Software Performance Engineering

SPEED
LIMIT

∞

PER ORDER OF SPE

LECTURE 19
GPU PROGRAMMING

Xuhao Chen

November 6, 2025

- https://accelerated-computing.github.io/

- Spring 2026!

- Advanced GPU programming!

- Let me know if you're in the waitlist

## CSE/CMSE 822:
## PARALLEL COMPUTING

**Instructor:** Xuhao Chen

**Interested in how GPU works?**

**Tired of GPU Out-of-memory?**

**Want to accelerate your training / simulation?**

Come learn about...

- Performance modeling & Optimization
- Parallel Algorithms and Theory
- Shared-memory & Multi-core Programming
- SIMD & GPU Programming 🔥🔥🔥
- Distributed Computing & MPI
- Memory-efficient Programming
- Domain-Specific Languages
- Machine Learning Systems 🔥🔥🔥

**Syllabus:**

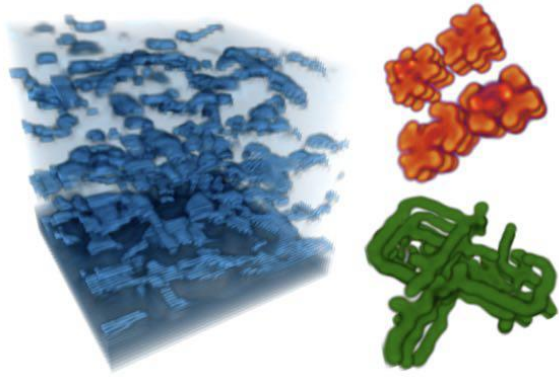https://accelerated-computing.github.io/spring26/syllabus/

**This course aims to teach you how to write fast code for parallel computers & hardware accelerators like GPUs**
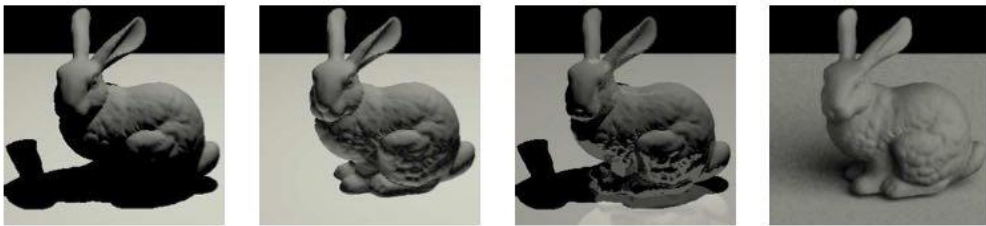
# Outline

- GPU Architecture
  - ◻ Three major ideas that make GPU processing cores run fast
  - ◻ Closer look at real GPU designs
  - ◻ The GPU memory hierarchy: moving data to processors

- **General Purpose GPU (GPGPU)**

- **GPU Programming Model**

- **Program a GPU with CUDA**
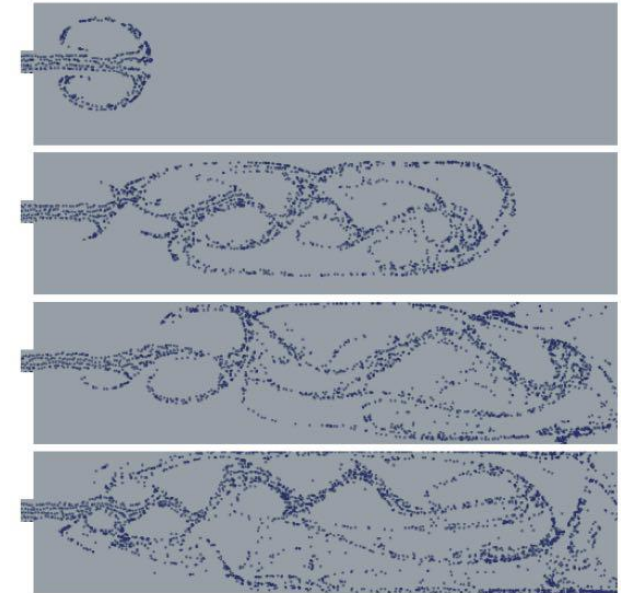
Coupled Map Lattice Simulation [Harris 02]



Ray Tracing on Programmable Graphics Hardware [Purcell 02]



Sparse Matrix Solvers [Bolz 03]

# Brook stream programming language (2004)

- Stanford graphics lab research project*

- Abstract GPU hardware as data-parallel processor

- Brook compiler translates generic stream program into graphics commands (e.g. drawTriangles) and a set of shader programs

```
kernel void scale(float amount, float a<>, out float b<>) {
        b = amount * a;
}
float scale_amount;
float input_stream<1000>; // stream declaration
float output_stream<1000>; // stream declaration
// omitting stream element initialization...
// map kernel function onto streams
scale(scale_amount, input_stream, output_stream);
```

* Buck, Ian, et al. "Brook for GPUs: stream computing on graphics hardware." *ACM transactions on graphics (TOG)* 23.3 (2004): 777-786.

# NVIDIA Tesla architecture (2007)

- First alternative, non-graphics-specific ("compute mode") interface to GPU hardware, e.g. GeForce 8800

- Let's say a user wants to run a non-graphics program on the GPU's programmable cores...
  - Application can allocate buffers in GPU memory and copy data to/from buffers
  - Application (via graphics driver) provides GPU a single kernel program binary
  - Application tells GPU to run the kernel in an SPMD fashion ("run N instances of this kernel") e.g. `launch(myKernel, N)`

- Interestingly, this is a far simpler operation than the graphics operation `drawPrimitives()`
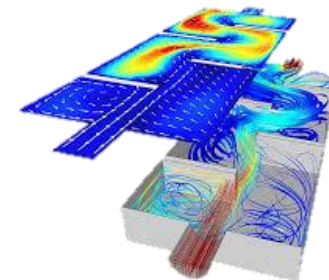
# CUDA programming language

- Introduced in 2007 with NVIDIA Tesla architecture

- "C-like" language to express programs that run on GPUs using the compute-mode hardware interface

- Relatively low-level abstractions: closely match the capabilities/performance characteristics of modern GPUs
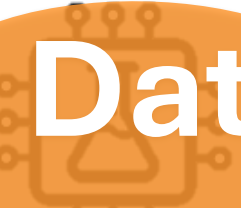
- Bioinformatics

- Computational Chemistry

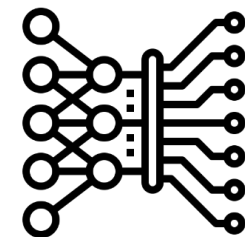- Computational Finance

- Computational Fluid Dynamics

- AI & Machine Learning

- Computer Vision

- Block Chain

- Data science, medical imaging,, weather and climate, …

**Data Parallel Applications**

# Why use GPGPU?

# Why can't CPUs Scale?

- # cores: Why Not Just Add More?
- SIMD width: Why Not Go Wider?
- Memory bandwidth: Why Not Boost It?

# Why GPGPU?

- **CPU**
  - ~10s cores
  - Low **latency**
  - Good for serial processing
  - Good for interactive tasks
  - Task parallelism
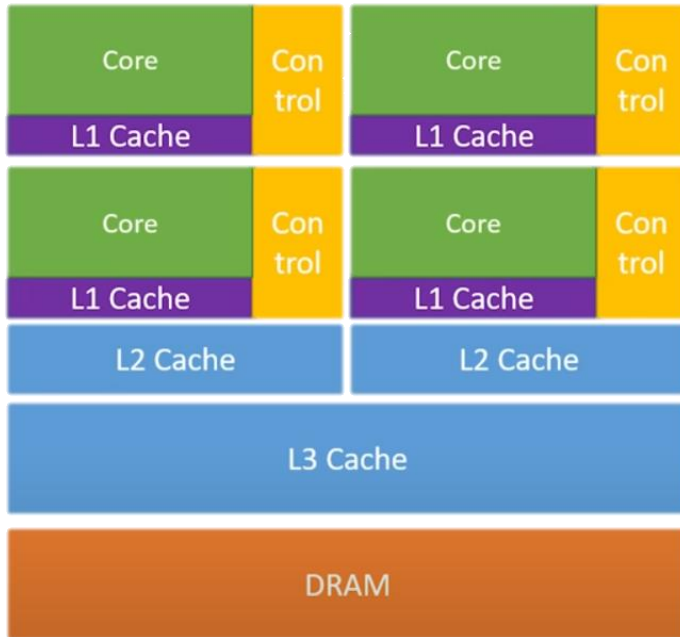
- **GPU**
  - 100s ~ 1000s cores
  - High **throughput**
  - Good for parallel processing
  - Good for big-data tasks
  - **Data** parallelism

# Why GPU?



CPU                    GPU

# What's Better?

- Scooter

- Sport car

# What's Better?

- Scooter

- Sport car



Deliver many packages within a reasonable timescale

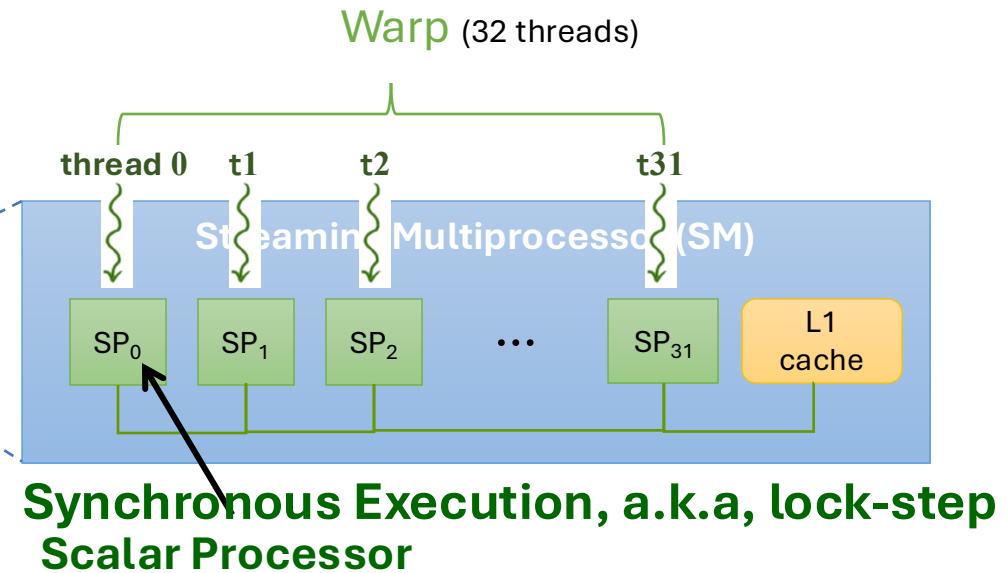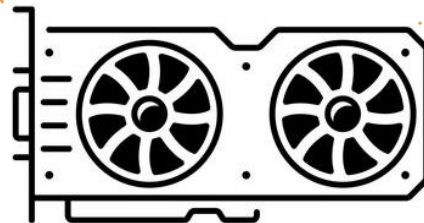Deliver a package as soon as possible

# Why can't CPUs Scale?

- # cores: Why Not Just Add More?

- SIMD width: Why Not Go Wider?

- Memory bandwidth: Why Not Boost It?

- The short answer is: they could, but they shouldn't.

- CPUs are latency oriented, GPUs are throughput oriented

- Power, Heat, and Cost Constraints

# GPU Architecture

## SIMT: single-instruction multiple threads

In A100 GPU each SM has:
- 64 FP32 cores
- 32 FP64 cores
- 64 KB registers
- 192 KB L1 + shared mem



Warp (32 threads)

thread 0    t1    t2    t31

Streaming Multiprocessor (SM)

$SP_0$    $SP_1$    $SP_2$    ...    $SP_{31}$    L1 cache

**Synchronous Execution, a.k.a, lock-step Scalar Processor**

GPU

SM    SM    SM

SM    SM    SM

L2 cache

Global Memory

# Warp Scheduler

# Warp Scheduler

# GPU Memory Hierarchy

# System Architecture With a GPU (2019)



GDDR5: 100s GB/s, 10s of GB
HBM2: ~1 TB/s, 10s of GB

GPU Memory (GDDR5, HBM2,...)

GPU

CPU

DDR4 2666 MHz
128 GB/s
100s of GB

I/O Hub (IOH)

NVMe

Network Interface

...

Host Memory (DDR4,...)

QPI/UPI
12.8 GB/s (QPI)
20.8 GB/s (UPI)

PCIe
16-lane PCIe Gen3: 16 GB/s
...

# Outline

- GPU Architecture
  - Three major ideas that make GPU processing cores run fast
  - Closer look at real GPU designs
  - The GPU memory hierarchy: moving data to processors

- General Purpose GPU (GPGPU)

- **GPU Programming Model**

- Program a GPU with CUDA

# Heterogeneous Programming (CPU+GPU)

- **CPU (Milan)**
  - □ 64 cores, 2 threads per core
  - □ 4 NUMA domains per socket
  - □ 2xAVX2 (256 bit), so 2x4 in double precision
  - □ ~512 way parallelism (64*2*4)

- **GPU (A100)**
  - □ 108 SM
  - □ 64 warps per SM (32 active at a time)
  - □ 32 threads per warp in double precision
  - □ ~200,000+ way parallelism (108*64*32) per GPU



CPU

CPU Memory

PCIe ~ 32 GB/Sec

PCI bus

GPU

GPU Memory

1,500 GB/Sec Memory Bandwidth

# Heterogeneous Programming (CPU+GPU)

# Data Movement



"kernel" function

```
data = open("input.dat");          # read the data on the CPU
copyToGPU(data);                   # copy the data to the GPU
matrix_inverse(data.gpu);          # perform a matrix operation on the GPU
copyFromGPU(data);                 # copy the resulting output to the CPU
write(data, "output.dat");         # write the output to file on the CPU
```

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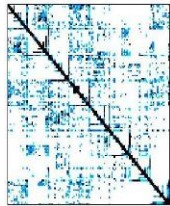  - CUDA, OpenMP, and TBB backends
  - Extensible and customizable
  - Integrates with existing software

- 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

```fortran
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

# Outline

- GPU Architecture
  - ◻ Three major ideas that make GPU processing cores run fast
  - ◻ Closer look at real GPU designs
  - ◻ The GPU memory hierarchy: moving data to processors

- General Purpose GPU (GPGPU)

- GPU Programming Model

- **Program a GPU with CUDA**

**CONCEPTS**

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

# PROGRAM A GPU WITH CUDA

- 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

CPU

Compute-Intensive Functions

Rest of Sequential CPU Code

Use GPU to Parallelize

+

# 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

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

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

- Here `__global__` is a CUDA C/C++ keyword meaning
  - `add()` will execute on the device
  - `add()` will be called from the host

**Vector Addition**



a     b     c

- 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()`

```cpp
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;
}
```

# RUNNING IN PARALLEL

**CONCEPTS**

Heterogeneous Computing

Blocks

Threads

Indexing

Shared memory

__syncthreads()

Asynchronous operation

Handling errors

Managing devices

# Multi-level Threading Model



CUDA Hierarchy of threads, blocks, and grids, with corresponding per-thread private, per-block shared, and per-application global memory spaces.

**Software**

Thread

Thread Block

Grid

**Hardware**

Scalar Processor

Multiprocessor

Device

Threads are executed by scalar processors

Thread blocks are executed on multiprocessors

Thread blocks do not migrate

Several concurrent thread blocks can reside on one multiprocessor - limited by multiprocessor resources (shared memory and register file)

A kernel is launched as a grid of thread blocks

# 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

# Hierarchical Programming Model

- Multi-level hierarchy
  - The thread is the smallest unit of program execution
  - Threads are organized in a thread block
  - Blocks form a grid
  - Block and grid have up to three dimensions

- Threads within a block are further implicitly divided into warps (32 or 64 threads)

- Choose a certain level of granularity given a problem

- 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()…

```
#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);
```

```cuda
// 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

CONCEPTS

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

# INTRODUCING 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()`...

```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 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;

}
```

**"kernel" function**

# CONCEPTS

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

# COMBINING THREADS AND BLOCKS

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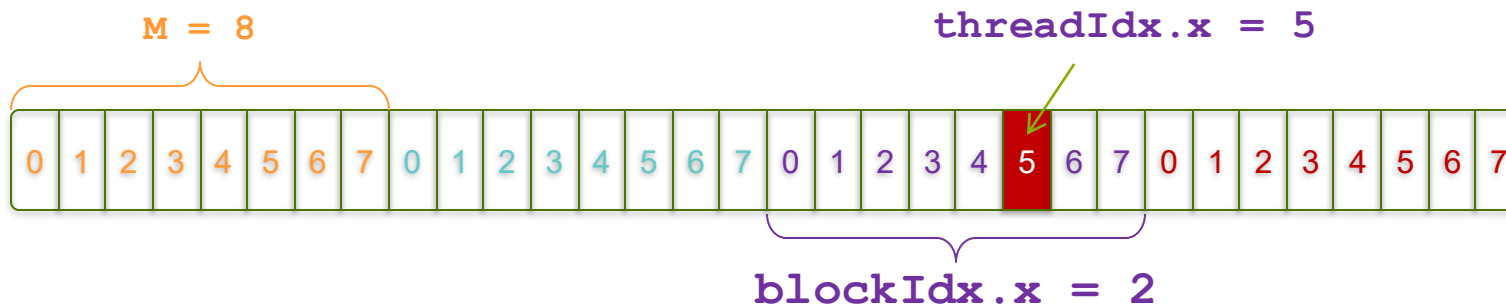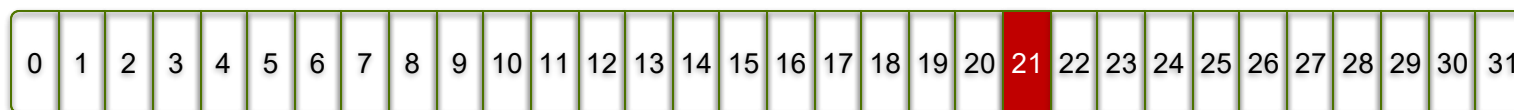
- With M threads/block a unique index for each thread is given by

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

- Which thread will operate on the red element?



$$\texttt{int index = threadIdx.x + blockIdx.x * M;}$$
$$= \quad 5 \quad + \quad 2 \quad * 8;$$
$$= 21;$$

# 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); 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];
}
```

**Check bound**

- Update the kernel launch:

```
add<<<(N-1) / M + 1,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…

**CONCEPTS**

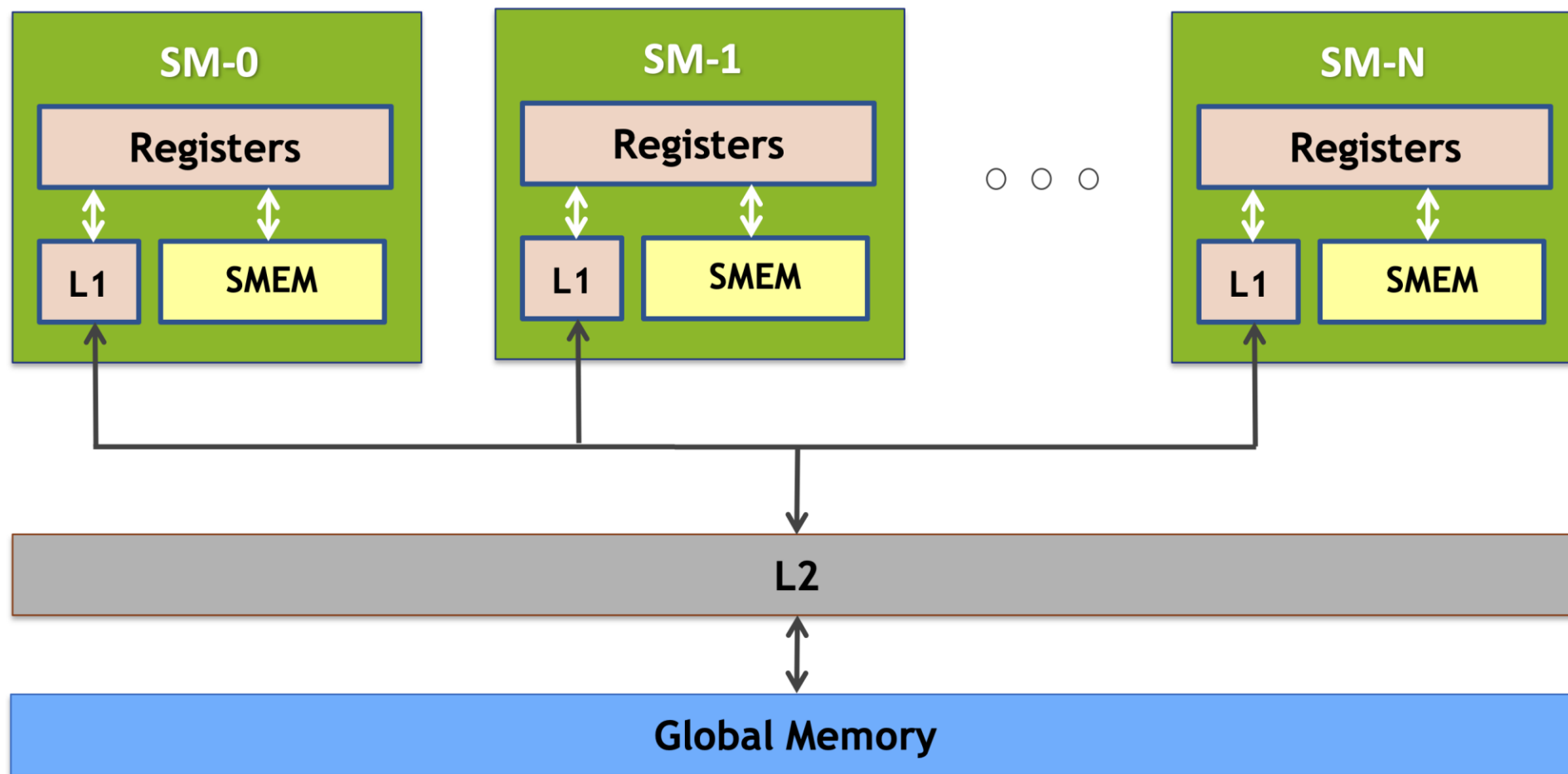Heterogeneous Computing

Blocks

Threads

Indexing

Shared memory

__syncthreads()

Asynchronous operation

Handling errors

Managing devices
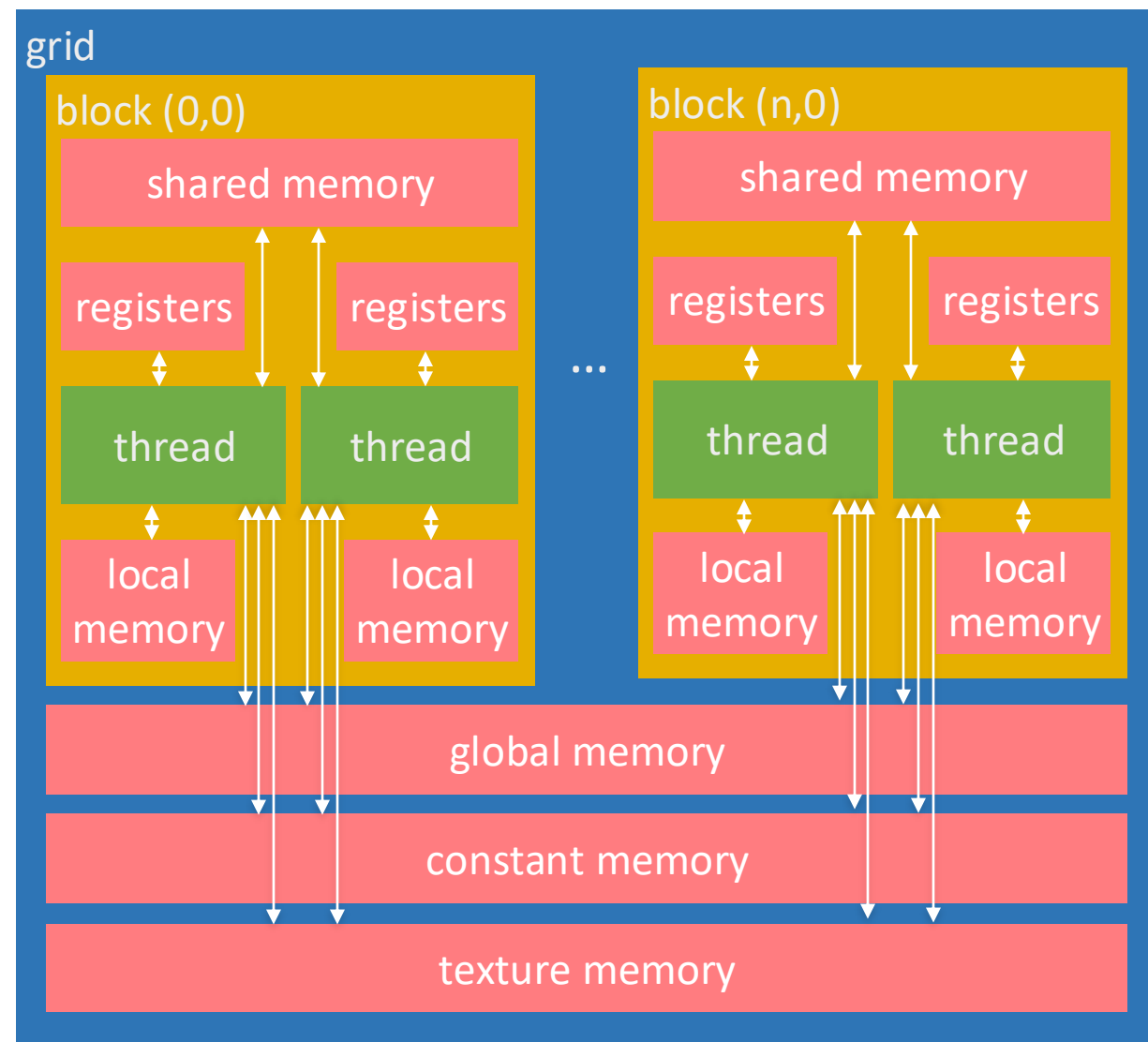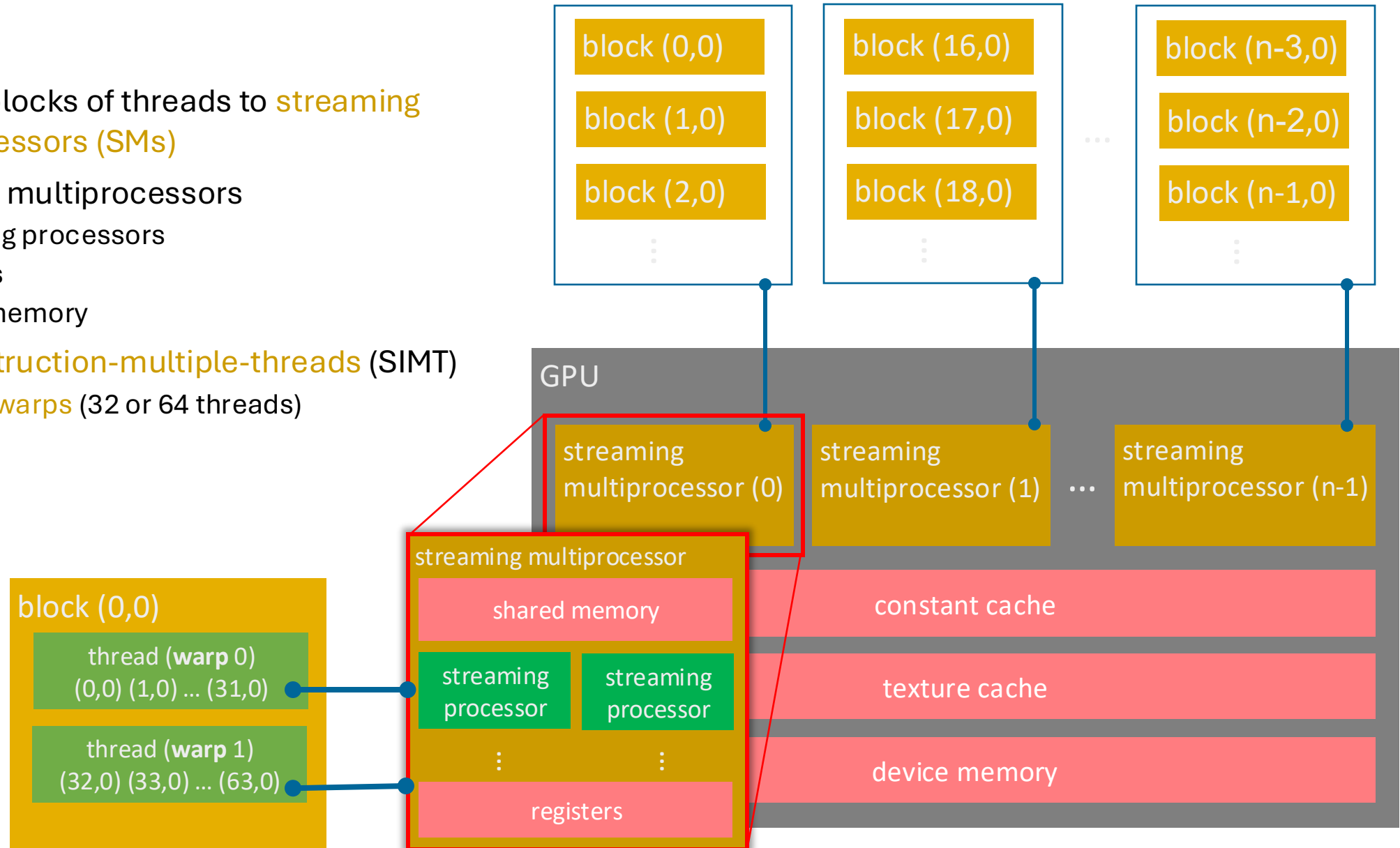
# COOPERATING THREADS

# GPU memory hierarchy

# Memory Model

- Define memory hierarchy
- Registers (VGPR, SGPR)
  - Local variables per thread
  - The fastest memory
- Local memory
  - local variables per thread
  - Slow off-chip memory (VRAM)
  - Register spills if not enough registers
- Shared memory (Local data share)
  - Shared between threads in a block
  - Faster on-chip memory
- Global memory
  - Shared with all blocks
  - Largest memory
  - Slow off-chip memory (VRAM)
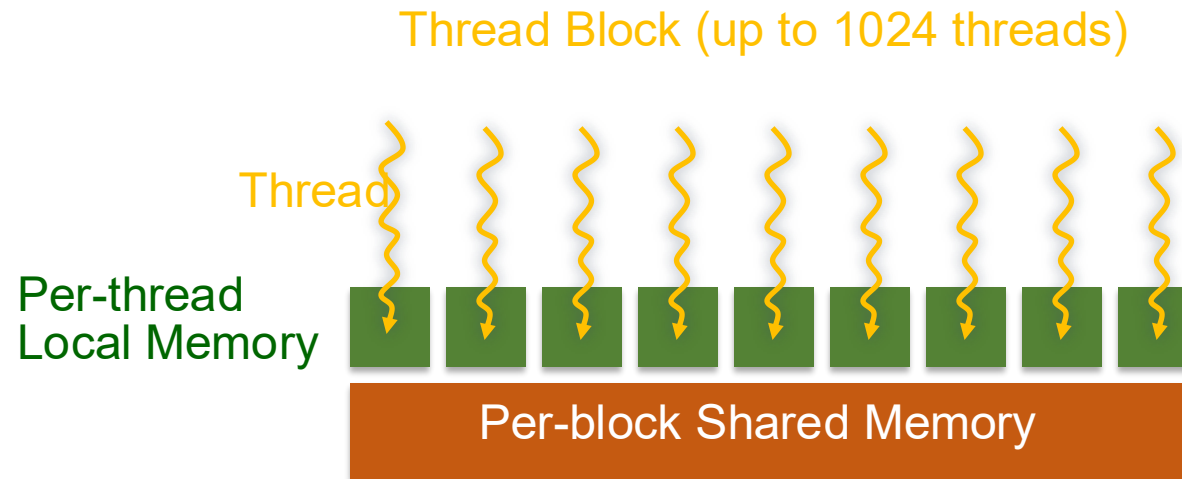- Texture and constant memory
  - Cached differently, on off-chip memory

# Execution Model

- Mapping blocks of threads to streaming multiprocessors (SMs)
- Streaming multiprocessors
  - Streaming processors
  - Registers
  - Shared memory
- Single-instruction-multiple-threads (SIMT)
  - Process warps (32 or 64 threads)

block (0,0)

block (1,0)

block (2,0)

⋮

block (16,0)

block (17,0)

block (18,0)

⋮

...

block (n-3,0)

block (n-2,0)

block (n-1,0)

⋮

GPU

streaming multiprocessor (0)

streaming multiprocessor (1)

···

streaming multiprocessor (n-1)

constant cache

texture cache

device memory

block (0,0)

thread (**warp** 0)
(0,0) (1,0) ... (31,0)

thread (**warp** 1)
(32,0) (33,0) ... (63,0)

streaming multiprocessor

shared memory

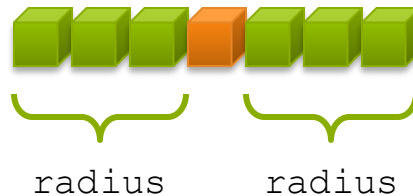streaming processor

streaming processor

⋮   ⋮

registers

# Shared (within a block) Memory

- Declare using __shared__, allocated per block
- Fast on-chip memory, user-managed
- Not visible to threads in other blocks

Thread Block (up to 1024 threads)
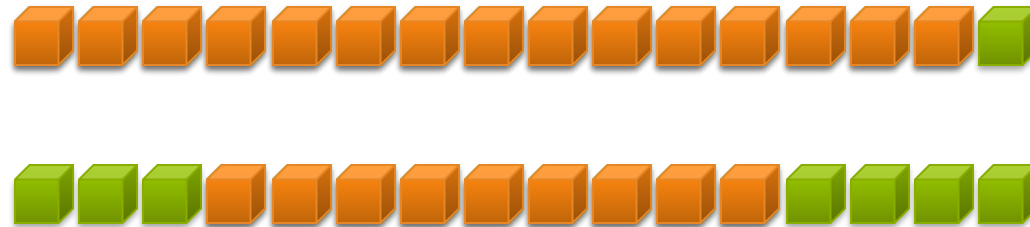
Thread

Per-thread
Local Memory

Per-block Shared Memory

- 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

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

- Input elements are read several times
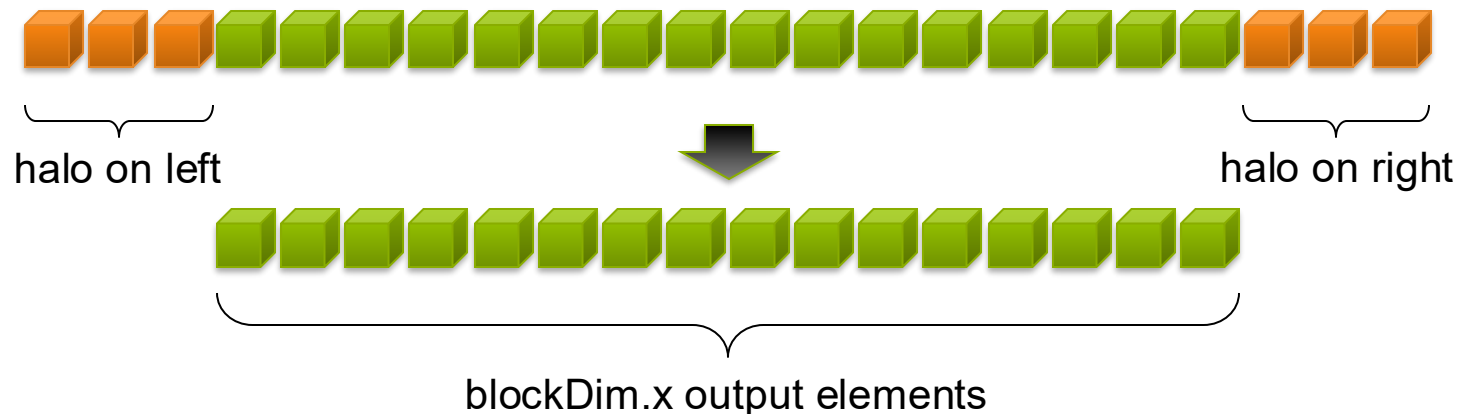  - With radius 3, each input element is read <span style="color:red">seven</span> times

# 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

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

  - Each block needs a <span style="color:orange">halo</span> 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];
  }
```



**BLOCK_SIZE**

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


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

# Any problem?

# 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

```
__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

```c
// 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 <span style="color:orange">asynchronous</span>
  - 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

# CUDA Summary

- GPU Architecture for Data Parallelism

- GPU Programming Model: SIMT

- 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()`