Computing Language) is portable language standard for general computing that can exploit capabilities of GPUs from any manufacturer.
- **OpenACC** (Open Accelerators) is portable like **OpenCL**, but features a directive syntax that is compatible with **OpenMP**

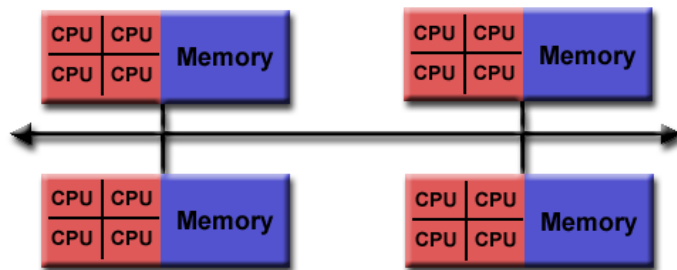
Each language provide extensions to **C** (as well as other languages) that enable the programmers to access the powerful computing capability for **general-purpose computing** on GPUs (**GPGPU**)

Today we will focus on the basics of **CUDA C programming**

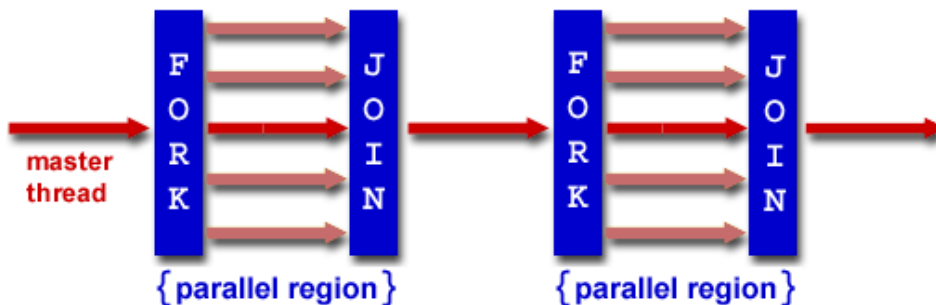
# Parallel Computing Architectures

Different computer architectures suggest three approaches to parallel computing:

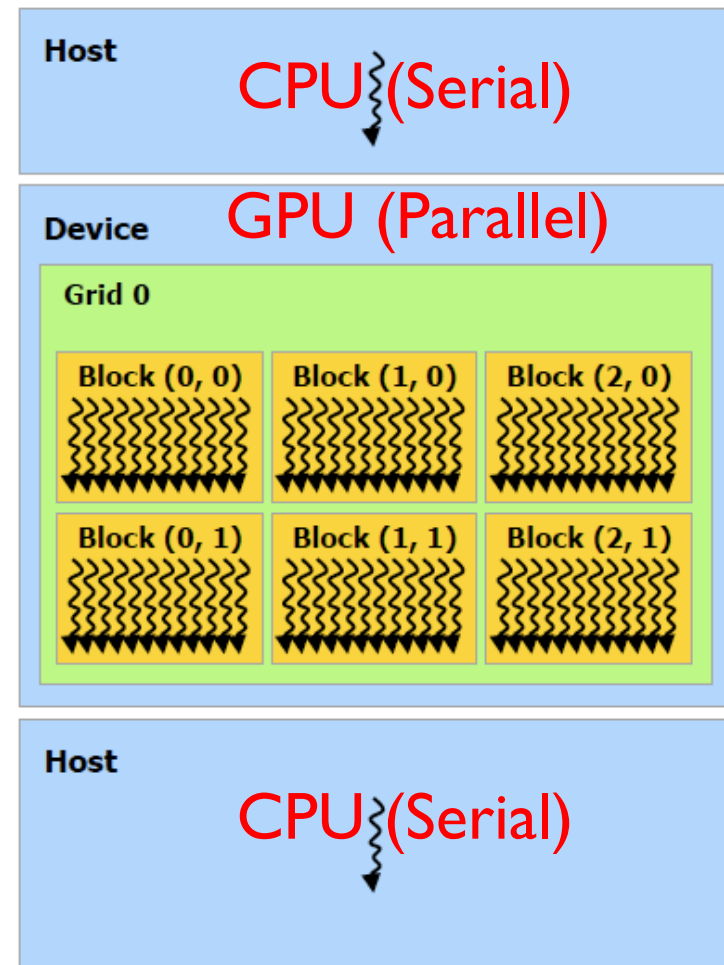
## 1) Message Passing (MPI)



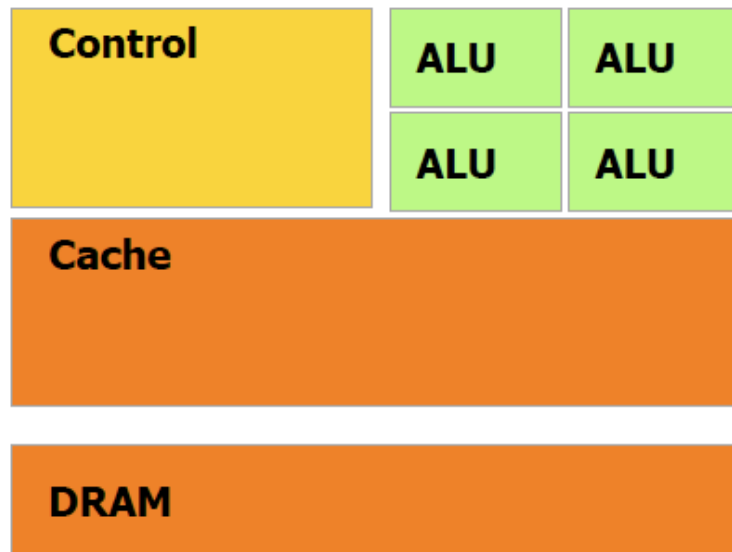
## 2) Multithreading (OpenMP)



## 3) GPU (CUDA)

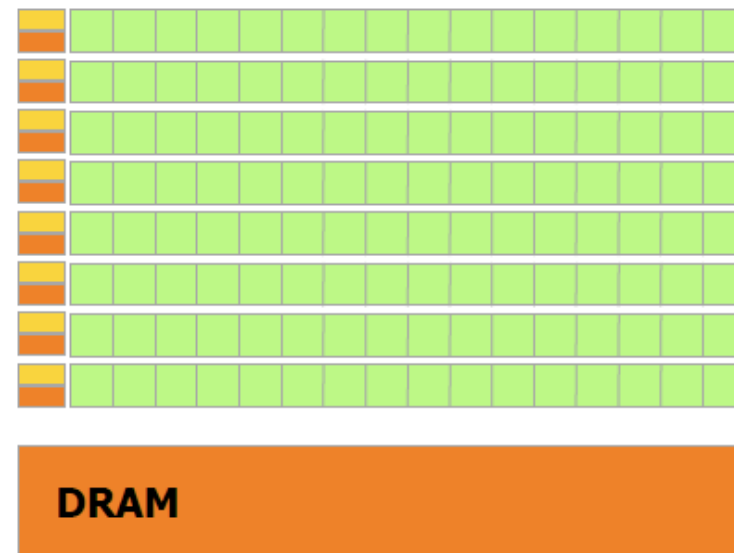


# More Transistors to Data Processing



**CPU**

CPUs devote a significant fraction of transistors to **data caching** and **flow control**



**GPU**

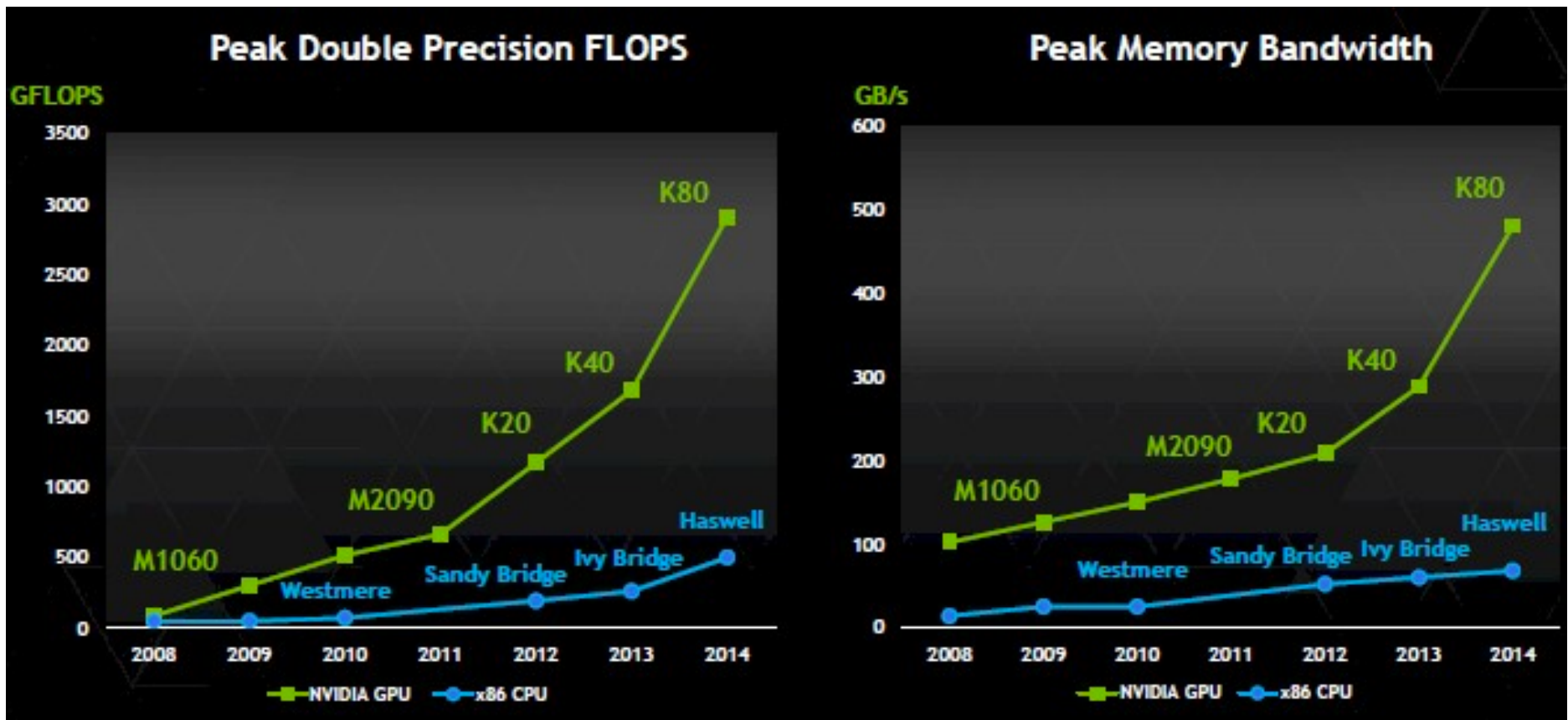
GPUs devote more transistors to **data processing** (arithmetic and logic units, ALU)

# GPU vs. CPU Peak Performance



**F**loating point **O**perations **P**er **S**econd  
(FLOPS)

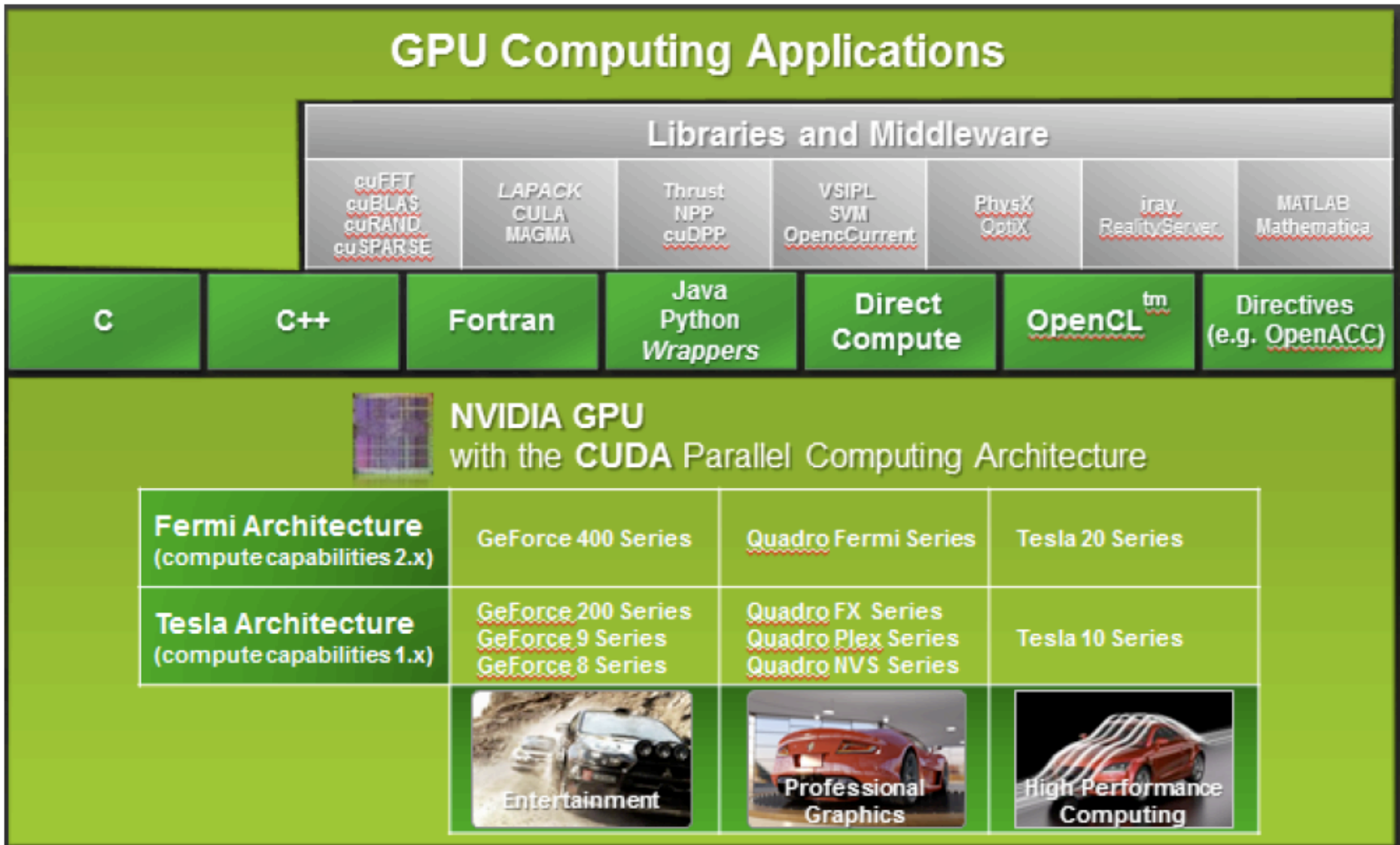
Memory Access Bandwidth  
(GB/s)



# Many Languages for GPU Computing



**NVIDIA**'s diverse offerings for GPU Computing





# Outline

---



- Concepts for GPU Computing
- Programming Model for GPU Computing using CUDA C
- CUDA C Programming
- Advanced CUDA Capabilities

# CUDA C Programming



## CUDA C Programming Language

- Minimal set of extensions to the C programming language
- Core concepts:
  - Hierarchy of thread groups
  - Shared memory
  - Barrier synchronization

### 1) Kernel:

- C function executed  $N$  times in parallel by  $N$  CUDA threads
- Called by the Host (CPU) but executed on the Device (GPU)

### 2) Thread Hierarchy:

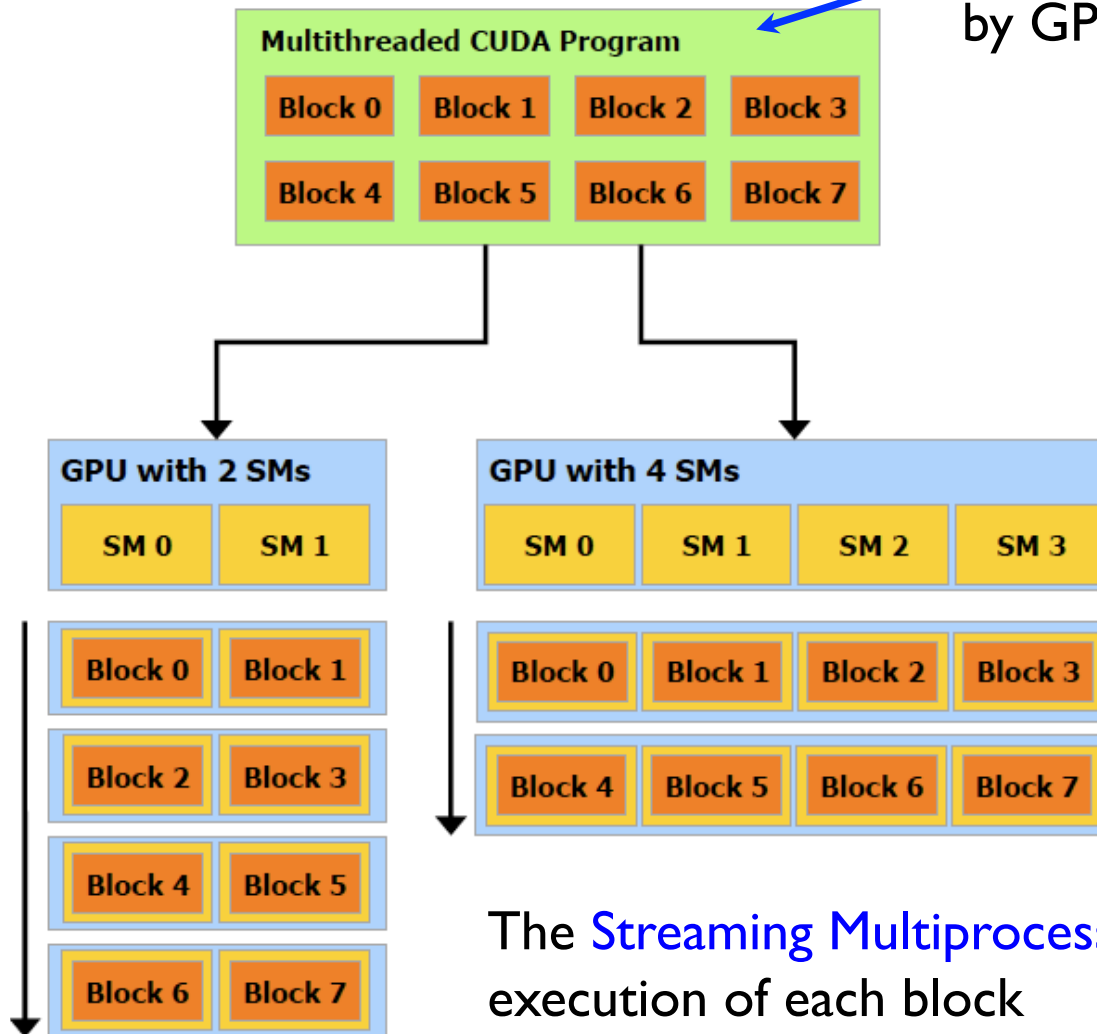
- Grid: Contains many blocks that can be solved independently in parallel
- Block: Contains many threads that can be solve cooperatively in parallel

### 3) Memory Functions:

- Allocate and Free memory space on Device (GPU),
- Copy data from Host (CPU) to Device (GPU), and vice versa

# Thread Hierarchy: Grid of Blocks

Grid: Full problem to be solved by GPU is broken into **blocks**

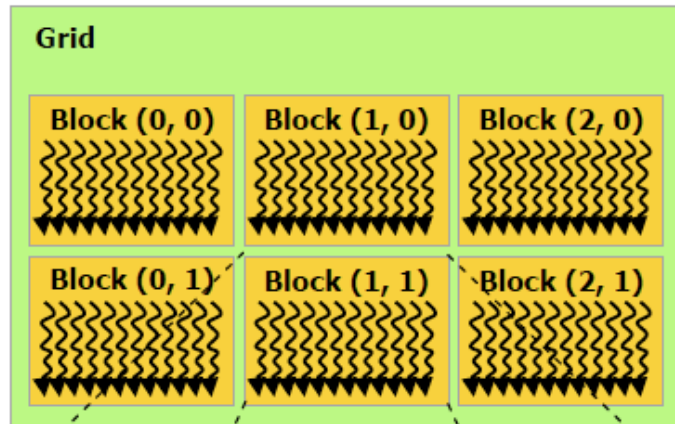


Each **block** is **independent**.  
The blocks can be executed, **in any order**, **concurrently** or **sequentially**

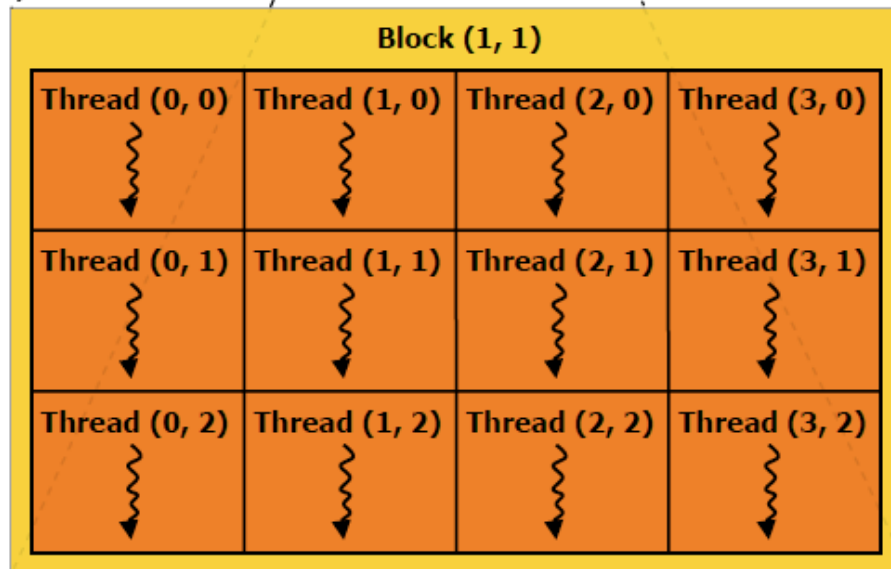
The **Streaming Multiprocessors (SMs)** control the execution of each block

This model enables excellent **scalability** for a varying number of cores per GPU

# Thread Hierarchy: Block of Threads



Each **block** contains many **threads**



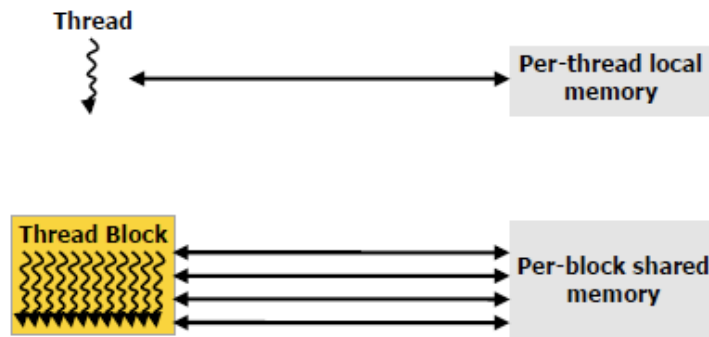
**Threads** within a **block** are executed in parallel, either **cooperatively** or **independently**

# Memory Hierarchy

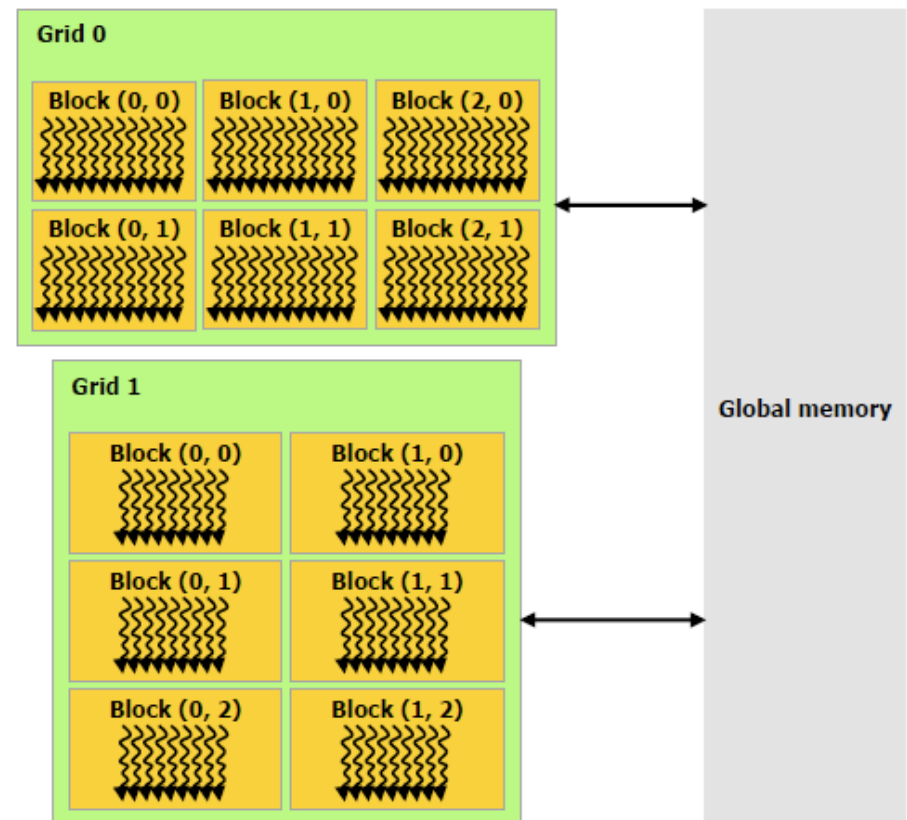
Very high memory bandwidth can be achieved using a hierarchy of memory

All threads have slower access to global memory

Each thread has private local memory

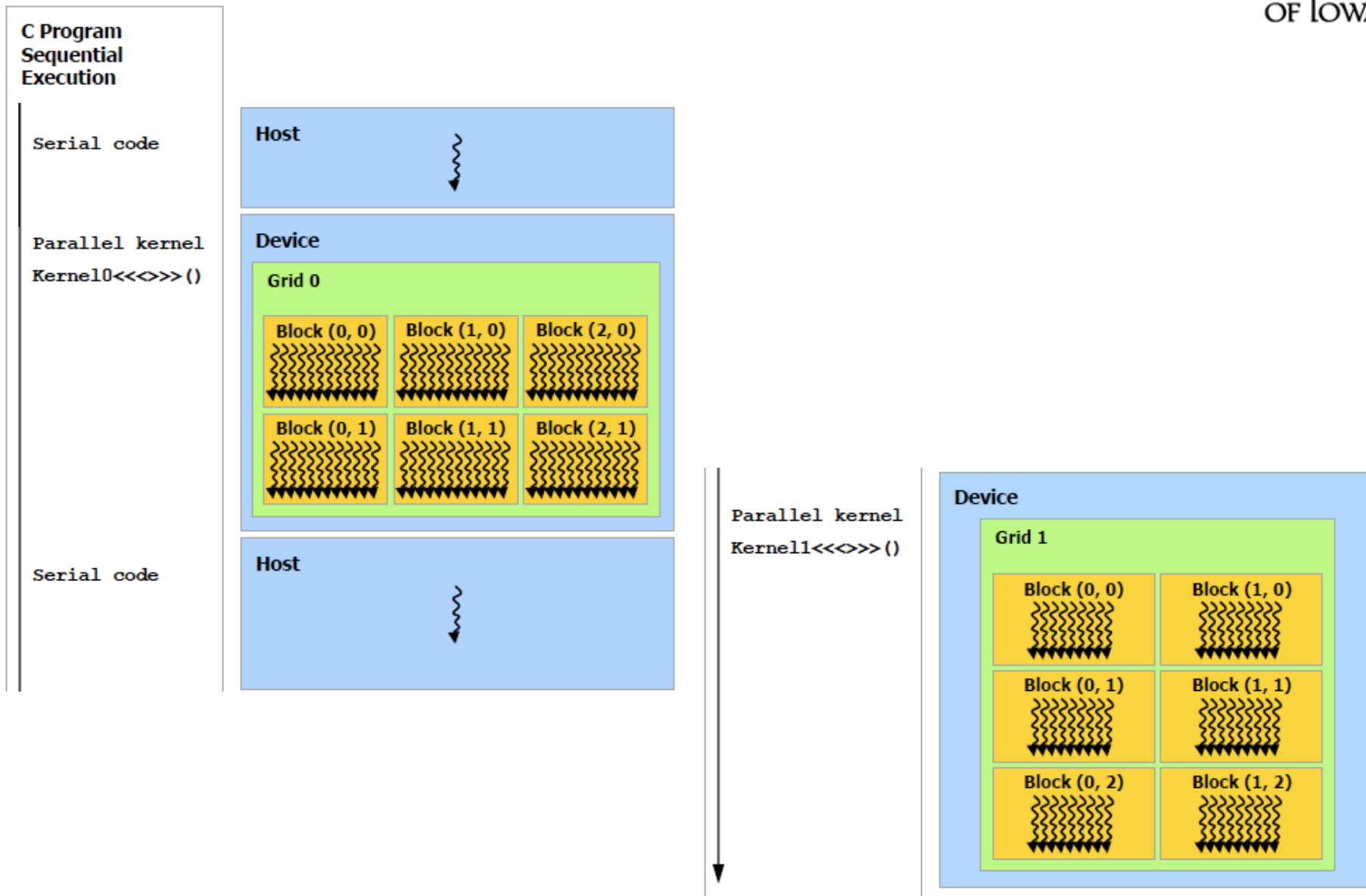


Each thread block has fast access to shared memory



All threads can also access read-only **Constant** and **Texture** memory, optimized for different memory usages

# General Flow of CUDA C Program



# Outline

---



- Concepts for GPU Computing
- Programming Model for GPU Computing using CUDA C
- **CUDA C Programming**
- Advanced CUDA Capabilities

# Programming in CUDA C

---



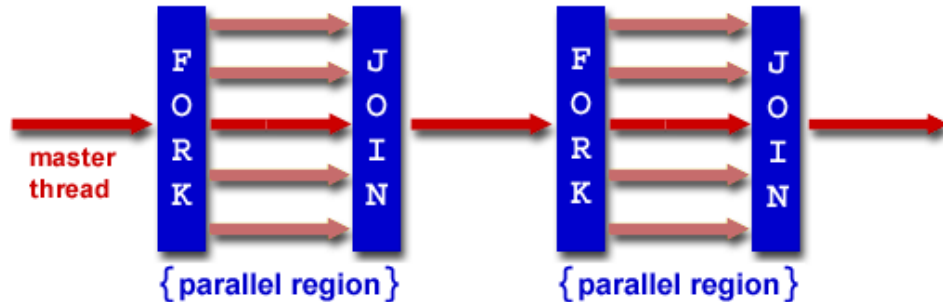
- Comparison of Multithreading and GPU Computing
- The Kernel
  - Thread Hierarchy
- Memory Functions
- Examples



# Comparison to OpenMP Multithreading

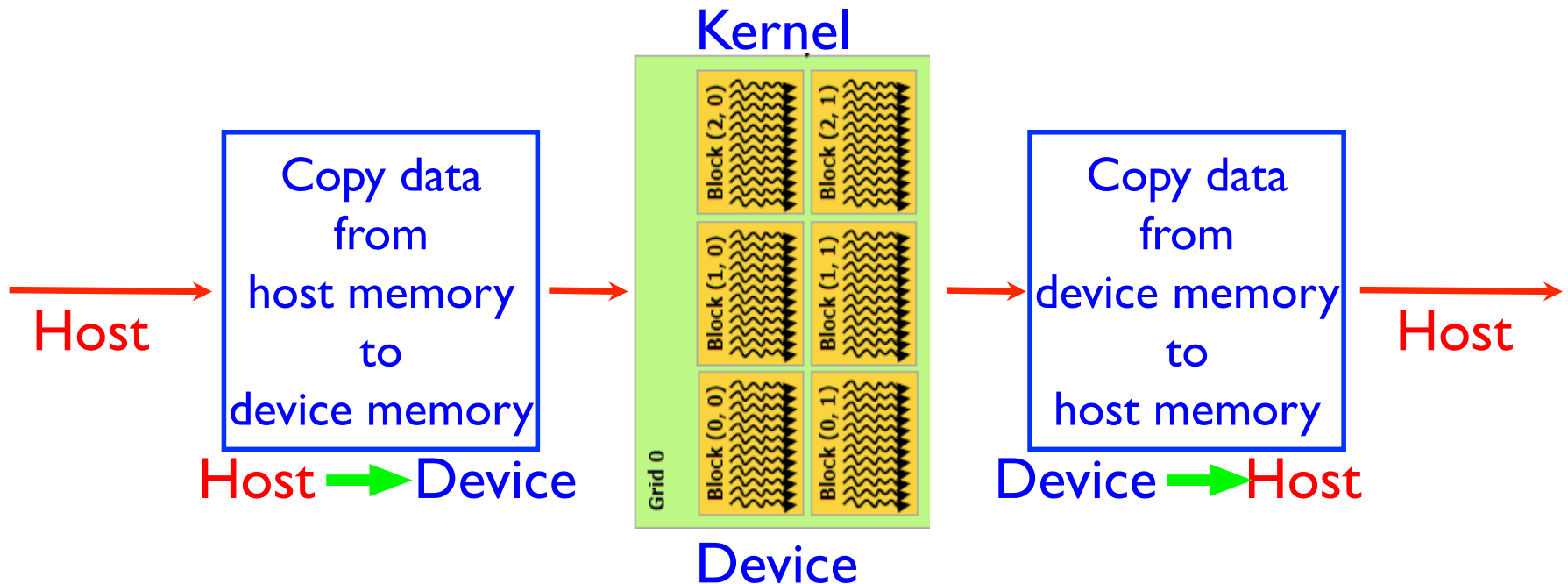
General Implementation:

- Multithreading using OpenMP



All of this executes on the multicore host CPU with access to shared memory

- Parallel Computation on GPU device



# The Kernel in CUDA C



## Kernels:

- CUDA C enables the programmer to define C functions, called **kernels**, that are executed N times in parallel on the GPU device
- Declaration specifier: Kernels are defined using `__global__`

```
__global__ void VecAdd(int *a, int *b, int *c)
```

- Kernels are called from the host (CPU) using a new **execution configuration** syntax,

```
<<<blocksPerGrid, threadsperBlock>>>
```

```
VecAdd<<<blocksPerGrid, threadsperBlock>>>(d_a, d_b, d_c);
```

**NOTE:** Kernel functions are called from the **host**,  
but executed on the **device**!

# Blocks and Threads



## Thread Hierarchy:

- The execution configuration specifies:

```
blocksPerGrid  
threadsperBlock
```

- These variables are 3-component integer vectors of type `dim3`

- For example,

```
dim3 threadsperBlock(N,N);  
defines 2D thread blocks of size  $N$  by  $N$ 
```

- In the kernel function, the thread ID is accessed through the variable

```
threadIdx, where the two dimensions are given by  
threadIdx.x and  
threadIdx.y
```

# Blocks and Threads



## Blocks:

- The block ID and block dimensions are similarly accessed through `blockIdx` giving `blockIdx.x` and `blockIdx.y`  
`blockDim` giving `blockDim.x` and `blockDim.y`
- A general formula for computing the appropriate index based on multiple blocks is

```
int i = blockIdx.x*blockDim.x + threadIdx.x;
```

## Limitations:

- The maximum number of threads per block is the number of GPU cores, 512 for the GeForce GTX580.
- The number of blocks per grid is unlimited, and is determined by the number of blocks required to do the entire calculation.
- For a 2D computation of size  $N$  by  $N$   
`dim3 numBlocks(N/threadsPerBlock.x, N/threadsPerBlock.y);`

# CUDA Kernel



## Standard serial C function

```
/* Function to compute y=a*x+y */
void saxpy_serial(int n, float a, float *x, float *y){
    for (int i=0; i < n; i++)
        y[i]= a*x[i] + y[i];
}

/* Call Function from main() */
saxpy_serial(4096*256, 2.0, x, y);
```

## Parallel CUDA C kernel

```
__global__ void saxpy_parallel(int n, float a, float *x, float *y)
{
    int i=blockIdx.x*blockDim.x + threadIdx.x;
    if (i < n)  y[i]= a*x[i] + y[i];
}

/* Call Function from main() */
saxpy_parallel<<<4096,256>>>(4096*256, 2.0, x, y);
```

# CUDA Kernel



## General Comments:

- The kernel contains only the commands within the loop
- The kernel call is asynchronous
- After the kernel is called, the host can continue processing before the GPU has completed the kernel computation

	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

- The computations in the kernel can only access data in device memory

Therefore, a critical part of CUDA programming is handling the transfer of data from host memory to device memory and back!

# CUDA Memory Functions



- Device memory is allocated and freed using

```
cudaMalloc()
```

```
cudaFree()
```

**Example:**

```
size_t size = N * sizeof(float);
```

```
float* d_A;
```

```
cudaMalloc(&d_A, size);
```

- Data is transferred using

```
cudaMemcpy()
```

**Example:**

```
/* Allocate array in host memory */
```

```
float* h_A = (float*)malloc(size);
```

```
/* Copy array from host memory to device memory */
```

```
cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
```

# Example CUDA code for Vector Addition



```
#include<stdio.h>
#include<cuda.h>

#define N 100 /* Size of vectors */

/* Define CUDA kernel */
__global__ void add( int *a, int *b, int *c ) {
    int tid = blockIdx.x*blockDim.x+threadIdx.x;
    if (tid < N)
        c[tid] = a[tid] + b[tid];
}
```



# Example CUDA code for Vector Addition



```
int main() {
    int a[N], b[N], c[N];
    int *dev_a, *dev_b, *dev_c;
    /* allocate the memory on the GPU */
    cudaMalloc( (void**)&dev_a, N * sizeof(int) );
    cudaMalloc( (void**)&dev_b, N * sizeof(int) );
    cudaMalloc( (void**)&dev_c, N * sizeof(int) );
    /* Copy the arrays 'a' and 'b' from CPU host to GPU device*/
    cudaMemcpy( dev_a, a, N * sizeof(int), cudaMemcpyHostToDevice );
    cudaMemcpy( dev_b, b, N * sizeof(int), cudaMemcpyHostToDevice );
    int threadsPerBlock=512;
    int blocksPerGrid=(N+threadsPerBlock-1)/threadsPerBlock;
    add<<<blocksPerGrid,threadsPerBlock>>>( dev_a, dev_b, dev_c );
    /* Copy the array 'c' back from GPU device to CPU host*/
    cudaMemcpy( c, dev_c, N * sizeof(int), cudaMemcpyDeviceToHost );
    /* Free the memory allocated on the GPU device*/
    cudaFree( dev_a );
    cudaFree( dev_b );
    cudaFree( dev_c );
}
```

# Outline

---



- Concepts for GPU Computing
- Programming Model for GPU Computing using CUDA C
- CUDA C Programming
- **Advanced CUDA Capabilities**

# Advanced CUDA Capabilities

---



- Shared Memory
- Concurrent memory copy and kernel execution
- Asynchronous concurrent execution
- Lower-level CUDA driver API
- Multiple devices on host system with peer-to-peer memory access
- Texture and surface memory
- Graphics functions with OpenGL and Direct3D Application Programming Interfaces (APIs)

# Starter Code for CUDA Vector Addition



```
#include <stdio.h>
#include <stdlib.h>
#include <math.h>

// CUDA kernel. Each thread takes care of one element of c
__global__ void vecAdd(double *a, double *b, double *c, int n)
{
    int id = ??? // Get our global thread ID
    // Make sure we do not go out of bounds
    if (id < n)
        c[id] = a[id] + b[id];
}

int main( int argc, char* argv[] )
{
    int n = 100000; // Size of vectors
    // Declare host vectors
    // Declare device input vectors
    size_t bytes = n*sizeof(double); // Size, in bytes, of each vector
    // Allocate memory for each vector on host
    // Allocate memory for each vector on GPU

    int i;
    // Initialize vectors on host
    for( i = 0; i < n; i++ ) {
        h_a[i] = sin(i)*sin(i);
        h_b[i] = cos(i)*cos(i);
    }
    cudaMemcpy( d_a, h_a, bytes, cudaMemcpyHostToDevice); // Copy host vectors to device

    int blockSize, gridSize;
    blockSize = 1024; // Number of threads in each thread block
    gridSize = (int)ceil((float)n/blockSize); // Number of thread blocks in grid
    // Execute the kernel
    // Copy array back to host
    // Sum up vector c and print result divided by n, this should equal 1 within error
    double sum = 0;
    for(i=0; i<n; i++)
        sum += h_c[i];
    printf("final result: %f\n", sum/n);

    // Release device memory
    // Release host memory
    return 0;
}
```

# Starter Code for CUDA Monte Carlo



```
#include <stdio.h>
#include <stdlib.h>
#include <cuda.h>
#include <curand.h>
#include <time.h>
__global__ void kernel(int* count_d, float* randomnums) {
}
int main(int argc, char* argv[])
{
    //NOTE: if threads and/or blocks is changed, niter needs to be changed to reflect
    //that change (niter=threads*blocks)
    int niter = 100000;
    float *randomnums;
    double pi;
    //Allocate the array for the random numbers
    cudaMalloc((void**)&randomnums, (2*niter)*sizeof(float));
    //Use CuRand to generate an array of random numbers on the device
    int status;
    curandGenerator_t gen;
    status = curandCreateGenerator(&gen, CURAND_RNG_PSEUDO_MRG32K3A);
    status |= curandSetPseudoRandomGeneratorSeed(gen, 4294967296ULL*time(NULL));
    status |= curandGenerateUniform(gen, randomnums, (2*niter));
    status |= curandDestroyGenerator(gen);
    if (status != CURAND_STATUS_SUCCESS) {
        printf("CuRand Failure\n");
        exit(EXIT_FAILURE);
    }
    int threads = 1000;
    int blocks = 100;
    int* count_d;
    int *count = (int*)malloc(blocks*threads*sizeof(int));
    unsigned int reducedcount = 0;
    //Allocate the array to hold a value (1,0) whether the point in is the circle (1) or not (0)
    cudaMalloc((void**)&count_d, (blocks*threads)*sizeof(int));
    //Launch the kernel
    kernel <<<blocks, threads>>> (count_d, randomnums);
    cudaDeviceSynchronize();
    //Copy the resulting array back
    int i = 0;
    // Reduce array into int
    // Free the cudaMalloc()'d arrays
    // Find the ratio
    pi = ((double)reducedcount/niter)*4.0;
    printf("Pi: %f\n", pi);
    return 0;
}
```