Heterogeneous
Parallel
Programming

Lesson 1.4

# Introduction to CUDA

– Data Parallelism and Threads

Wen-mei Hwu · University of Illinois at Urbana-Champaign
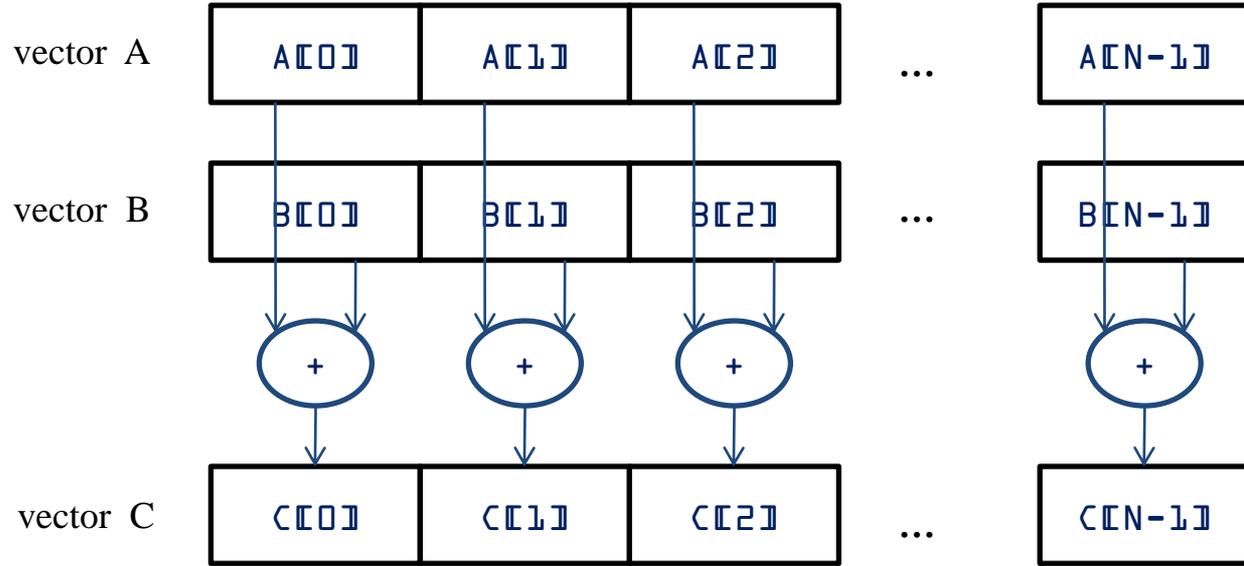
- To learn about data parallelism and the basic features of CUDA C, a heterogeneous parallel programming interface that enables exploitation of data parallelism
  - Hierarchical thread organization
  - Main interfaces for launching parallel execution
  - Thread index to data index mapping

# Data Parallelism - Vector Addition Example

| vector A | A[0] | A[1] | A[2] | ... | A[N-1] |
|---|---|---|---|---|---|

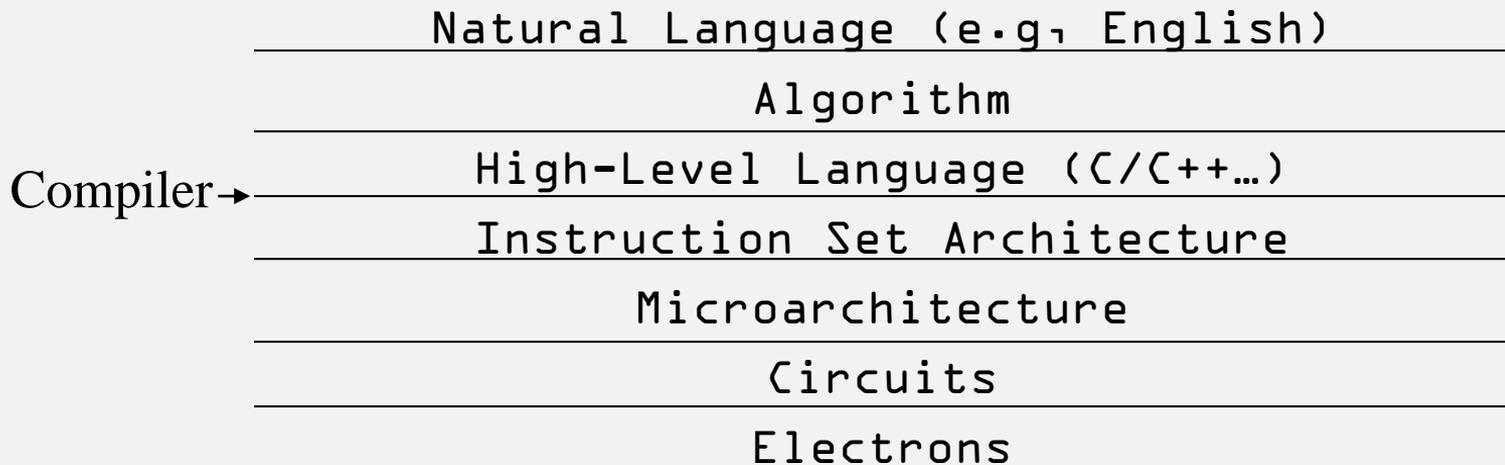| vector B | B[0] | B[1] | B[2] | ... | B[N-1] |
|---|---|---|---|---|---|

+    +    +    +

| vector C | C[0] | C[1] | C[2] | ... | C[N-1] |
|---|---|---|---|---|---|

# CUDA /OpenCL - Execution Model

- Heterogeneous host+device application C program
  - Serial parts in **host** C code
  - 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);

# From Natural Language to Electrons

| |
|---|
| Natural Language (e.g, English) |
| Algorithm |
| High-Level Language (C/C++…) |
| Instruction Set Architecture |
| Microarchitecture |
| Circuits |
| Electrons |

Compiler →

©Yale Patt and Sanjay Patel, *From bits and bytes to gates and beyond*

- An Instruction Set Architecture (ISA) is a contract between the hardware and the software.

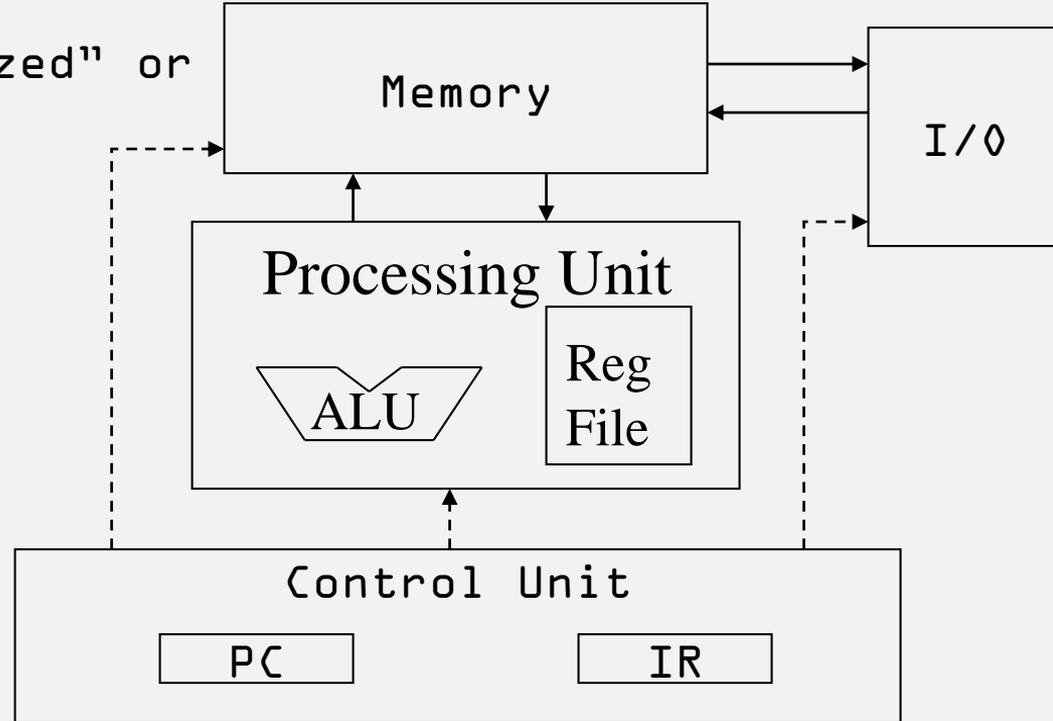- As the name suggests, it is a set of instructions that the architecture (hardware) can execute.

# A program at the ISA level

- A program is a set of instructions stored in memory that can be read, interpreted, and executed by the hardware.

- Program instructions operate on data stored in memory or provided by Input/Output (I/O) device.
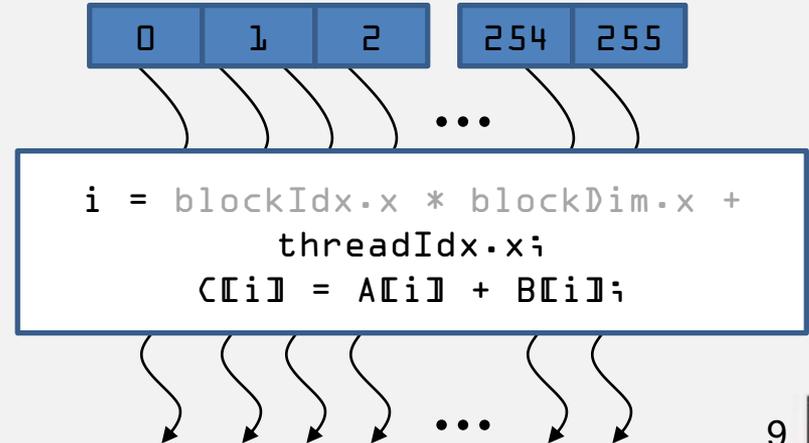
A thread is a "virtualized" or "abstracted"
Von-Neumann Processor

**Memory**

**I/O**

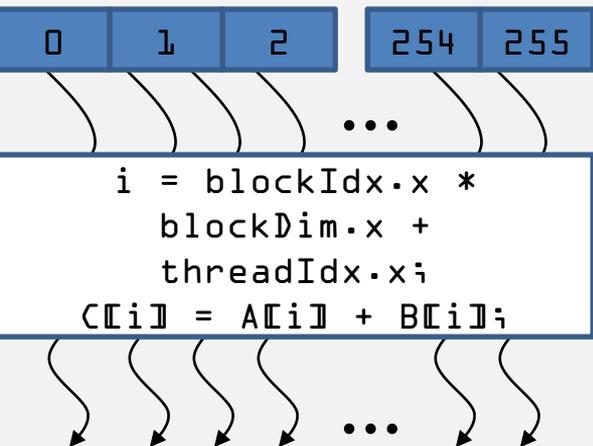**Processing Unit**

ALU

Reg File

**Control Unit**

PC

IR

# Arrays of Parallel Threads

- A CUDA kernel is executed by a grid (array) of threads
    - All threads in a grid run the same kernel code (SPMD)
    - Each thread has indexes that it uses to compute memory addresses and make control decisions
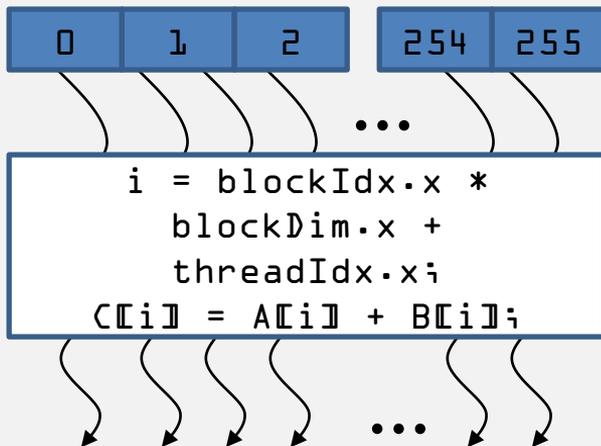
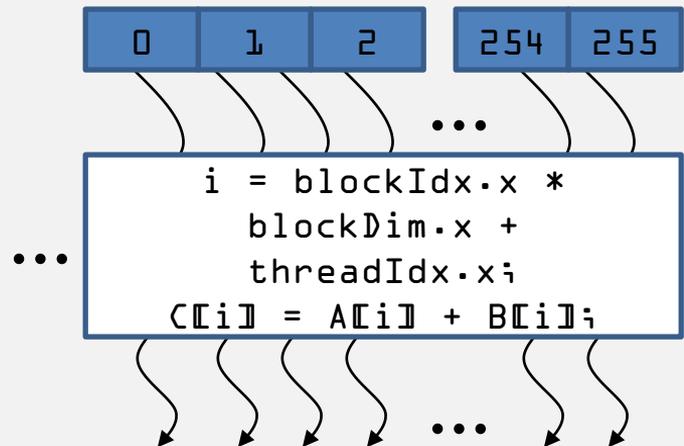| 0 | 1 | 2 | 254 | 255 |
|---|---|---|-----|-----|

• • •

```
i = blockIdx.x * blockDim.x +
          threadIdx.x;
      C[i] = A[i] + B[i];
```

• • •

# Thread Blocks: Scalable Cooperation

## Thread Block 0

| 0 | 1 | 2 | | 254 | 255 |

...

```
i = blockIdx.x *
   blockDim.x +
   threadIdx.x;
C[i] = A[i] + B[i];
```

...

## Thread Block 1

| 0 | 1 | 2 | | 254 | 255 |

...

```
i = blockIdx.x *
   blockDim.x +
   threadIdx.x;
C[i] = A[i] + B[i];
```

...

...

## Thread Block N-1

| 0 | 1 | 2 | | 254 | 255 |

...

```
i = blockIdx.x *
   blockDim.x +
   threadIdx.x;
C[i] = A[i] + B[i];
```
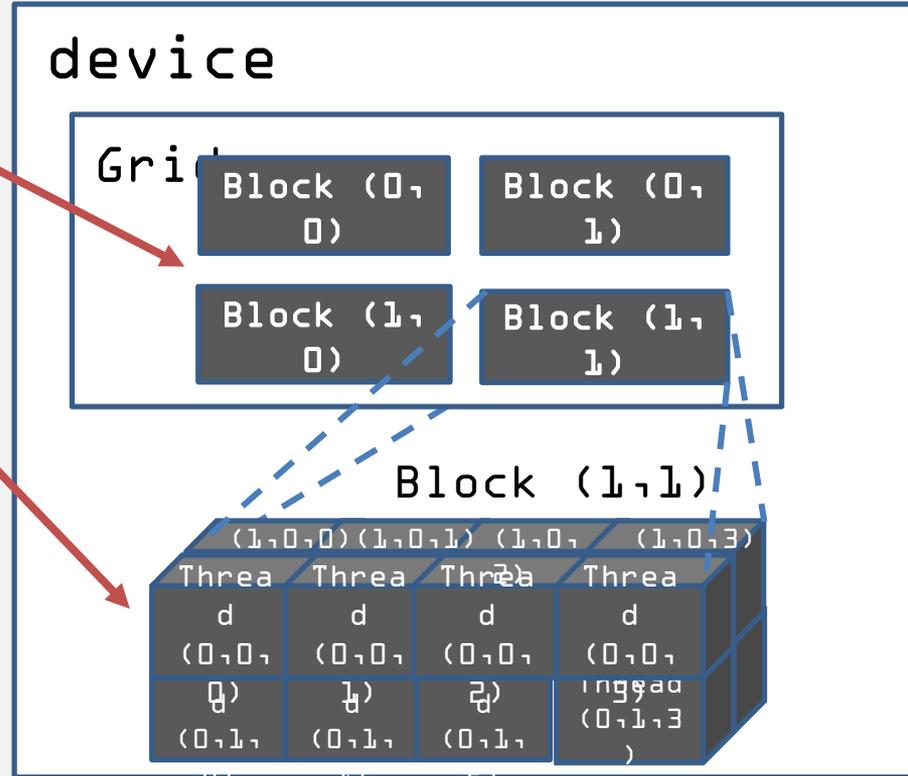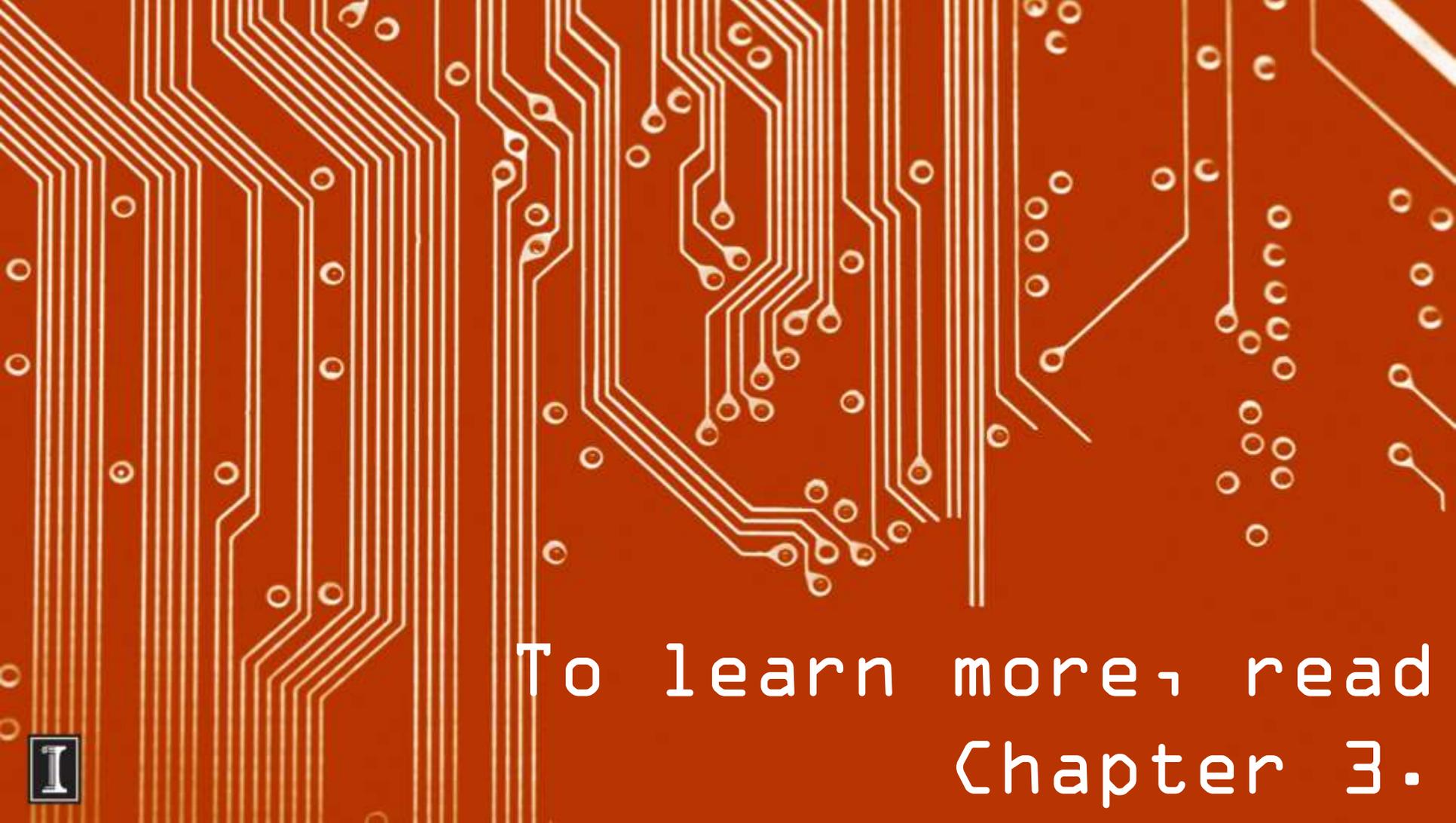
...

- Divide thread array into multiple blocks
  - Threads within a block cooperate via shared memory, atomic operations and barrier synchronization
  - Threads in different blocks do not interact

- Each thread uses indices to decide what data to work on
  - blockIdx: 1D, 2D, or 3D (CUDA 4.0)
  - threadIdx: 1D, 2D, or 3D

- Simplifies memory addressing when processing multidimensional data
  - Image processing
  - Solving PDEs on volumes
  - …

device

Grid

Block (0, 0)  Block (0, 1)

Block (1, 0)  Block (1, 1)

Block (1,1)

(1,0,0) (1,0,1) (1,0, (1,0,3)

Thread (0,0, 0)  Thread (0,0, 1)  Thread (0,0, 2)  Thread (0,0,

Thread (0,1,  Thread (0,1,  Thread (0,1,  Thread (0,1,3 )

11

To learn more, read Chapter 3.