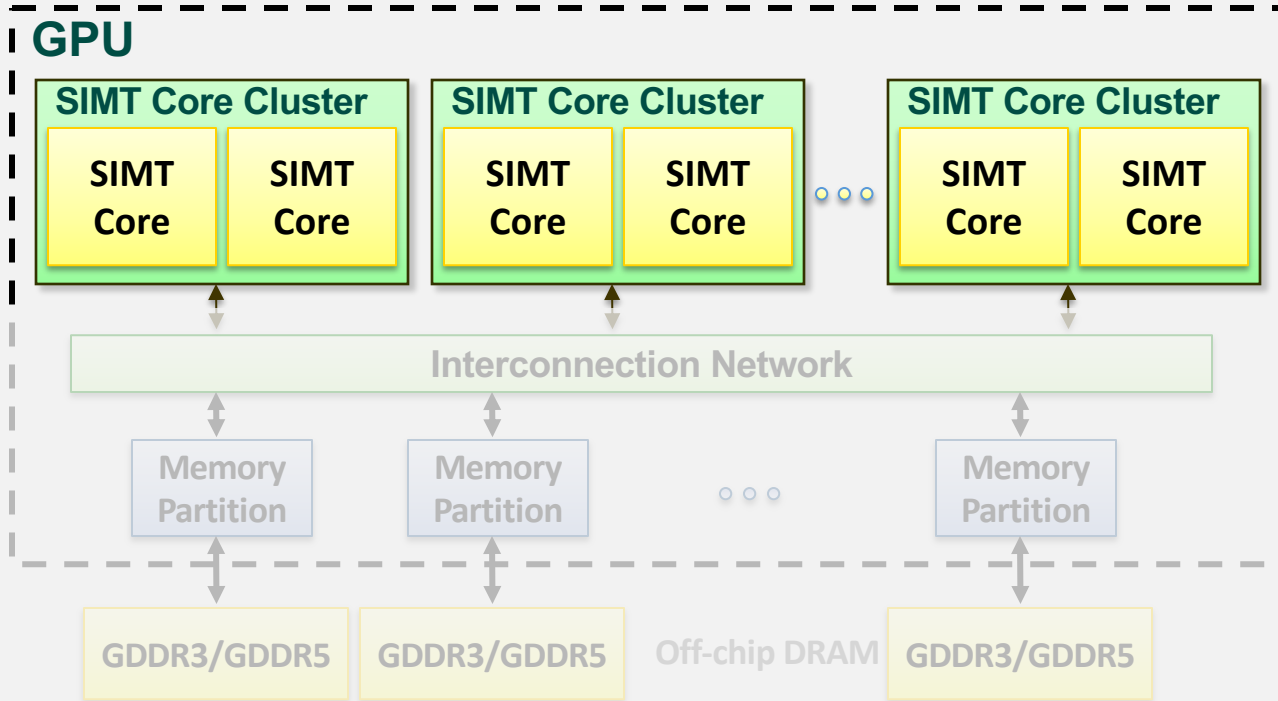# Introduction to NVIDIA CUDA Programming

2024 NSF CyberTraining Workshop

Jan. 8, 2024 – Jan. 19, 2024
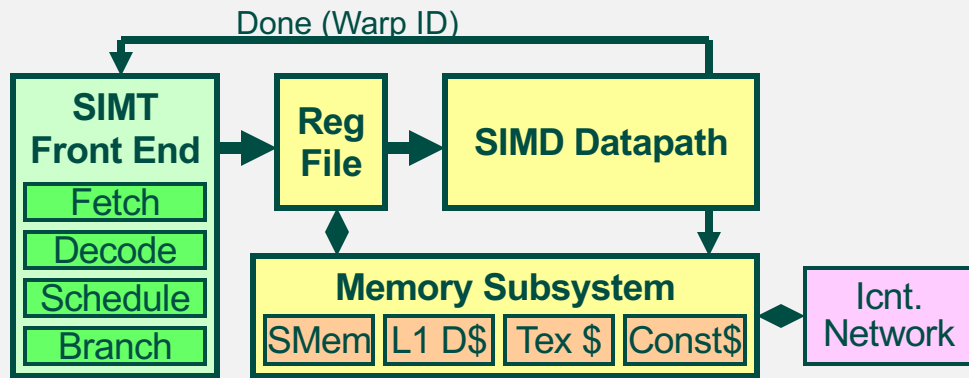
Clarkson University

Note: The lecture slides are adapted from the tutorial of CUDA programing from NVIDIA

# GPU Microarchitecture Overview

# Inside a SIMT Core


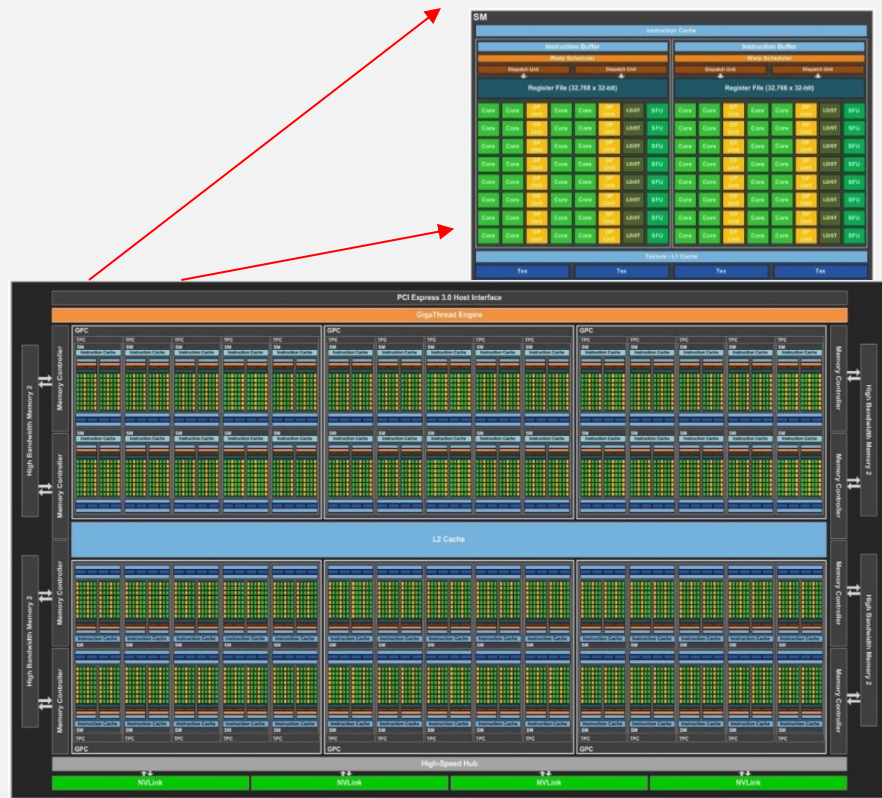
- Fine-grained multithreading
  - Interleave warp execution to hide latency
  - Register values of all threads stays in core

# Nvidia Pascal GP100 GPU

## Architecture

- 15.3 B Transistors @1.4 GHz clock speed
- Up to 60 "SM" units
- 32 "cuda cores" each
- Up to 5.7 TFlop/s peak
- 4 MB L2 Cache
- 4096-bit HBM2
- MemBW ~ 732 GB/s (theoretical)
- MemBW ~ 510 GB/s (measured)

# GPU vs. CPU

GPU vs. CPU

- ▪ Both are shared memory based arch.
- ▪ light speed estimate (per device)

MemBW  ~ 5-10x
Peak      ~ 6-15x



CPU

GPU

| | 2x Intel Xeon E5-2697v4 "Broadwell" | Intel Xeon Phi 7250 "Knights Landing" | NVidia Tesla P100 "Pascal" |
|---|---|---|---|
| Cores@Clock | 2 x 18 @ ≥2.3 GHz | 68 @ 1.4 GHz | 56 SMs @ ~1.3 GHz |
| SP Performance/core | ≥73.6 GFlop/s | 89.6 GFlop/s | ~166 GFlop/s |
| Threads@STREAM | ~12 | ~60 | >25000 |
| SP peak | ≥2.6 TFlop/s | 6.1 TFlop/s | ~9.3 TFlop/s |
| Stream BW (meas.) | 2 x 62.5 GB/s | 450 GB/s (HBM) | 510 GB/s |
| Transistors / TDP | ~2x7 Billion / 2x145 W | 8 Billion / 215W | 14 Billion/300W |

# What is CUDA?

- CUDA Architecture
  - Expose GPU parallelism for general-purpose computing
  - Boost performance

- CUDA C/C++
  - Based on industry-standard C/C++
  - Small set of extensions to enable parallel programming
  - Straightforward APIs to manage devices, memory etc.

- This session introduces CUDA C

Note: this lecture is adapted from the NVIDIA training course

Clarkson
UNIVERSITY
*defy*.convention

# Introduction to CUDA C

- What will you learn in this session?
  - Start from "Hello World!"
  - Write and launch CUDA C kernels
  - Manage GPU memory
  - Manage communication and synchronization

# Part I: Heterogenous Computing

# HELLO WORLD!

**CONCEPTS**

- Heterogeneous Computing
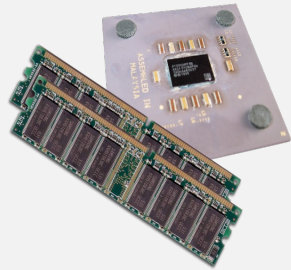- Blocks
- Threads
- Indexing
- Shared memory
- __syncthreads()
- Asynchronous operation
- Handling errors
- Managing devices

Clarkson
UNIVERSITY
*defy.*convention

# Heterogeneous Computing

- Terminology:
  - *Host* — The CPU and its memory (host memory)
  - *Device* — The GPU and its memory (device memory)



Host



Device

# Heterogeneous Computing

# Simple Processing Flow



1. Copy input data from CPU memory to GPU memory

# Simple Processing Flow



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute

Labels in figure: CPU, Bridge, CPU Memory, PCI Bus, GigaThread™, Interconnect, L2, DRAM

# Simple Processing Flow



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute
3. Copy results from GPU memory to CPU memory

# Hello World!

```c
int main(void) {
        printf("Hello World!\n");
        return 0;
}
```

Output:

```
$ nvcc
hello_world.cu
$ ./a.out
Hello World!
$
```

- Standard C that runs on the host

- NVIDIA compiler (nvcc) can be used to compile programs with no *device* code

# Hello World! with Device Code

```
__global__ void mykernel(void) {
}

int main(void) {
    mykernel<<<1,1>>>();
    printf("Hello World!\n");
    return 0;
}
```

- Two new syntactic elements…

# Hello World! with Device Code

```
__global__ void mykernel(void) {
}
```

- CUDA C/C++ keyword `__global__` indicates a function that:
  - Runs on the device
  - Is called from host code

- `nvcc` separates source code into host and device components
  - Device functions (e.g. `mykernel()`) processed by NVIDIA compiler
  - Host functions (e.g. `main()`) processed by standard host compiler
    - **e.g., gcc**

# Hello World! with Device Code

```
mykernel<<<1,1>>>();
```

- Triple angle brackets mark a call from *host* code to *device* code
    - Also called a "kernel launch"
    - We'll return to the parameters (1,1) in a moment

- That's all that is required to execute a function on the GPU!

# Hello World! with Device Code

```
__global__ void mykernel(void){

}


int main(void) {
        mykernel<<<1,1>>>();
        printf("Hello World!\n");
        return 0;

}
```

Output:

```
$ nvcc
hello.cu
$ ./a.out
Hello World!
$
```

- **mykernel()** does nothing, somewhat anticlimactic!

# Parallel Programming in CUDA C

- But wait... GPU computing is about massive parallelism!

- We need a more interesting example...

- We'll start by adding two integers and build up to vector addition



a    b         c

# Addition on the Device

- A simple kernel to add two integers

```
__global__ void add(int *a, int *b, int *c) {
        *c = *a + *b;
}
```

- As before __global__ is a CUDA C keyword meaning
  - add() will execute on the device
  - add() will be called from the host

# Addition on the Device

- Note that we use pointers for the variables

```
__global__ void add(int *a, int *b, int *c) {
        *c = *a + *b;
}
```

- `add()` runs on the device, so `a`, `b` and `c` must point to device memory

- We need to allocate memory on the GPU

Clarkson
UNIVERSITY
*defy* convention

# Memory Management

- Host and device memory are separate entities
  - *Device* pointers point to GPU memory
  - *Host* pointers point to CPU memory

- Simple CUDA API for handling device memory
  - `cudaMalloc()`, `cudaFree()`, `cudaMemcpy()`
  - Similar to the C equivalents `malloc()`, `free()`, `memcpy()`

# **Addition on the Device:** `add()`

- Returning to our `add()` kernel

```
__global__ void add(int *a, int *b, int *c) {
        *c = *a + *b;
}
```

- Let's take a look at main()…

# Addition on the Device: `main()`

```
int main(void) {
        int a, b, c;                        // host copies of a, b, c
        int *d_a, *d_b, *d_c;               // device copies of a, b, c
        int size = sizeof(int);

        // Allocate space for device copies of a, b, c
        cudaMalloc((void **)&d_a, size);
        cudaMalloc((void **)&d_b, size);
        cudaMalloc((void **)&d_c, size);

        // Setup input values
        a = 2;
        b = 7;
```

# Addition on the Device: `main()`

```
    // Copy inputs to device
    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

    // Launch add() kernel on GPU
    add<<<1,1>>>(d_a, d_b, d_c);

    // Copy result back to host
    cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

    // Cleanup
    cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
    return 0;
}
```

# Part II: Blocks

# Moving to Parallel

- GPU computing is about massive parallelism
  - So how do we run code in parallel on the device?

```
add<<< 1, 1 >>>();

add<<< N, 1 >>>();
```

- Instead of executing `add()` once, execute N times in parallel

# Vector Addition on the Device

- With `add()` running in parallel we can do vector addition

- Terminology: each parallel invocation of `add()` is referred to as a block
  - Each invocation can refer to its block index using `blockIdx.x`

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

- By using `blockIdx.x` to index into the array, each block handles a different index

# Vector Addition on the Device

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

- On the device, each block can execute in parallel:

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];
```

# Vector Addition on the Device: `add()`

- Returning to our parallelized `add()` kernel

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

- Let's take a look at main()…

# Vector Addition on the Device: main()

```c
#define N 512
int main(void) {
    int *a, *b, *c;          // host copies of a, b, c
    int *d_a, *d_b, *d_c;    // device copies of a, b, c
    int size = N * sizeof(int);

    // Alloc space for device copies of a, b, c
    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    // Alloc space for host copies of a, b, c and setup input values
    a = (int *)malloc(size);
    b = (int *)malloc(size);
    c = (int *)malloc(size);
```

# Vector Addition on the Device:

## main()

```
// Copy inputs to device
cudaMemcpy(d_a, a, size, cudaMemcpyHostToDevice);
cudaMemcpy(d_b, b, size, cudaMemcpyHostToDevice);

// Launch add() kernel on GPU with N blocks
add<<<N,1>>>(d_a, d_b, d_c);

// Copy result back to host
cudaMemcpy(c, d_c, size, cudaMemcpyDeviceToHost);

// Cleanup
free(a); free(b); free(c);
cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
return 0;
}
```

# Review (1 of 2)

- Difference between *host* and *device*
  - *Host*    CPU
  - *Device*    GPU


- Using `__global__` to declare a function as device code
  - Executes on the device
  - Called from the host


- Passing parameters from host code to a device function

# Review (2 of 2)

- Basic device memory management
  - `cudaMalloc()`
  - `cudaMemcpy()`
  - `cudaFree()`

- Launching parallel kernels
  - Launch `N` copies of `add()` with `add<<<N,1>>>(…);`
  - Use `blockIdx.x` to access block index

# Part III: Threads

# CUDA Threads

- Terminology: a block can be split into parallel threads

- Let's change `add()` to use parallel *threads* instead of parallel *blocks*

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

- We use `threadIdx.x` instead of `blockIdx.x`

- Need to make one change in `main()`...

# Vector Addition Using Threads:
## main()

```
#define N 512
int main(void) {
    int *a, *b, *c;                          // host copies of a, b, c
    int *d_a, *d_b, *d_c;          // device copies of a, b, c
    int size = N * sizeof(int);

    // Alloc space for device copies of a, b, c
    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    // Alloc space for host copies of a, b, c and setup input values
    a = (int *)malloc(size);
    b = (int *)malloc(size);
    c = (int *)malloc(size);
```

# Vector Addition Using Threads:

## `main()`

```
  // Copy inputs to device
cudaMemcpy(d_a, a, size, cudaMemcpyHostToDevice);
cudaMemcpy(d_b, b, size, cudaMemcpyHostToDevice);

// Launch add() kernel on GPU with N threads
add<<<1,N>>>(d_a, d_b, d_c);

// Copy result back to host
cudaMemcpy(c, d_c, size, cudaMemcpyDeviceToHost);

// Cleanup
free(a); free(b); free(c);
cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
return 0;
}
```
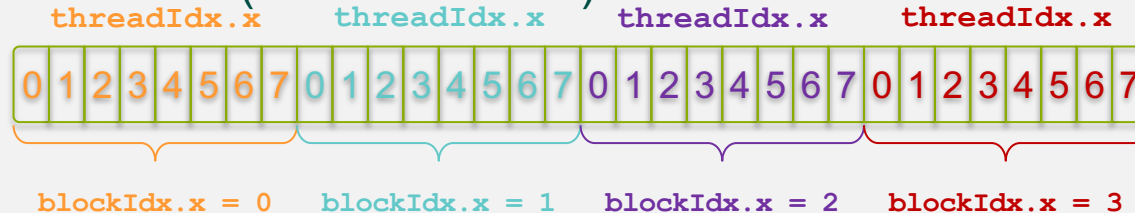
# Part IV: Indexing

# Combining Blocks and Threads

- We've seen parallel vector addition using:
  - Many blocks with one thread each
  - One block with many threads

- Let's adapt vector addition to use both blocks and threads

- Why? We'll come to that…

- First let's discuss data indexing…

# Indexing Arrays with Blocks and Threads

- No longer as simple as using `blockIdx.x` and `threadIdx.x`

  - Consider indexing an array with one element per thread (8 threads/block)

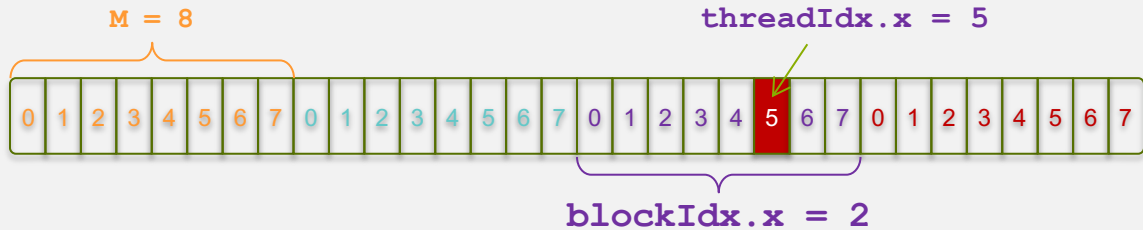| `threadIdx.x` | `threadIdx.x` | `threadIdx.x` | `threadIdx.x` |
|---|---|---|---|
| 0 1 2 3 4 5 6 7 | 0 1 2 3 4 5 6 7 | 0 1 2 3 4 5 6 7 | 0 1 2 3 4 5 6 7 |
| `blockIdx.x = 0` | `blockIdx.x = 1` | `blockIdx.x = 2` | `blockIdx.x = 3` |

- With M threads/block a unique index for each thread is given by:

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

# Indexing Arrays: Example

- Which thread will operate on the red element?



M = 8

threadIdx.x = 5

blockIdx.x = 2

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

# Vector Addition with Blocks and Threads

- Use the built-in variable `blockDim.x` for threads per block

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

- Combined version of `add()` to use parallel threads *and* parallel blocks

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

- What changes need to be made in `main()`?

# Addition with Blocks and Threads: `main()`

```c
#define N (2048*2048)
#define THREADS_PER_BLOCK 512
int main(void) {
    int *a, *b, *c;                    // host copies of a, b, c
    int *d_a, *d_b, *d_c;         // device copies of a, b, c
    int size = N * sizeof(int);

    // Alloc space for device copies of a, b, c
    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    // Alloc space for host copies of a, b, c and setup input values
    a = (int *)malloc(size);
    b = (int *)malloc(size);
    c = (int *)malloc(size);
```

# Addition with Blocks and Threads: `main()`

```
// Copy inputs to device
cudaMemcpy(d_a, a, size, cudaMemcpyHostToDevice);
cudaMemcpy(d_b, b, size, cudaMemcpyHostToDevice);

// Launch add() kernel on GPU
add<<<N/THREADS_PER_BLOCK,THREADS_PER_BLOCK>>>(d_a, d_b, d_c);

// Copy result back to host
cudaMemcpy(c, d_c, size, cudaMemcpyDeviceToHost);

// Cleanup
free(a); free(b); free(c);
cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
return 0;
}
```