**Performance Engineering of Software Systems**

SPEED LIMIT

∞

PER ORDER OF 6.106
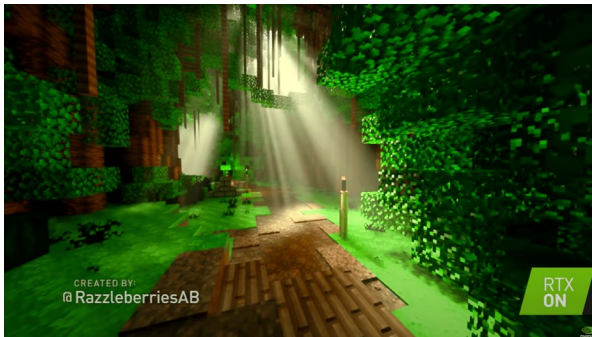
LECTURE 19
**GPU PROGRAMMING**

**Xuhao Chen**

**September 20, 2024**

# What is a GPU?

- Graphics Processing Units



ray tracing

gaming

3D rendering

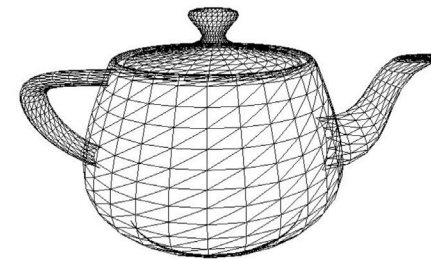Image credit: Henrik Wann Jensen

# Why GPU?

- CPU
  - ~10s cores
  - Low Latency
  - Good for serial processing
  - Good for interactive tasks
  - Task parallelism

- GPU
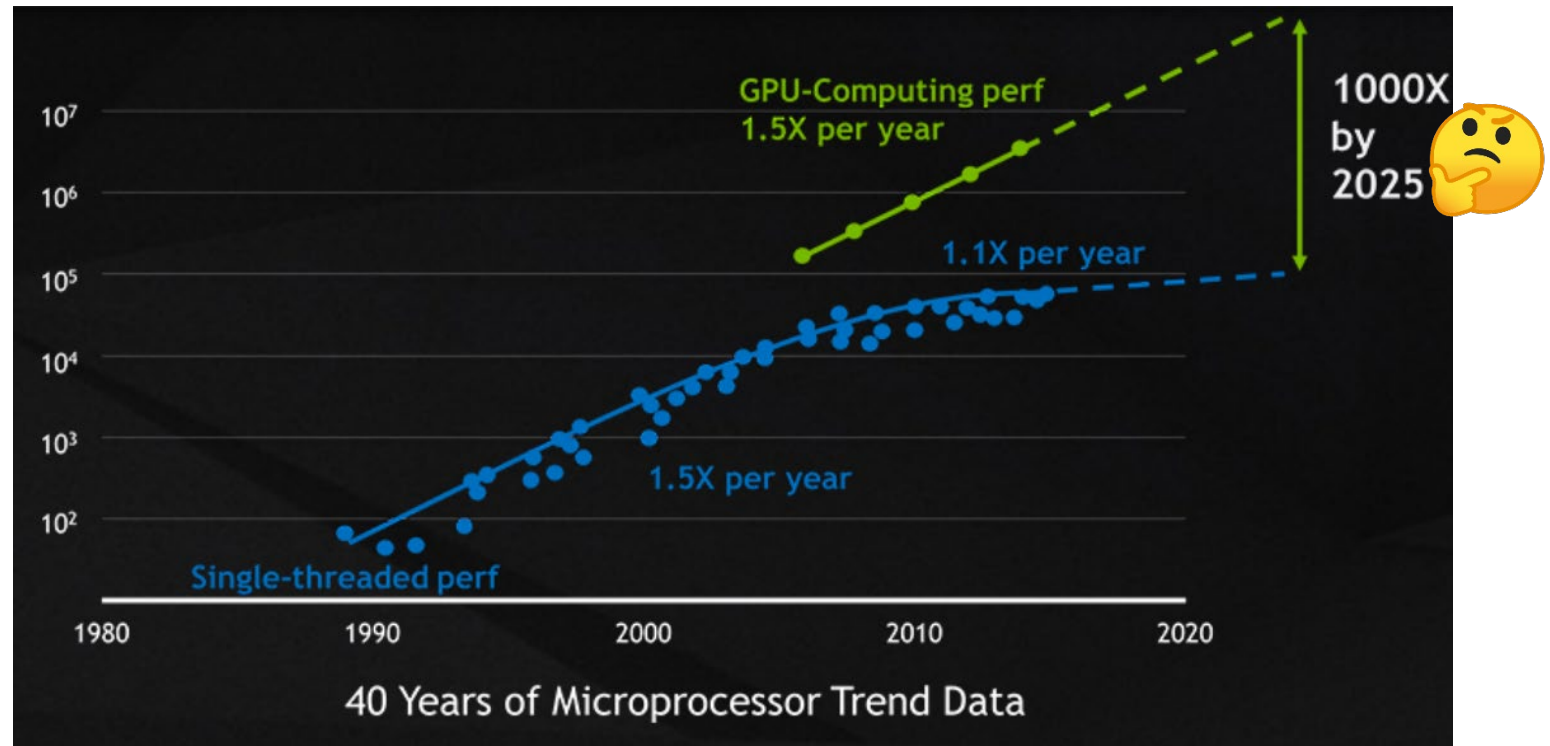  - 100s ~ 1000s cores
  - High throughput
  - Good for parallel processing
  - Good for big-data tasks
  - Data parallelism

# Why GPU?

| | Throughput | Power | Throughput/Power |
|---|---|---|---|
| Intel Skylake | 128 SP GFLOPS/4 Cores | 100+ Watts | ~1 GFLOPS/Watt |
| NVIDIA V100 | 15 TFLOPS | 200+ Watts | ~75 GFLOPS/Watt |

Also,



40 Years of Microprocessor Trend Data

# Compute Intensive Applications
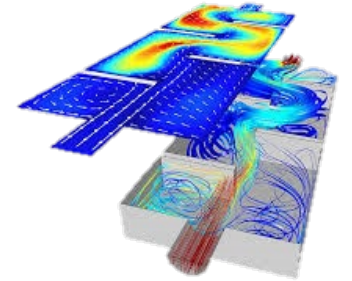
- Bioinformatics
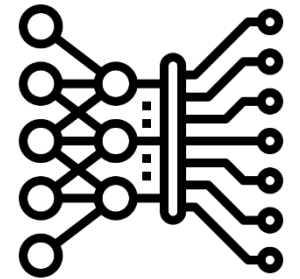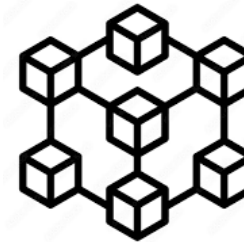
- Computational Chemistry

- Computational Finance

- Computational Fluid Dynamics

- AI & Machine Learning

- Block Chain

- Data Science, Medical Imaging, Imaging & Computer Vision, Weather and Climate, ⋯

# GPU Architecture

## SIMT: single-instruction multiple threads

# 3 Ways of GPU Acceleration

**Applications**

| GPU-accelerated libraries | OpenACC Directives | Programming Languages |
|---|---|---|
| Seamless linking to GPU-enabled libraries. | Simple directives for easy GPU-acceleration of new and existing applications | Most powerful and flexible way to design GPU accelerated applications |
| cuFFT, cuBLAS, Thrust, NPP, IMSL, CULA, cuRAND, etc. | PGI Accelerator | C/C++, Fortran, Python, Java, etc. |

# 3 Ways of GPU Acceleration

Applications

| Libraries | OpenACC Directives | Programming Languages |
|---|---|---|
| "Drop-in" Acceleration | Easily Accelerate Applications | Maximum Flexibility |

# GPU Accelerated Libraries



NVIDIA cuBLAS

NVIDIA cuRAND

NVIDIA NPP

OpenCV

NVIDIA cuSPARSE

CUSP
Sparse Linear Algebra

Thrust
C++ STL Features for CUDA

NVIDIA cuFFT

# Thrust: Rapid Parallel C++ Development
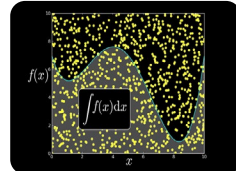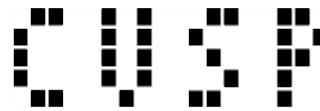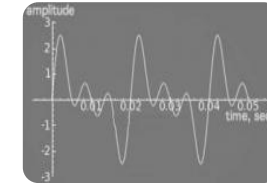
- Resembles C++ STL

- High-level interface
  - Enhances developer productivity
  - Enables performance portability between GPUs and multicore CPUs

- Flexible
  - CUDA, OpenMP, and TBB backends
  - Extensible and customizable
  - Integrates with existing software

- Open source

```cpp
// generate 32M random numbers on host
thrust::host_vector<int> h_vec(32 << 20);
thrust::generate(h_vec.begin(),
                 h_vec.end(),
                 rand);

// transfer data to device (GPU)
thrust::device_vector<int> d_vec = h_vec;

// sort data on device
thrust::sort(d_vec.begin(), d_vec.end());

// transfer data back to host
thrust::copy(d_vec.begin(),
             d_vec.end(),
             h_vec.begin());
```

http://developer.nvidia.com/thrust   or   http://thrust.googlecode.com

# Libraries: Easy, High-Quality Acceleration

- Ease of use:   Using libraries enables GPU acceleration without in-depth knowledge of GPU programming

- "Drop-in":   Many GPU-accelerated libraries follow standard APIs, thus enabling acceleration with minimal code changes

- Quality:   Libraries offer high-quality implementations of functions encountered in a broad range of applications

- Performance:   NVIDIA libraries are tuned by experts

# 3 Ways of GPU Acceleration

Applications

Libraries

OpenACC
Directives

Programming
Languages

"Drop-in"
Acceleration

Easily Accelerate
Applications

Maximum
Flexibility

# OpenACC Directives

CPU            GPU

```
Program myscience
   ... serial code ...
!$acc kernels
  do k = 1,n1
    do i = 1,n2
      ... parallel code ...
    enddo
  enddo
!$acc end kernels
   ...
End Program myscience
```

**Your original**
**Fortran or C code**

OpenACC
compiler
Hint

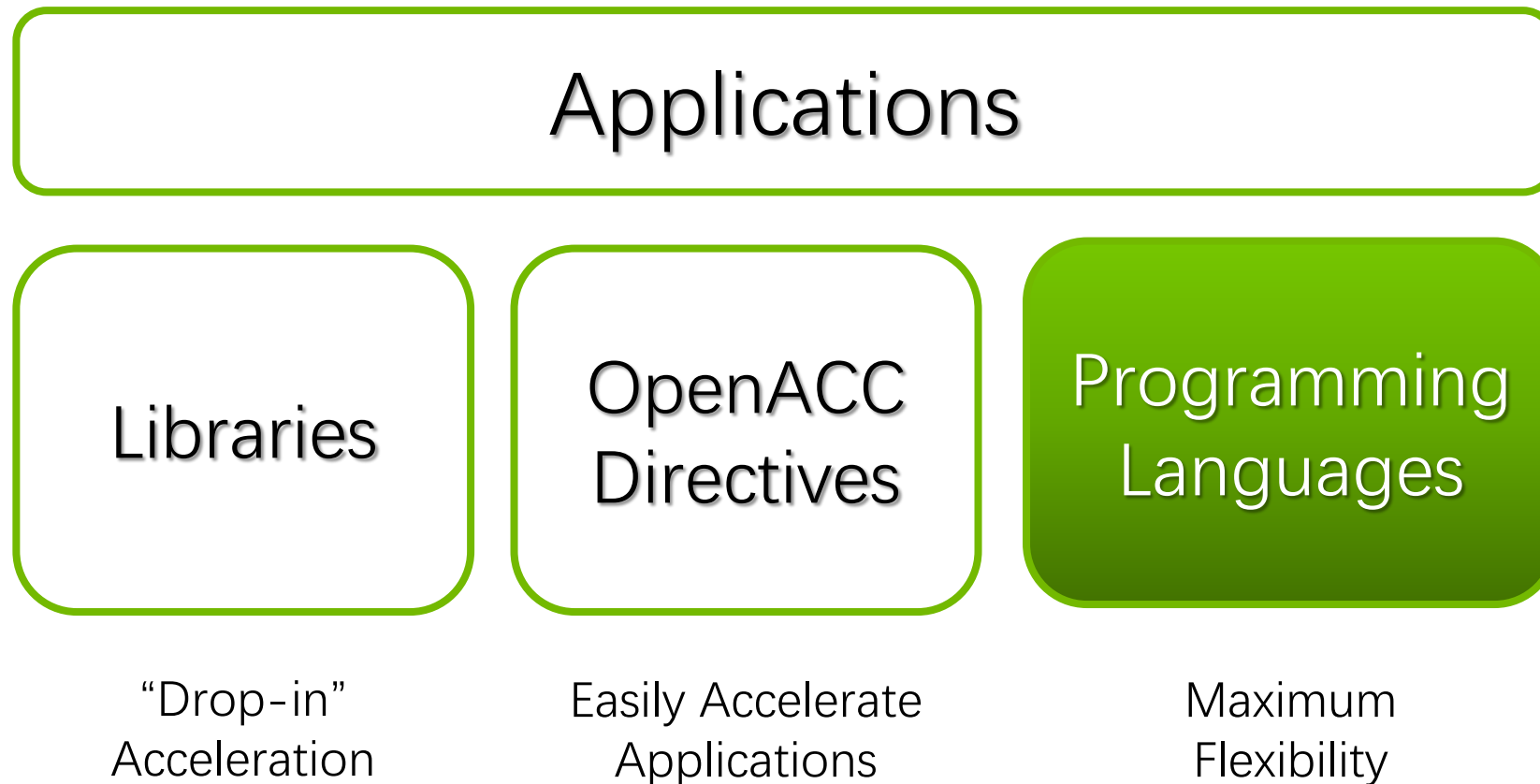- Simple Compiler hints
- Compiler Parallelizes code
- Works on many-core GPUs & multicore CPUs

Easy      Open      Powerful

# 3 Ways of GPU Acceleration



Applications

| Libraries | OpenACC Directives | Programming Languages |

"Drop-in" Acceleration

Easily Accelerate Applications

Maximum Flexibility

# GPU Programming Languages

| | |
|---|---|
| C | OpenACC, CUDA C |
| C++ | Thrust, CUDA C++ |
| Fortran | OpenACC, CUDA Fortran |
| Python | PyCUDA, PyOpenCL, Numba |
| Numerical analytics | MATLAB, Mathematica, LabVIEW |
| Machine Learning | Theano, Tensorflow, Caffe, Torch, etc. |

**CONCEPTS**

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

# PROGRAM A GPU WITH CUDA

# Heterogeneous Computing

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

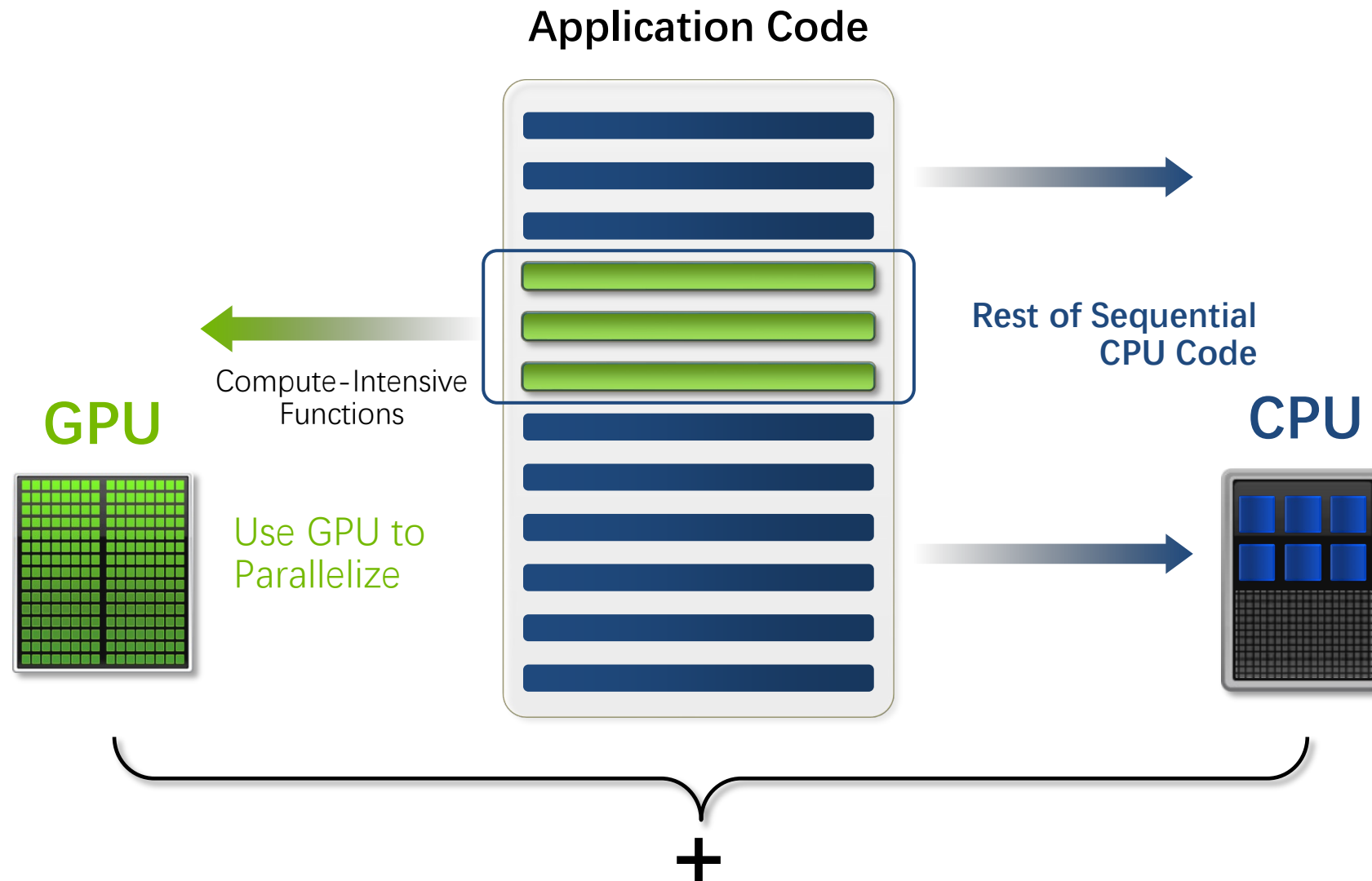**Host:** the CPU and its memory

**Device:** the GPU and its memory

# CPU-GPU Heterogeneous Computing

**Application Code**

**GPU**

Compute-Intensive
Functions
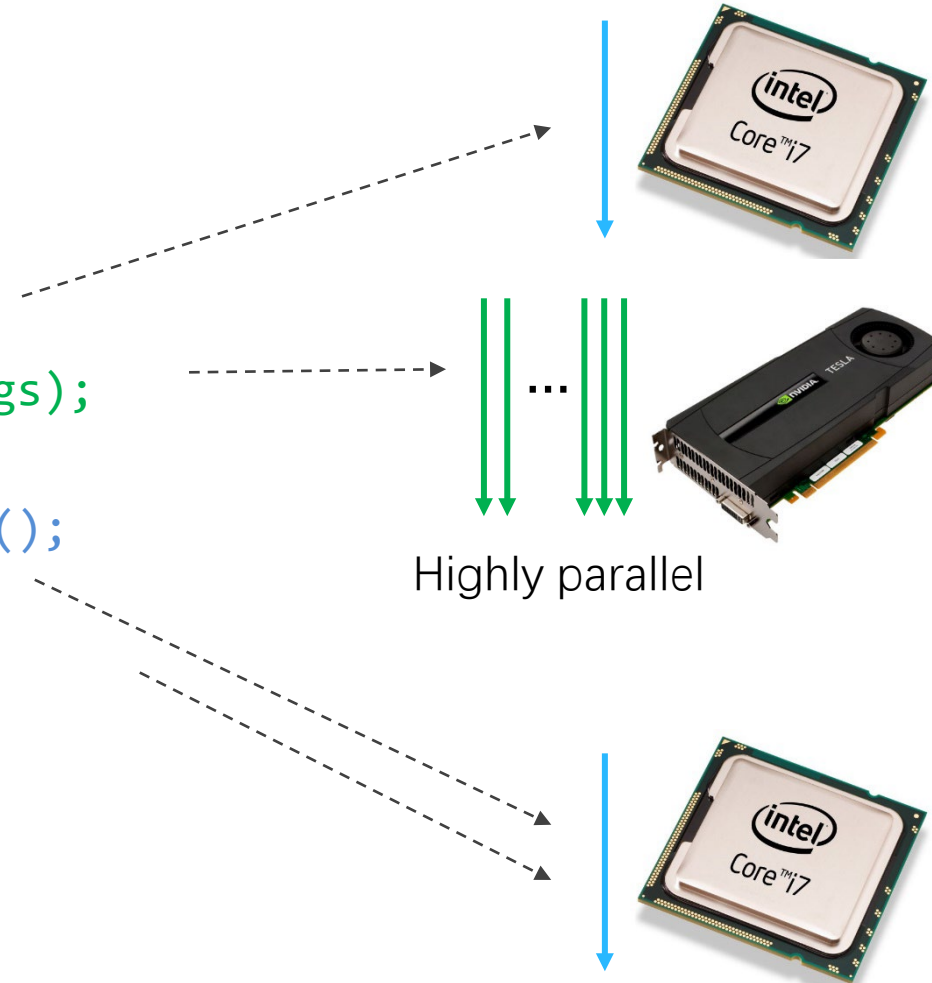
Use GPU to
Parallelize

**Rest of Sequential
CPU Code**

**CPU**

+

# Heterogeneous Computing with CUDA

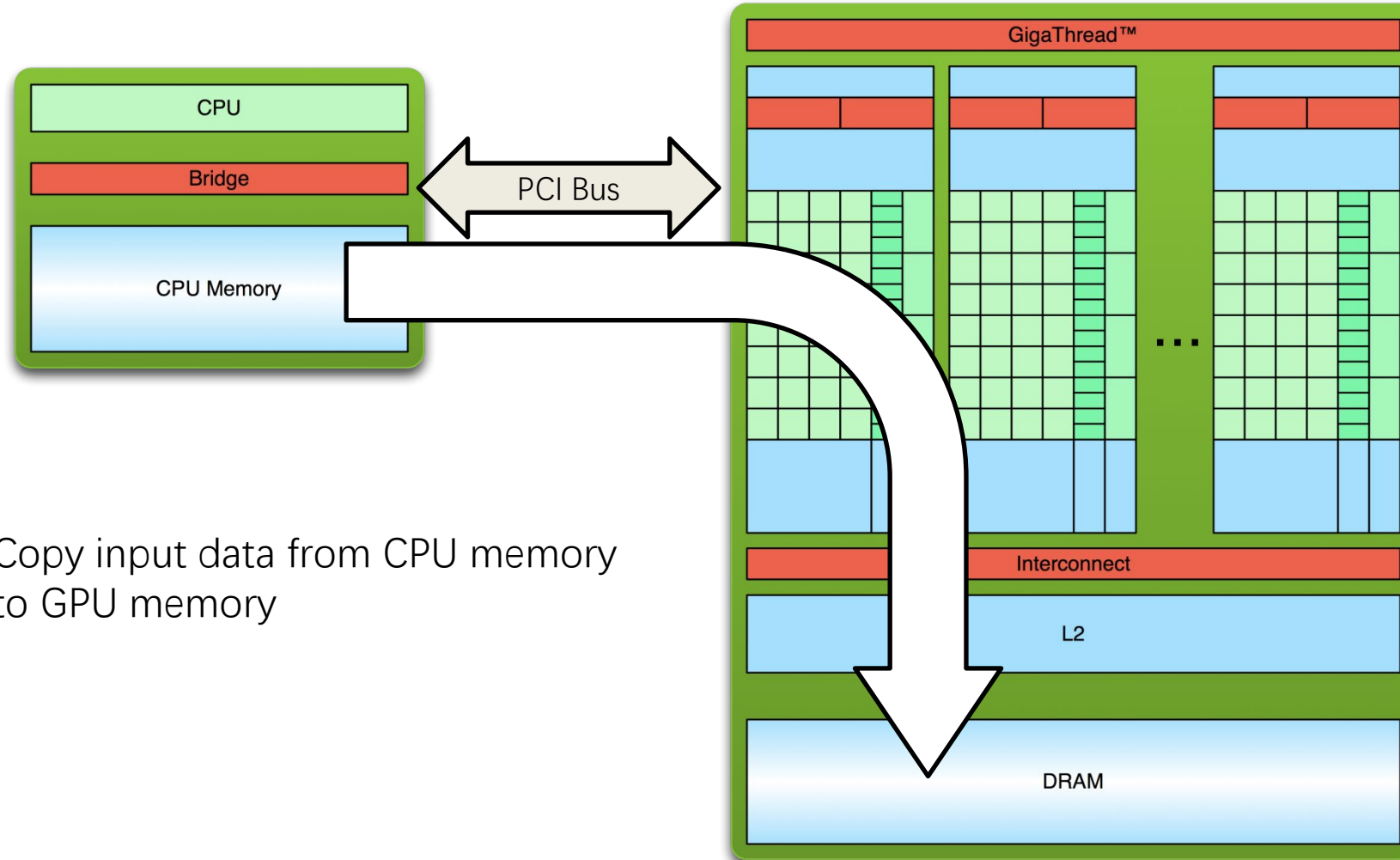- CUDA Compute Unified Device Architecture

```
do_something_on_host();
kernel<<<nBlk, nTid>>>(args);
cudaDeviceSynchronize();
do_something_else_on_host();
```

Highly parallel

# Simple Processing Flow
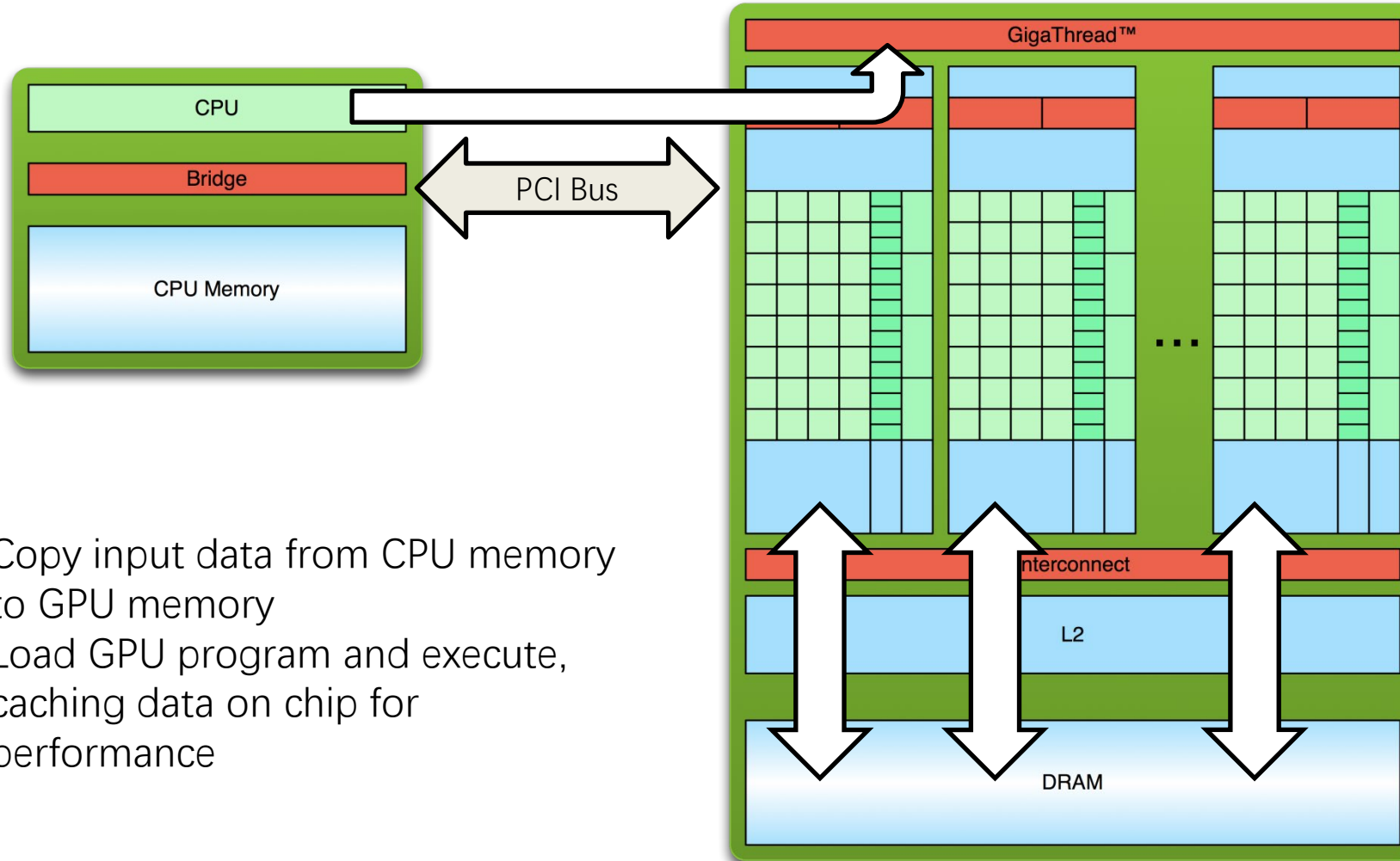


1. Copy input data from CPU memory to GPU memory

# Simple Processing Flow



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance

# Simple Processing Flow



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance
3. Copy results from GPU memory to CPU memory

# Heterogeneous Computing with CUDA C

- Let's start with simply adding two integers

**Vector Addition**

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

- Here **__global__** is a CUDA C/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

# Memory Management

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

    May be passed to/from host code

    May *not* be dereferenced in host code
  - *Host* pointers point to CPU memory

    May be passed to/from device code

    May *not* be dereferenced in device code

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

# Addition on the Device: `main()`

```c
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;
```

```
// 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;
}
```

# Moving to Parallel Execution

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

# Thread Batching: Grids and Blocks

- A kernel is executed as a **grid** of **thread blocks**
  - All threads within a thread block share a portion of data memory
  - Threads/blocks have 1D/2D/3D IDs

- A **thread block** is a batch of threads that can **cooperate** with each other by:
  - Synchronizing their execution
    - For hazard-free common memory accesses
  - Efficiently sharing data through a low latency **shared memory**

- Two threads from two different thread blocks cannot directly cooperate

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

- Each parallel invocation of **add()** is referred to as a block
  - The set of blocks is referred to as a grid
  - 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];
```

- 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()···

```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); random_ints(a, N);
    b = (int *)malloc(size); random_ints(b, N);
    c = (int *)malloc(size);
```

```
// 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;
}
```

- 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

- 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

# INTRODUCING THREADS

**CONCEPTS**

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

# 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()`...

```
#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); random_ints(a, N);
        b = (int *)malloc(size); random_ints(b, N);
        c = (int *)malloc(size);
```

# Vector Addition Using Threads: `main()`

```c
// 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;
}
```
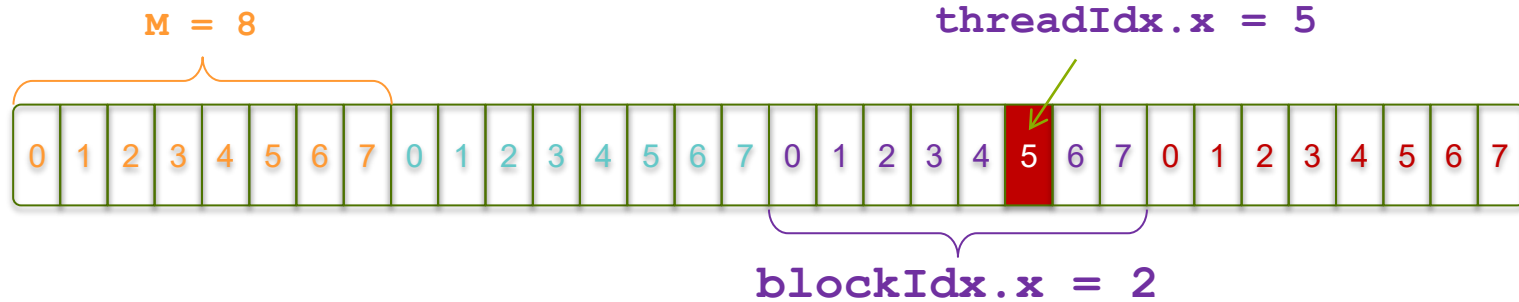
# Indexing Arrays with Blocks and Threads

- No longer as simple as using `blockIdx.x` and `threadIdx.x`
  - Consider indexing an array with one element per thread (8 threads/block)



- With M threads/block a unique index for each thread is given by
  `int index = threadIdx.x + blockIdx.x * M;`

# 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()**?

```
#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); random_ints(a, N);
        b = (int *)malloc(size); random_ints(b, N);
        c = (int *)malloc(size);
```

```
// 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;
    }
```

# Handling Arbitrary Vector Sizes

- Typical problems are not friendly multiples of `blockDim.x`

- Avoid accessing beyond the end of the arrays:

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

- Update the kernel launch:

```
add<<<(N + M-1) / M,M>>>(d_a, d_b, d_c, N);
```

# Why Bother with Threads?

- Threads seem unnecessary
  - They add a level of complexity
  - What do we gain?

- Unlike parallel blocks, threads have mechanisms to:
  - Communicate
  - Synchronize

- To look closer, we need a new example⋯

# COOPERATING THREADS

**CONCEPTS**

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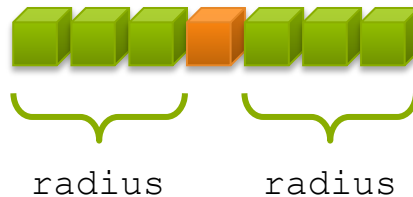- Shared memory
- __syncthreads()
- Asynchronous operation
- Handling errors
- Managing devices

# 1D Stencil

- Consider applying a 1D stencil to a 1D array of elements
  - Each output element is the sum of input elements within a radius

- If radius is 3, then each output element is the sum of 7 input elements:



radius          radius

# Implementing Within a Block

- Each thread processes one output element
  - blockDim.x elements per block


- Input elements are read several times
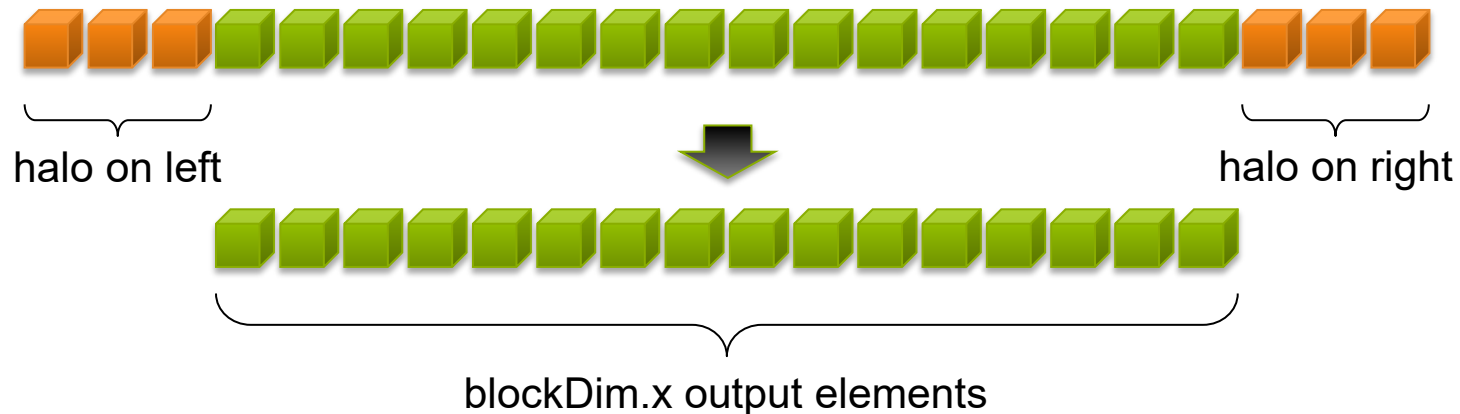  - With radius 3, each input element is read seven times

# Sharing Data Between Threads

- Terminology: within a block, threads share data via shared memory

- Extremely fast on-chip memory, user-managed

- Declare using `__shared__`, allocated per block

- Data is not visible to threads in other blocks

# Implementing With Shared Memory

- Cache data in shared memory
  - Read (blockDim.x + 2 * radius) input elements from global memory to shared memory
  - Compute blockDim.x output elements
  - Write blockDim.x output elements to global memory

  - Each block needs a **halo** of radius elements at each boundary



halo on left          halo on right

blockDim.x output elements

# Stencil Kernel

```
__global__ void stencil_1d(int *in, int *out) {
  __shared__ int temp[BLOCK_SIZE + 2 * RADIUS];
  int gindex = threadIdx.x + blockIdx.x * blockDim.x;
  int lindex = threadIdx.x + RADIUS;

  // Read input elements into shared memory
  temp[lindex] = in[gindex];
  if (threadIdx.x < RADIUS) {
    temp[lindex - RADIUS] = in[gindex - RADIUS];
    temp[lindex + BLOCK_SIZE] =
      in[gindex + BLOCK_SIZE];
  }
```

# Stencil Kernel

```
// Apply the stencil
int result = 0;
for (int offset = -RADIUS ; offset <= RADIUS ; offset++)
  result += temp[lindex + offset];

// Store the result
out[gindex] = result;
}
```

# Data Race!

- The stencil example will not work…

- Suppose thread 15 reads the halo before thread 0 has fetched it…

```
temp[lindex] = in[gindex];                          Store at temp[18]
if (threadIdx.x < RADIUS) {
  temp[lindex - RADIUS = in[gindex - RADIUS];          Skipped, threadIdx > RADIUS
  temp[lindex + BLOCK_SIZE] = in[gindex + BLOCK_SIZE];
}

int result = 0;
result += temp[lindex + 1];                         Load from temp[19]
```

# __syncthreads()

- `void __syncthreads();`

- Synchronizes all threads within a block
  - Used to prevent RAW / WAR / WAW hazards

- All threads must reach the barrier
  - In conditional code, the condition must be uniform across the block

# Stencil Kernel

```
__global__ void stencil_1d(int *in, int *out) {
    __shared__ int temp[BLOCK_SIZE + 2 * RADIUS];
    int gindex = threadIdx.x + blockIdx.x * blockDim.x;
    int lindex = threadIdx.x + radius;


    // Read input elements into shared memory
    temp[lindex] = in[gindex];
    if (threadIdx.x < RADIUS) {
        temp[lindex - RADIUS] = in[gindex - RADIUS];
        temp[lindex + BLOCK_SIZE] = in[gindex + BLOCK_SIZE];
    }


    // Synchronize (ensure all the data is available)
    __syncthreads();
```

# Stencil Kernel

```
// Apply the stencil
int result = 0;
for (int offset = -RADIUS ; offset <= RADIUS ; offset++)
    result += temp[lindex + offset];

// Store the result
out[gindex] = result;
}
```

- Launching parallel threads
  - Launch `N` blocks with `M` threads per block with `kernel<<<N,M>>>(…);`
  - Use `blockIdx.x` to access block index within grid
  - Use `threadIdx.x` to access thread index within block

- Allocate elements to threads:

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

- Use `__shared__` to declare a variable/array in shared memory
  - Data is shared between threads in a block
  - Not visible to threads in other blocks

- Use `__syncthreads()` as a barrier
  - Use to prevent data hazards

**CONCEPTS**

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

# MANAGING THE DEVICE

# Coordinating Host & Device

- Kernel launches are asynchronous
  - Control returns to the CPU immediately

- CPU needs to synchronize before consuming the results

| | |
|---|---|
| **cudaMemcpy()** | Blocks the CPU until the copy is complete<br>Copy begins when all preceding CUDA calls have completed |
| **cudaMemcpyAsync()** | Asynchronous, does not block the CPU |
| **cudaDeviceSynchronize()** | Blocks the CPU until all preceding CUDA calls have completed |

# Reporting Errors

- All CUDA API calls return an error code (`cudaError_t`)
  - Error in the API call itself

    OR
  - Error in an earlier asynchronous operation (e.g. kernel)

- Get the error code for the last error:

  ```
  cudaError_t cudaGetLastError(void)
  ```
- Get a string to describe the error:

  ```
  char *cudaGetErrorString(cudaError_t)
  ```

  ```
  printf("%s\n", cudaGetErrorString(cudaGetLastError()));
  ```

# Device Management

- Application can query and select GPUs

  `cudaGetDeviceCount``(int *count)`

  `cudaSetDevice``(int device)`

  `cudaGetDevice``(int *device)`

  `cudaGetDeviceProperties``(cudaDeviceProp *prop, int device)`

- Multiple threads can share a device

- A single thread can manage multiple devices

  `cudaSetDevice``(i)` to select current device

  `cudaMemcpy``(…)` for peer-to-peer copies[†]

[†] requires OS and device support

# Summary: What have we learned?

- Write and launch CUDA C/C++ kernels
  - **`__global__`, `blockIdx.x`, `threadIdx.x`, `<<<>>>`**

- Manage GPU memory
  - **`cudaMalloc(), cudaMemcpy(), cudaFree()`**

- Manage communication and synchronization
  - **`__shared__, __syncthreads()`**
  - **`cudaMemcpy()`** vs. **`cudaMemcpyAsync()`**
  - **`cudaDeviceSynchronize()`**

# Getting Started

- Download CUDA Toolkit & SDK: www.nvidia.com/getcuda

- Nsight IDE (Eclipse or Visual Studio): www.nvidia.com/nsight

- Programming Guide/Best Practices: www.docs.nvidia.com

- Questions:
  - NVIDIA Developer forums: devtalk.nvidia.com
  - Search or ask on: www.stackoverflow.com/tags/cuda

- General: www.nvidia.com/cudazone

# Learn More

- These languages are supported on all CUDA-capable GPUs.
- You might already have a CUDA-capable GPU in your laptop or desktop PC!

CUDA C/C++
http://developer.nvidia.com/cuda-toolkit

GPU.NET
http://tidepowerd.com

Thrust C++ Template Library
http://developer.nvidia.com/thrust

MATLAB
http://www.mathworks.com/discovery/
matlab-gpu.html

CUDA Fortran
http://developer.nvidia.com/cuda-toolkit

PyCUDA (Python)
http://mathema.tician.de/software/pycuda

Mathematica
http://www.wolfram.com/mathematica/new
-in-8/cuda-and-opencl-support/