



# Introduction to Parallel Computing and GPU programming with CUDA

Simone Balducci, **Felice Pantaleo**, Aurora Perego

CERN Experimental Physics Department

# Content of the theoretical session



- Introduction to Parallel computing
- Heterogeneous Parallel computing systems
- CUDA Basics
- Parallel kernels
- Shared Memory and Synchronization
- Device Management

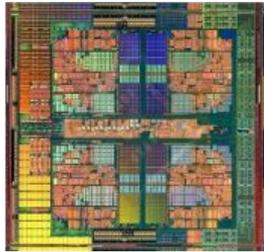


# Content of the tutorial session

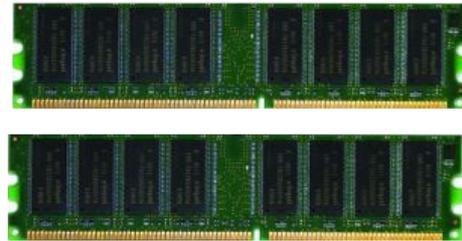
- Write and launch kernels
- Manage GPU memory
- Manage communication and synchronization
- From C++ standard algorithms to GPU execution using thrust

# Past systems

- Computing landscape is very different from 10-20 years ago
- Every component and its interfaces are being re-examined



Microprocessor



Main Memory



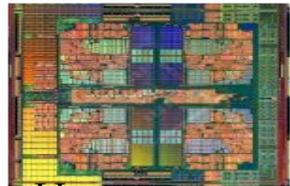
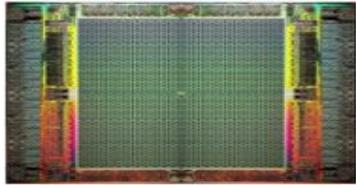
Storage (SSD/HDD)

# Modern systems



- Applications and technology demand novel architectures
  - Driven by huge hunger for data (Big Data), new applications (ML/AI, graph analytics, genomics), ever-greater realism
  - We can easily collect more data than we can analyze/understand
  - Five walls: Energy, reliability, complexity, security, scalability

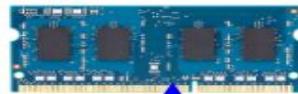
FPGA



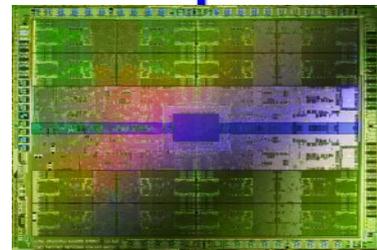
Heterogeneous Processors and Accelerators



Hybrid Memory



Persistent memory/Storage

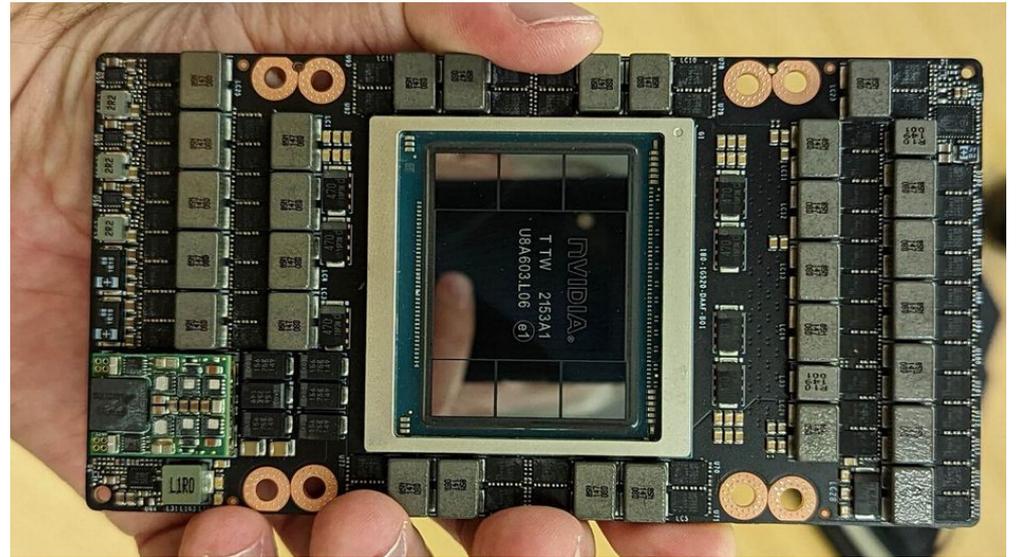


GPU

# Accelerators



- Exceptional raw power and memory bandwidth wrt CPUs
- Lower energy to solution
- Massively parallel architecture
- Low Memory/core

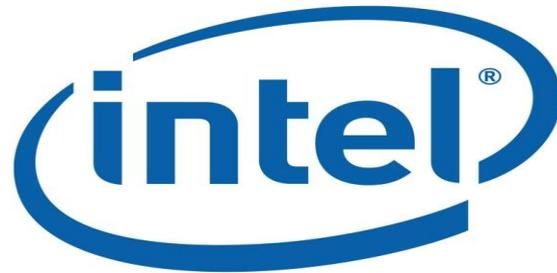
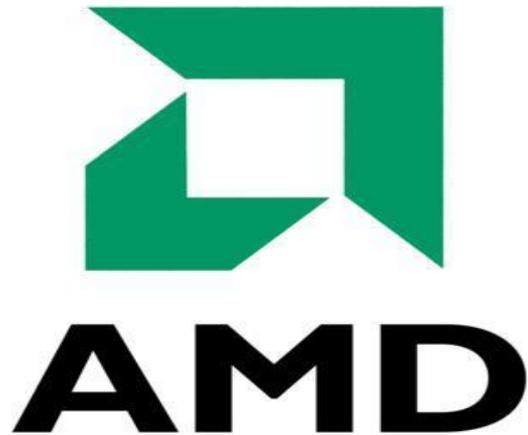


# Accelerators

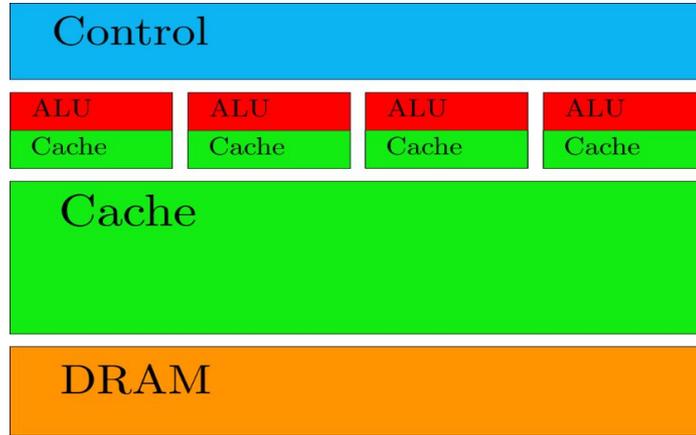


GPUs were traditionally used for real-time rendering/gaming.

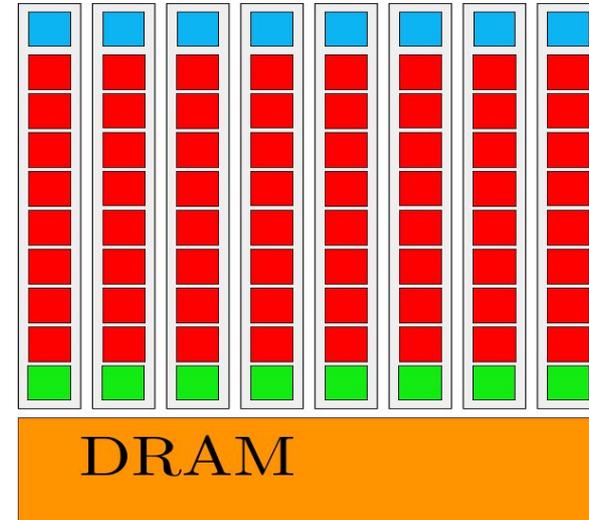
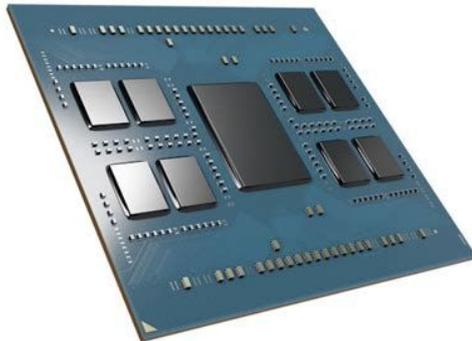
AMD and NVIDIA main manufacturers for discrete GPUs, Intel for integrated ones



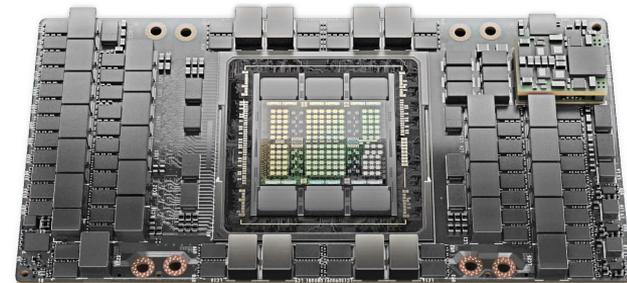
# CPU vs GPU architectures



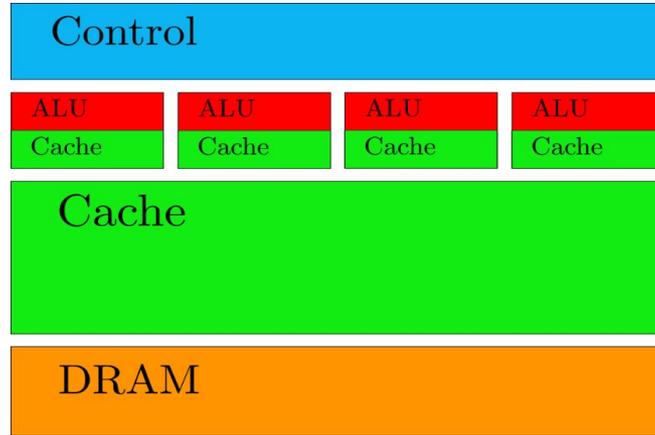
**CPU**



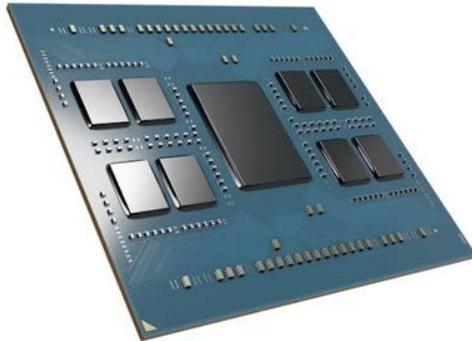
**GPU**



# CPU vs GPU architectures



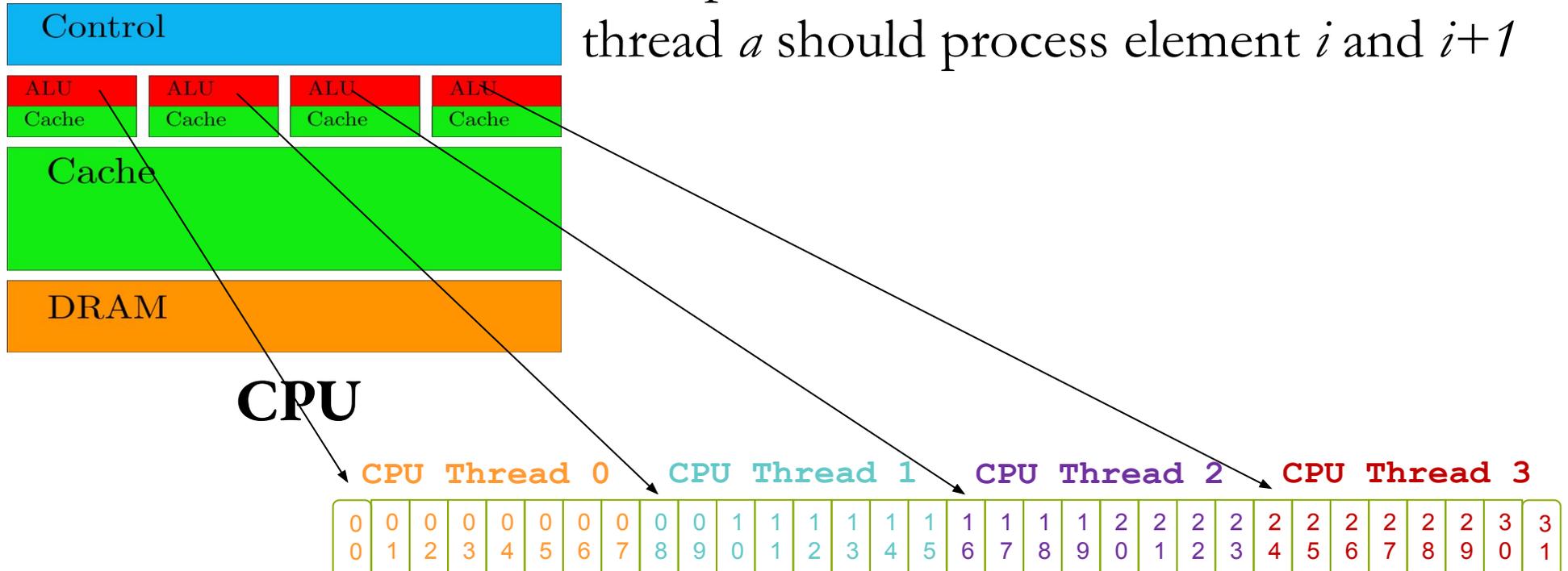
**CPU**



- Large caches (slow memory accesses to quick cache accesses)
- SIMD
- Branch prediction/speculative
- Powerful ALU
- Pipelining

# Memory access patterns: cached

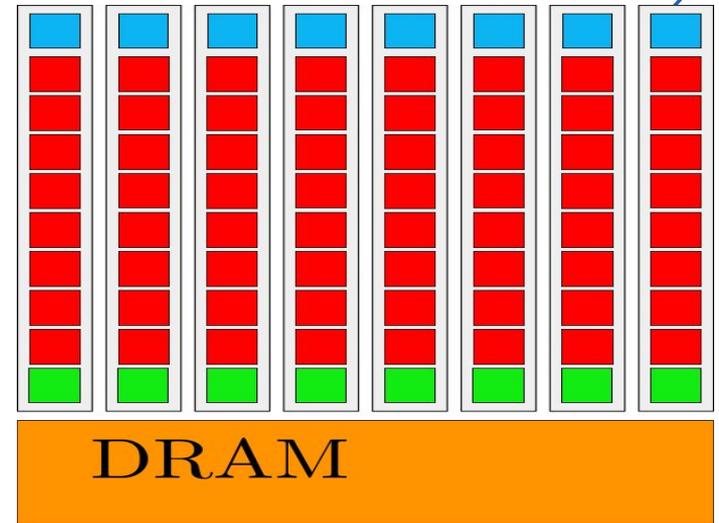
For optimal CPU cache utilization, the thread  $a$  should process element  $i$  and  $i+1$



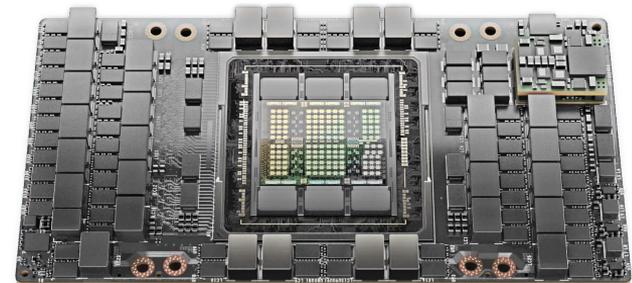
# CPU vs GPU architectures



- Hundreds of “cores” (e.g. streaming multiprocessors, Xe cores, compute units)
- SIMT (Single-Instruction, Multiple-Thread) with hundreds of SIMD-like warps in fly
- Instructions pipelined
- Thread-level parallelism
- Instructions issued in order
- Branch predication

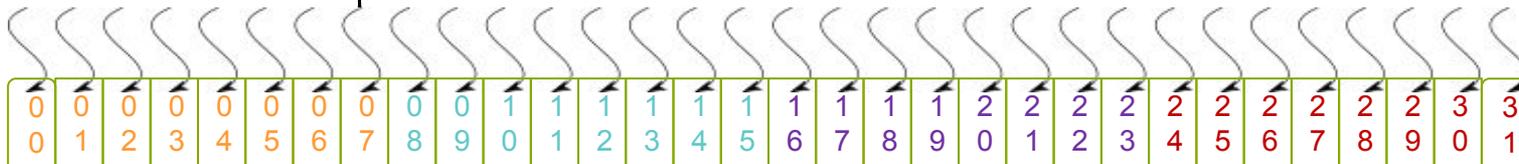


**GPU**



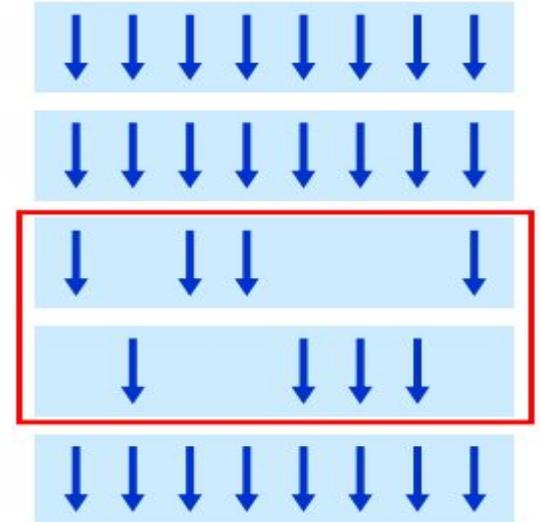
# Inside a GPU SM: coalesced

- L1 data cache shared among ALUs
- ALUs work in SIMD mode in groups called warps
  - Think about it as vectors on the same CPU core
- If a *load* is issued by each thread, they have to wait for all the loads in the same warp to complete before the next instruction can execute
- Coalesced memory access pattern optimal for GPUs: thread  $a$  should process element  $i$ , thread  $a+1$  the element and  $i+1$ 
  - Lose an order of magnitude in performance if cached access pattern used on GPU



# Warps

- Once a block is assigned to an SM, it is divided into units called warps.
- Thread IDs within a warp are consecutive and increasing
- Threads within a warp are executed in a SIMD fashion
- If an operand is not ready the warp will stall
- Context switch between warps when stalled
- Context switch must be very fast
- Typical values of warp size 16, 32, 64 depending on vendor





# Heterogeneous Parallel Computing Systems

# Heterogeneous Computing



- Terminology
  - Host The CPU and its memory space
  - Device The GPU and its memory space

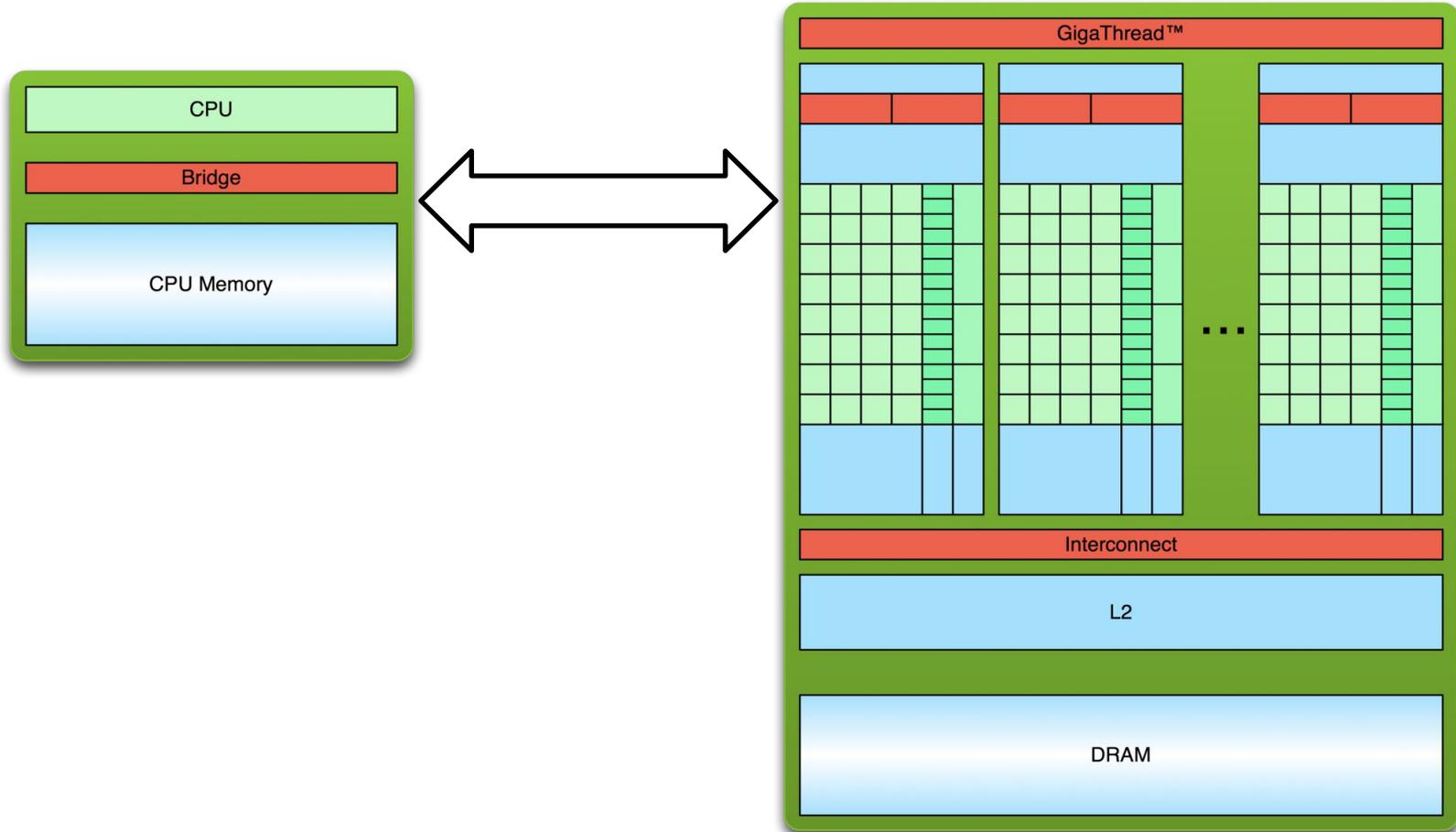


Host

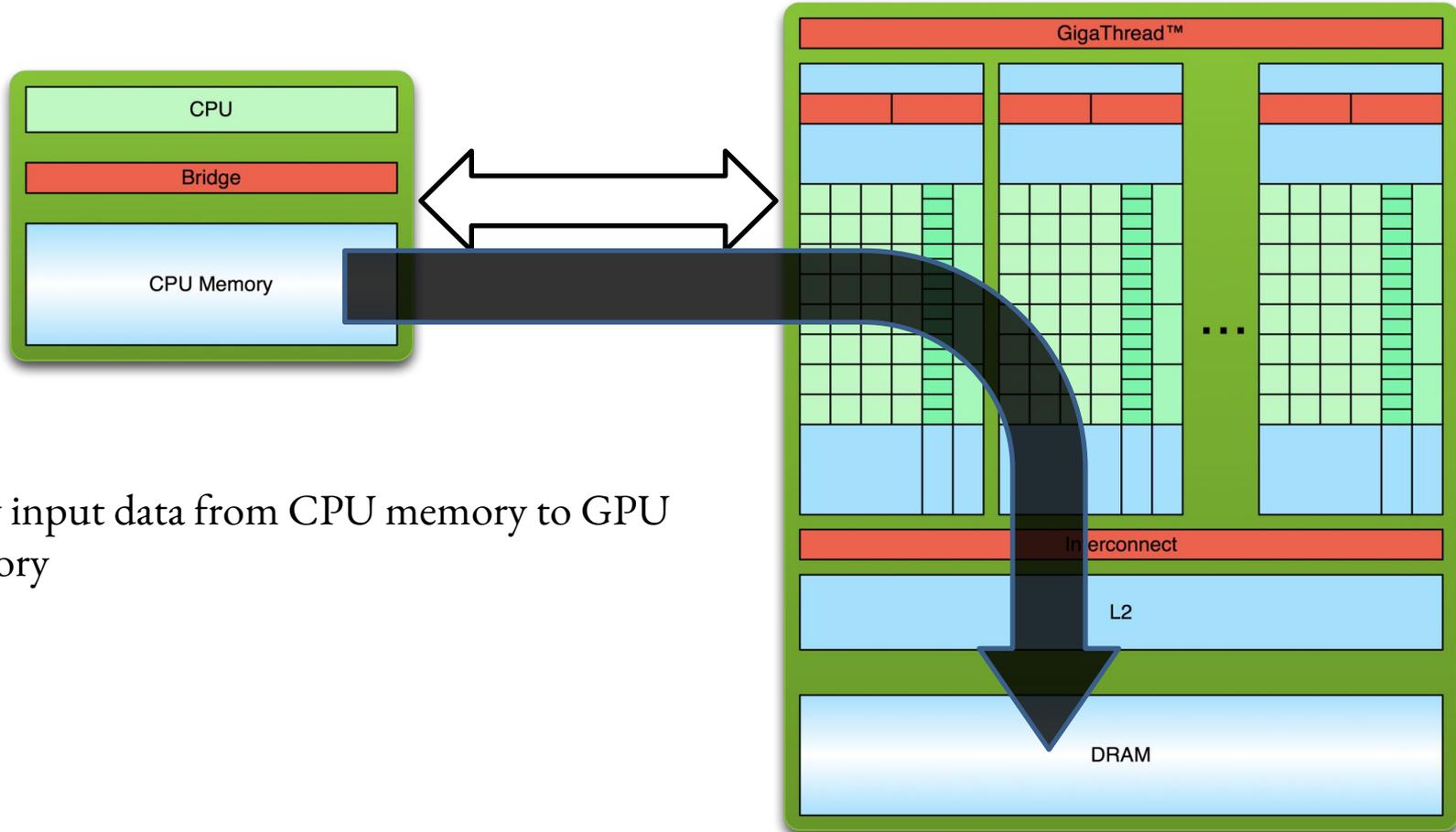


Device

# Simple Processing Flow

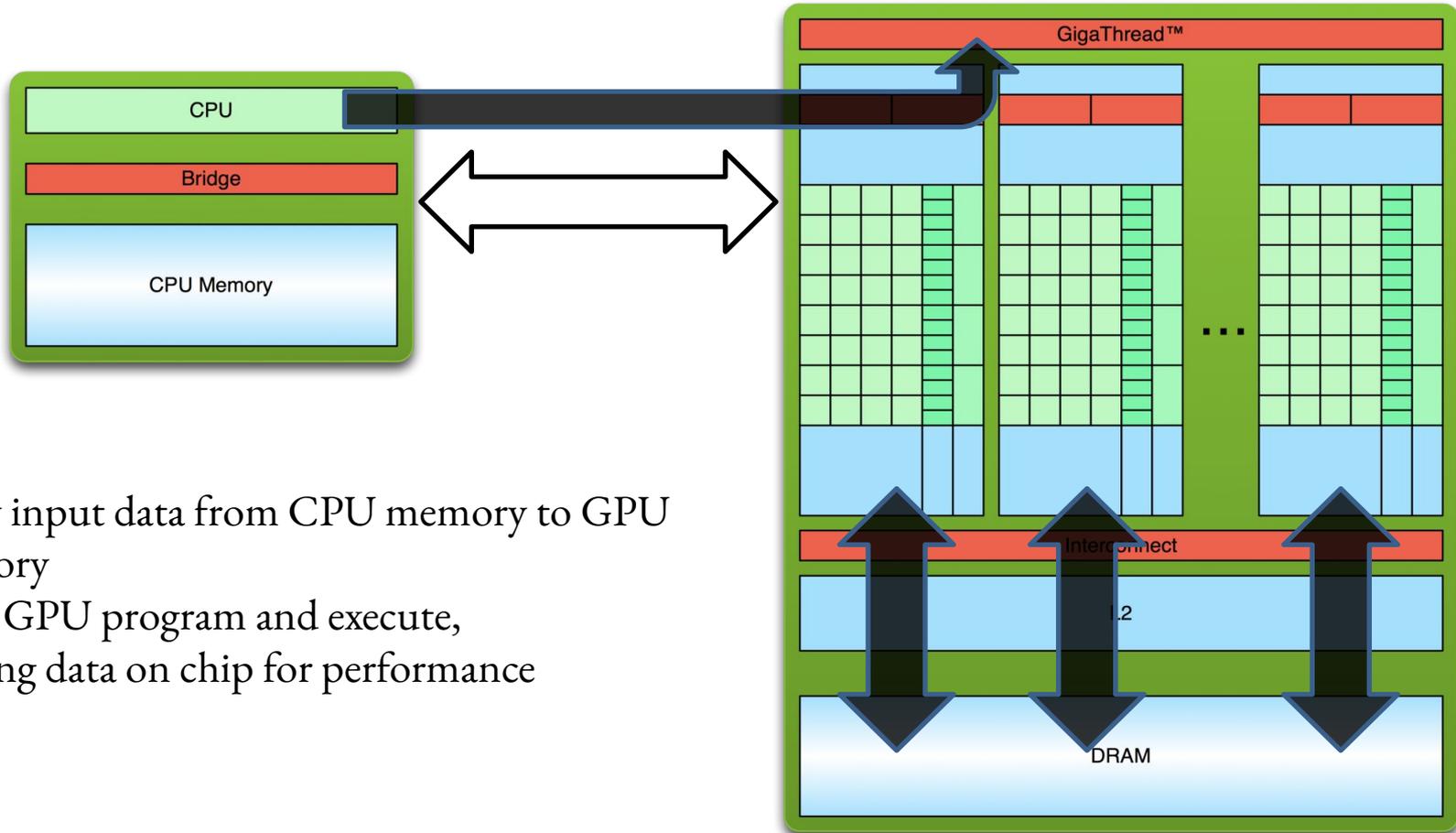


# Simple Processing Flow



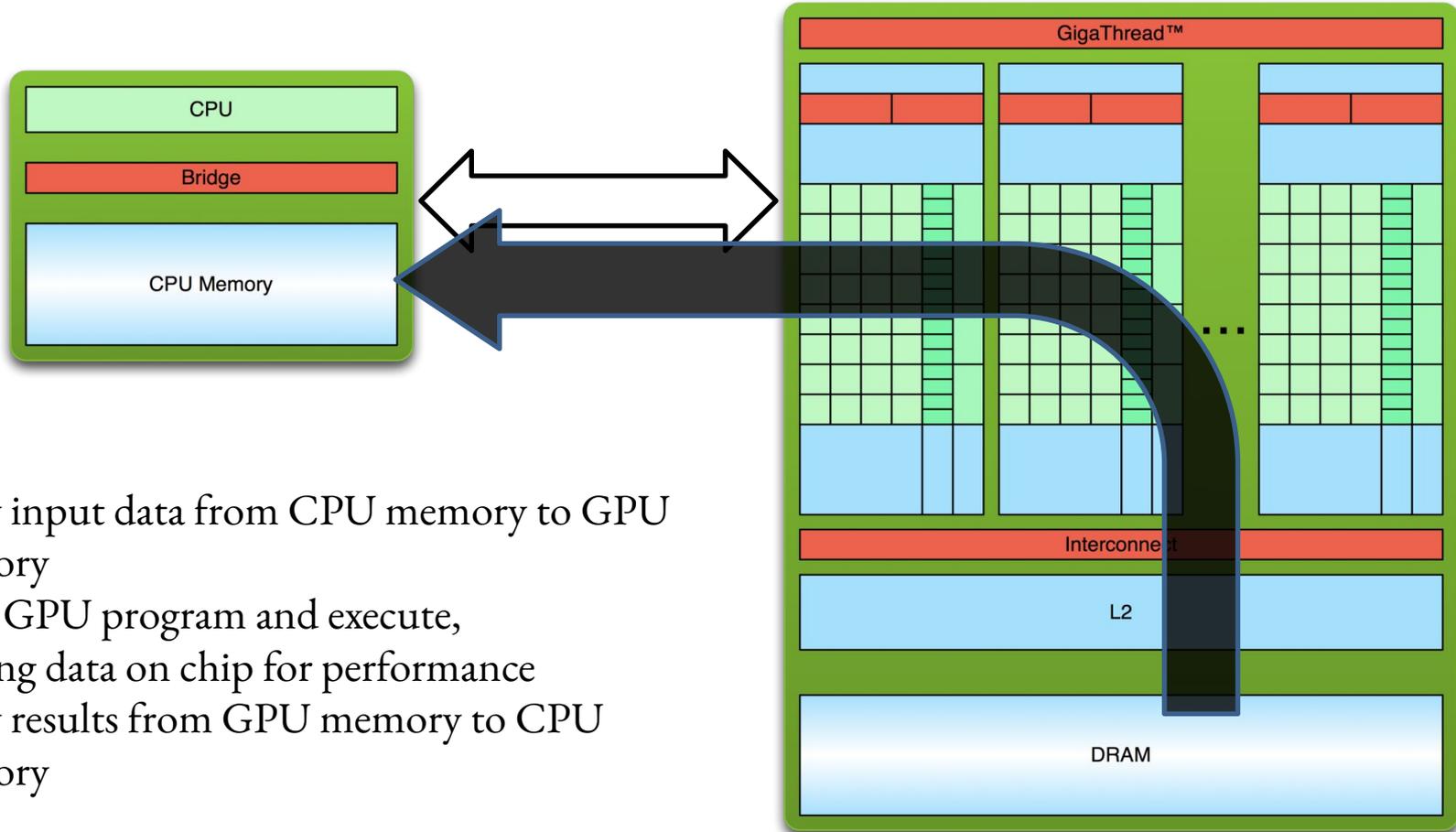
1. Copy input data from CPU memory to GPU memory

# 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



# Basics

- Small set of extensions to enable asynchronous heterogeneous computing using NVIDIA GPUs
- Straightforward APIs to manage devices, memory etc.
- Concepts learned in CUDA apply to other parallel frameworks (OpenMP, SYCL, HIP, alpaka).
- Learning Curve
  - Initial effort focuses on understanding new APIs and language constructs.
  - Core programming logic often remains unchanged.

# SPMD Phases

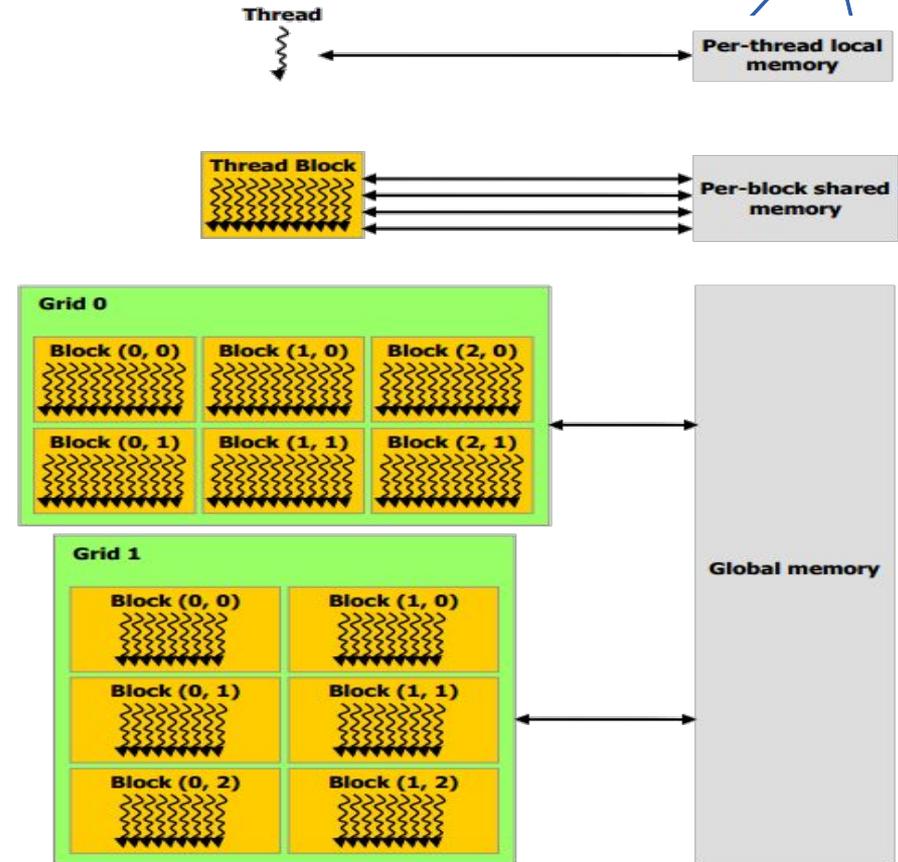


- Initialize
  - Establish localized data structure and communication channels
- Obtain a unique identifier
  - Each thread acquires a unique identifier, typically range from 0 to  $N-1$ , where  $N$  is the number of threads
- Distribute Data
  - Decompose global data into chunks and localize them, or
  - Sharing/replicating major data structure using thread ID to associate subset of the data to threads
- Run the core computation
- Finalize
  - Reconcile global data structure, prepare for the next major iteration

# Memory Hierarchy in GPU programming



- Registers/Shared memory:
  - Fast
  - Only accessible by the thread/block
  - Lifetime of the thread/block
- Global memory:
  - Potentially 150x slower than register or shared memory
  - Accessible from either the host or device
  - Lifetime of the application



# Hello World!



```
1 #include <iostream>
2
3 int main() {
4     std::cout << "Hello World!\n";
5 }
```

Standard C++ that runs on the host  
nvcc can be used to compile programs with no *device* code

Output:

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

# Hello World! with Device Code



```
1 #include <iostream>
2
3 __global__ void mykernel() {}
4
5 int main() {
6     cudaStream_t stream; cudaStreamCreate(&stream);
7     mykernel<<<1,1,0,stream>>>();
8     std::cout << "Hello World!\n";
9     cudaStreamSynchronize(stream);
10    cudaStreamDestroy(stream);
11 }
```

# Hello World! with Device Code



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

- CUDA 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 `nvcc` compiler
  - Host functions (e.g. `main()`) processed by `gcc`

# Hello World! with Device Code



```
mykernel<<<1, 1, 0, stream>>> ();
```

- Triple angle brackets mark a call from host code to device code
  - Also called a “kernel launch”
  - We’ll return to the parameters in a moment
- That’s all that is required to execute a function on the GPU!



# Parallel constructs in CUDA

# Addition on the Device



- A simple kernel to add two integers

```
1 __global__ void add(const int *a, const int *b, int *c) {  
2     *c = *a + *b;  
3 }
```

- As before `__global__` is a CUDA 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

```
1 __global__ void add(const int *a, const int *b, int *c) {  
2     *c = *a + *b;  
3 }
```

- `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 `malloc()`, `free()`, `memcpy()`



# Addition on the Device: add()



- Returning to our add () kernel

```
1 __global__ void add(const int *a, const int *b, int *c) {  
2     *c = *a + *b;  
3 }
```

- Let's take a look at main () ...



A stream is a queue of operations to be executed on the device

```
1 int main() {
2     cudaStream_t stream;
3     cudaStreamCreate(&stream);
4     int *a, *b, *c;           // host copies of a, b, c
5     int *d_a, *d_b, *d_c;    // device copies of a, b, c
6     int size = sizeof(int);
7     // Allocate space for device copies of a, b, c
8     cudaMallocHost(&a, size);
9     cudaMallocHost(&b, size);
10    cudaMallocHost(&c, size);
11    *a = 2; *b = 7;
12    // Allocate memory on device
13    cudaMallocAsync(&d_a, size, stream);
14    cudaMallocAsync(&d_b, size, stream);
15    cudaMallocAsync(&d_c, size, stream);
16    // Copy inputs to device
17    cudaMemcpyAsync(d_a, a, size, cudaMemcpyHostToDevice, stream);
18    cudaMemcpyAsync(d_b, b, size, cudaMemcpyHostToDevice, stream);
19    // Launch add() kernel on GPU
20    add<<<1,1,0,stream>>>(d_a, d_b, d_c);
21    // Copy result back to host
22    cudaMemcpyAsync(c, d_c, size, cudaMemcpyDeviceToHost, stream);
23    cudaFreeAsync(d_a, stream);
24    cudaFreeAsync(d_b, stream);
25    cudaFreeAsync(d_c, stream);
26    // Synchronize to be able to use c
27    cudaStreamSynchronize(stream);
28    cudaStreamDestroy(stream);
29    cudaFreeHost(a); cudaFreeHost(b); cudaFreeHost(c);
30 }
```

# Coordinating Host & Device



- Kernel launches are asynchronous
  - control is returned to the host thread before the device has completed the requested task
  - CPU needs to synchronize before consuming the results

**cudaMemcpy() , cudaMemcpyAsync()**

Blocks the CPU thread until the copy/allocation is complete  
Copy/allocation begins when all preceding CUDA calls have completed

**cudaMemcpyAsync() ,  
cudaMallocAsync()**

Asynchronous, does not block the CPU thread

**cudaDeviceSynchronize()**

Blocks the CPU thread until all preceding CUDA calls have completed

**cudaStreamSynchronize(stream)**

Blocks the CPU thread until all preceding CUDA calls in the stream have completed

# Moving to Parallel



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

```
add<<< 1, 1, 0, stream >>> ();
```



```
add<<< N, 1, 0, stream >>> ();
```

- 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
  - The set of blocks is referred to as a grid
  - Each invocation can refer to its block index using `blockIdx.x`

```
1 __global__ void add(const int *a, const int *b, int *c)
2 {
3     c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];
4 }
```

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

# Remember SPMD?



```
1 __global__ void add(const int *a, const int *b, int *c)
2 {
3     c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];
4 }
```

- 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

```
1 __global__ void add(const int *a, const int *b, int *c)
2 {
3     c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];
4 }
```

- Let's take a look at main () ...



```
1 int main() {
2     cudaStream_t stream; cudaStreamCreate(&stream);
3     int N = 512;
4     std::vector<int> a, b, c;
5     a.resize(N); b.resize(N); c.resize(N);
6     int *d_a, *d_b, *d_c; // device copies of a, b, c
7     int size = N * sizeof(int);
8     // Alloc space for host copies of a, b, c and
9     // setup input values
10    my_favorite_random_ints(a, N);
11    my_favorite_random_ints(b, N);
12    // Alloc memory for device copies of a, b, c
13    cudaMallocAsync(&d_a, size, stream);
14    cudaMallocAsync(&d_b, size, stream);
15    cudaMallocAsync(&d_c, size, stream);
16    // Copy inputs to device
17    cudaMemcpyAsync(d_a, a.data(), size, cudaMemcpyHostToDevice, stream);
18    cudaMemcpyAsync(d_b, b.data(), size, cudaMemcpyHostToDevice, stream);
19    // Launch add() kernel on GPU with N blocks
20    add<<<N, 1, 0, stream>>>(d_a, d_b, d_c);
21    // Copy result back to host
22    cudaMemcpyAsync(c.data(), d_c, size, cudaMemcpyDeviceToHost, stream);
23    // Cleanup
24    cudaFreeAsync(d_a, stream);
25    cudaFreeAsync(d_b, stream);
26    cudaFreeAsync(d_c, stream);
27    cudaStreamSynchronize(stream);
28    // Now you can use content of the c vector...
29    cudaStreamDestroy(stream);
30 }
```

# CUDA Threads



- Terminology: a block can be split into parallel threads
- Let's change `add()` to use parallel threads instead of parallel blocks

```
1 __global__ void add(const int *a, const int *b, int *c) {  
2     c[threadIdx.x] = a[threadIdx.x] + b[threadIdx.x];  
3 }
```

- We use `threadIdx.x` instead of `blockIdx.x`
- Need to make one change in `main()` ...

# 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

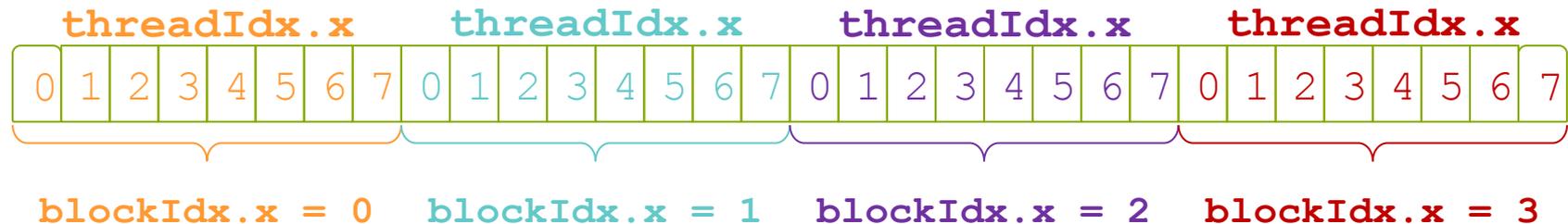
Why? We'll come to that...

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)



With `blockDim.x` threads/block a unique index for each thread is given by:  
`auto index = threadIdx.x + blockIdx.x * blockDim.x;`

# Vector Addition with Blocks and Threads



- Use the built-in variable `blockDim.x` for threads per block  
`auto index = threadIdx.x + blockIdx.x * blockDim.x;`
- Combined version of `add()` to use parallel threads *and* parallel blocks

```
1 __global__ void add(const int *a, const int *b, int *c) {  
2     auto index = threadIdx.x + blockIdx.x * blockDim.x;  
3     c[index] = a[index] + b[index];  
4 }
```

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



```
1 int main() {
2     cudaStream_t stream; cudaStreamCreate(&stream);
3     int N = 2048*2048;
4     int threads_per_block = 512;
5     std::vector<int> a, b, c;
6     a.resize(N); b.resize(N); c.resize(N);
7     int *d_a, *d_b, *d_c; // device copies of a, b, c
8     int size = N * sizeof(int);
9     // Alloc space for host copies of a, b, c and
10    // setup input values
11    my_favorite_random_ints(a, N);
12    my_favorite_random_ints(b, N);
13    // Alloc memory for device copies of a, b, c
14    cudaMallocAsync(&d_a, size, stream);
15    cudaMallocAsync(&d_b, size, stream);
16    cudaMallocAsync(&d_c, size, stream);
17    // Copy inputs to device
18    cudaMemcpyAsync(d_a, a.data(), size, cudaMemcpyHostToDevice, stream);
19    cudaMemcpyAsync(d_b, b.data(), size, cudaMemcpyHostToDevice, stream);
20    // Launch add() kernel on GPU with N blocks
21    add<<<N/threads_per_block, threads_per_block, 0, stream>>>(d_a, d_b, d_c);
22    // Copy result back to host
23    cudaMemcpyAsync(c.data(), d_c, size, cudaMemcpyDeviceToHost, stream);
24    // Cleanup
25    cudaFreeAsync(d_a, stream); cudaFreeAsync(d_b, stream); cudaFreeAsync(d_c, stream);
26    cudaStreamSynchronize(stream);
27    // Now you can use content of the c vector...
28    cudaStreamDestroy(stream);
29 }
```

# Handling Arbitrary Vector Sizes



- Typical problems are not friendly multiples of `blockDim.x`
- Avoid accessing beyond the end of the arrays:

```
1 __global__ void add(const int *a, const int *b, int *c, int n) {  
2     auto index = threadIdx.x + blockIdx.x * blockDim.x;  
3     if (index < n)  
4         c[index] = a[index] + b[index];  
5 }
```

Update the kernel launch:

```
add<<<std::ceil((float)n / nThPerBlock), nThPerBlock>>>(d_a, d_b, d_c, n);
```



Time for the first three exercises!

<https://infn-esc.github.io/sesame26/gpu/cuda.html>



# Shared Memory with CUDA

# Why Bother with Threads?

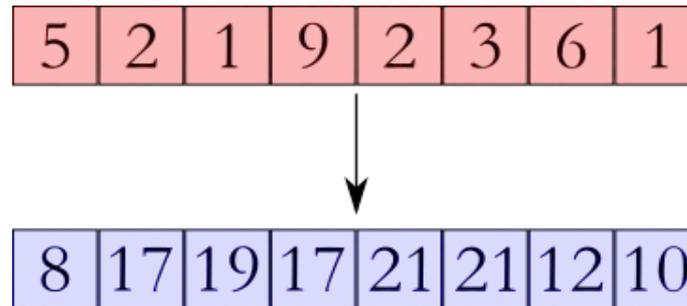


- Threads seem unnecessary
  - They add a level of complexity
  - What do we gain?
- Unlike parallel blocks, threads have mechanisms to:
  - Communicate
  - Synchronize
- To understand the gain, we need a new example...

# 1D Stencil



- Consider applying a 1D stencil sum to a 1D array of elements
  - Each output element is the sum of input elements within a radius
  - Example of stencil with radius 2:



# Sharing Data Between Threads



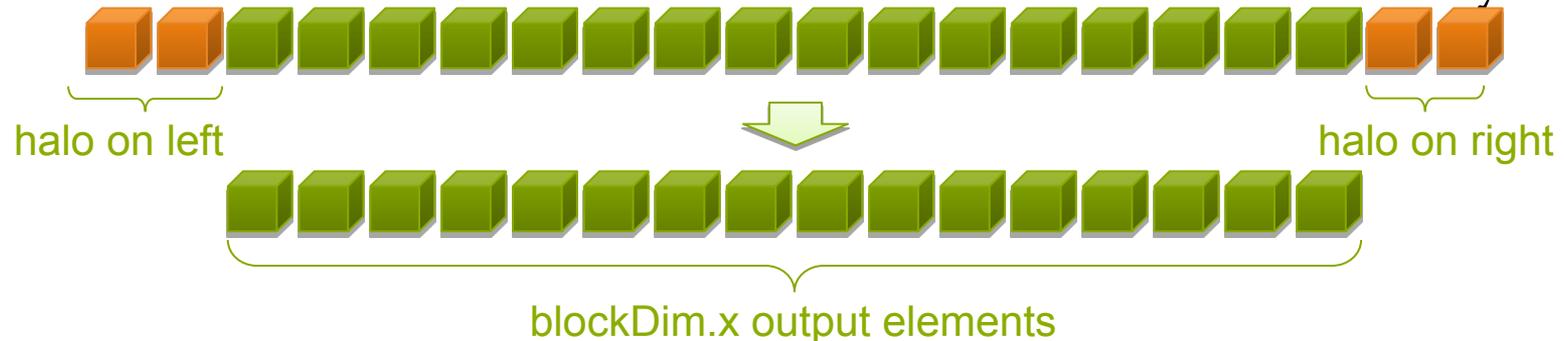
- Terminology: within a block, threads share data via shared memory
- Extremely fast on-chip memory, user-managed
- Declare using `__shared__`, allocated per block
- Data is not visible to threads in other blocks

# Implementing With Shared Memory



## Cache data in shared memory

- Read  $(\text{blockDim.x} + 2 * \text{radius})$  input elements from global memory to shared memory
- Compute  $\text{blockDim.x}$  output elements
- Write  $\text{blockDim.x}$  output elements to global memory
- Each block needs a halo of radius elements at each boundary



# Stencil Kernel



```
1 __global__ void stencil_1d(const int *in, int *out, int n) {
2   __shared__ int temp[BLOCK_SIZE + 2 * RADIUS];
3   auto g_index = threadIdx.x + blockIdx.x * blockDim.x;
4   if (g_index < n) {
5     auto s_index = threadIdx.x + RADIUS;
6
7     // Read input elements into shared memory
8     temp[s_index] = in[g_index];
9     if (threadIdx.x < RADIUS) {
10      temp[s_index - RADIUS] = g_index - RADIUS < 0? 0 :
11        in[g_index - RADIUS];
12      temp[s_index + BLOCK_SIZE] = g_index + BLOCK_SIZE < n ?
13        in[g_index + BLOCK_SIZE] : 0;
14    }
15
16    // Apply the stencil
17    int result = 0;
18    for (int offset = -RADIUS ; offset <= RADIUS ; offset++)
19      result += temp[s_index + offset];
20
21    // Store the result
22    out[g_index] = result;
23  }
24 }
```



# Race condition



- The stencil example will not work...
- A race condition occurs when multiple tasks read from and write to the same memory without proper synchronization.
- The “race” may finish correctly sometimes and therefore complete without errors, and at other times it may finish incorrectly.
- If a data race occurs, the behavior of the program is undefined.

# \_\_syncthreads ()



```
void __syncthreads ();
```

Synchronizes all threads within a block

- Ensuring correct execution order when threads share data.
- Used to prevent race conditions

All threads must reach the barrier

- In conditional code, the condition must be uniform across the block

# Stencil Kernel, fixed



```
1 __global__ void stencil_1d(const int *in, int *out, int n) {
2   __shared__ int temp[BLOCK_SIZE + 2 * RADIUS];
3   auto g_index = threadIdx.x + blockIdx.x * blockDim.x;
4   if (g_index < n) {
5     auto s_index = threadIdx.x + RADIUS;
6
7     // Read input elements into shared memory
8     temp[s_index] = in[g_index];
9     if (threadIdx.x < RADIUS) {
10      temp[s_index - RADIUS] = g_index - RADIUS < 0? 0 :
11        in[g_index - RADIUS];
12      temp[s_index + BLOCK_SIZE] = g_index + BLOCK_SIZE < n ?
13        in[g_index + BLOCK_SIZE] : 0;
14    }
15    __syncthreads();
16
17    // Apply the stencil
18    int result = 0;
19    for (int offset = -RADIUS ; offset <= RADIUS ; offset++)
20      result += temp[s_index + offset];
21
22    // Store the result
23    out[g_index] = result;
24  }
25 }
```



# Atomic operations



- When we need to modify a variable that is shared among many threads we use atomic operations
- All atomics take as input the address of the shared variable and the value needed for the operation
  - Ex. `atomicAdd(&data, increment)`
- Atomics grant access to the shared variable to only one thread at a time
  - Hence the name, they are indivisible
- There are many different atomic operations
  - `atomicInc`, `atomicDec`, `atomicMax`, `atomicMin`, `atomicAdd`, ...
- If we want block-level synchronization, we call the atomic with the `_block` suffix
  - Ex. `atomicAdd_block`

## Launching parallel threads

- Launch  $N$  blocks with  $M$  threads per block with `kernel<<<N, M, 0, stream>>> (...);`
- Use `blockIdx.x` to access block index within grid
- Use `threadIdx.x` to access thread index within block

## Allocate elements to threads:

```
auto 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 to prevent data hazards



# Device Management

# CUDA Runtime system

- Threads assigned to execution resources on a block-by-block basis.
- CUDA runtime automatically reduces number of blocks assigned to each SM until resource usage is under limit.
- Runtime system:
  - maintains a list of blocks that need to execute
  - assigns new blocks to SM as they compute previously assigned blocks
- Example of SM resources:
  - threads/block or threads/SM or blocks/SM
  - number of threads that can be simultaneously tracked and scheduled
  - shared memory

# Context Switching

- Registers and shared memory are allocated for a block as long as that block is active
- Once a block is active it will stay active until all threads in that block have completed
- Context switching is very fast because registers and shared memory do not need to be saved and restored
- Goal: Have enough transactions in flight to saturate the memory bus
- Latency can be hidden by having more transactions in flight
- Increase active threads or Instruction Level Parallelism

# Maximizing asynchronous operations



- If I call the synchronous API functions, the host will needlessly wait for the GPU to finish executing
- That time could be used to launch other operations on the CPU
- In general, always prefer asynchronous API functions
  - `cudaMemcpyAsync`, `cudaMallocAsync`, `cudaFreeAsync`
- Only synchronize when needed

# Pinned memory



- Pinned memory is a main memory area that is not pageable by the operating system
- Ensures faster transfers (the DMA engine can work without CPU intervention)
- The only way to get closer to PCI peak bandwidth
- Needed in order for CUDA asynchronous operations to work correctly

```
1 // allocate pinned memory
2 cudaMallocHost(&area, sizeof(double) * N);
3 // free pinned memory
4 cudaFreeHost(area);
```

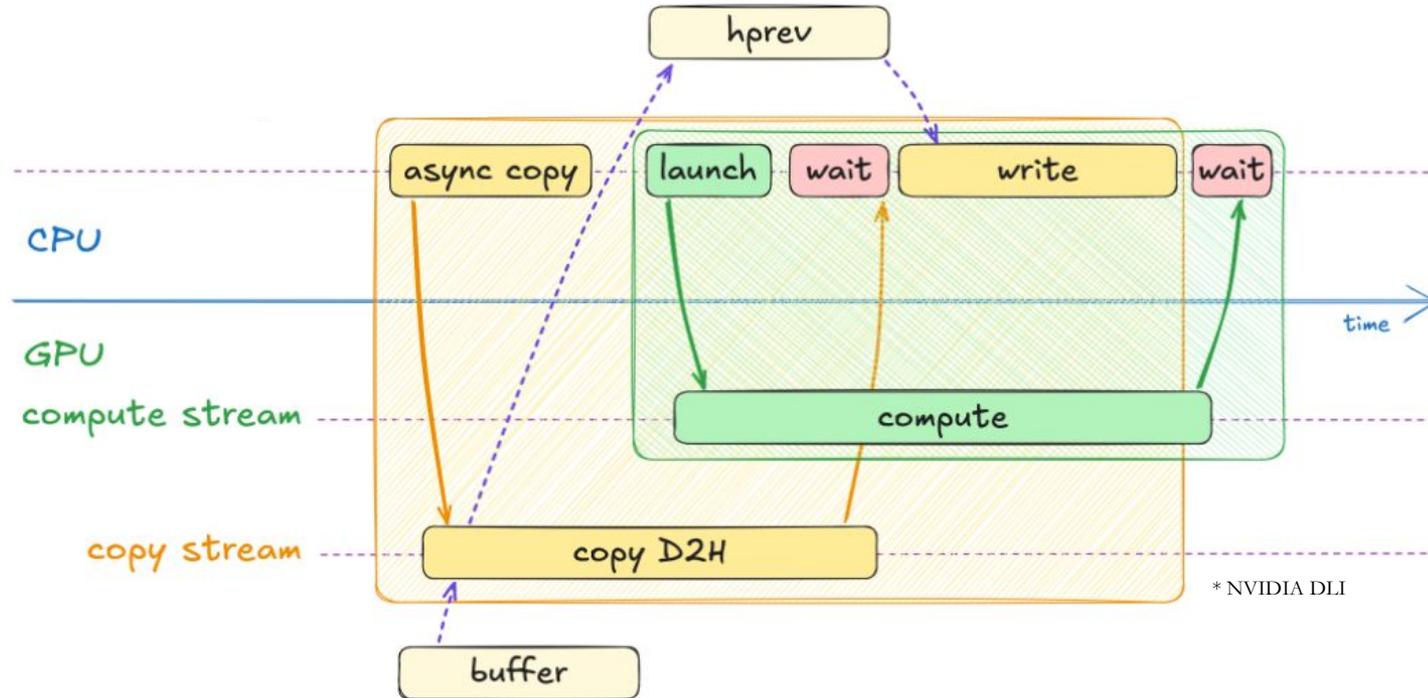
# CUDA streams enable concurrency



- Simultaneous support:
  - CUDA kernels on GPU
  - 2 `cudaMemcpyAsync` (in opposite directions)
  - Computation on the CPU
- Requirements for Concurrency:
  - CUDA operations must be in different, non-0, streams
  - `cudaMemcpyAsync` with host from 'pinned' memory

# Concurrency for overlapping independent operations

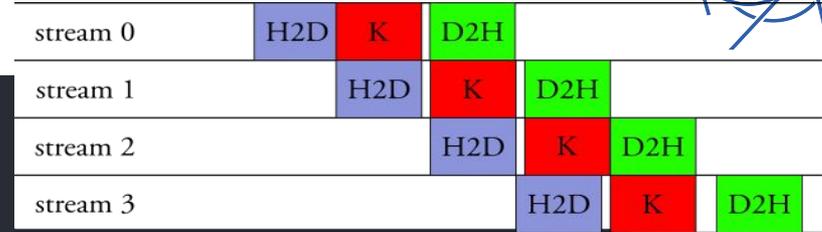
- Using multiple streams allows to execute concurrently pieces of code that are independent
- This allows to maximize the use of the machine



# CUDA Streams



```
1 std::vector<cudaStream_t> streams(4);
2 // create a set of streams
3 for (auto& s: streams)
4   cudaStreamCreate(&s);
5 // allocate data buffers on the host
6 std::vector<float*> hPtrs(4); std::vector<float*> dPtrs(4);
7 for (int i=0; i<4; ++i)
8   cudaMallocHost(&hPtrs[i], memSize);
9 // allocate on device, copy data and launch kernels
10 for (int i=0; i<4; ++i) {
11   cudaMallocAsync(&dPtrs[i], memSize, streams[i]);
12   cudaMemcpyAsync(dPtrs[i], hPtrs[i], memSize, cudaMemcpyHostToDevice, streams[i]);
13   kernelA<<<100,512,0,streams[i]>>>(dPtrs[i]);
14   kernelB<<<100,512,0,streams[i]>>>(dPtrs[i]);
15   cudaMemcpyAsync(hResults[i], dPtrs[i], memSize, cudaMemcpyDeviceToHost, streams[i]);
16 }
17 // synchronize and destroy the streams
18 for (auto& s: streams) {
19   cudaStreamSynchronize(s);
20   cudaStreamDestroy(s); // if the stream is not needed any longer
21 }
```



# 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)
```

```
cudaGetErrorString(cudaGetLastError());
```

# Timing



- You can use the standard timing facilities (host side) in an almost standard way...  
**...but remember that in general you want to operate on the GPU as asynchronously as possible**

# Performance portability

- Started effort to make CMS online and offline event reconstruction heterogeneous in 2016
- Ability to write code that can target different hardware (NVIDIA, AMD, Intel GPUs, CPUs).
- Reduces vendor lock-in and expands deployment options.
  - While keeping more than an eye on SYCL, we ported our CUDA code to [alpsaka](#) portability library
- Fortunately GPUs all work in very similar ways and once you learn one programming model and know how/if to map logical names to the hardware you can program any GPU
  - <https://github.com/CHIP-SPV/chipStar>
  - <https://github.com/ROCm/HIPIFY>



HOW STANDARDS PROLIFERATE:  
(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC.)



# Emphasizing Parallel Algorithm Design



- Algorithm Adaptability
  - Focusing on data partitioning strategies that work across models.
  - Designing algorithms that minimize inter-thread communication and maximize concurrency.
- Optimization Techniques
  - Memory coalescing, minimizing divergence, and efficient use of shared memory are universal concerns.
  - Profiling and tuning performance using model-specific tools.

# CUDA in Python with CuPy



- CuPy is the GPU-accelerated version of NumPy / SciPy
  - Reuses the same APIs
  - Computation executed on NVIDIA GPUs using CUDA
    - 🔥 experimental support for AMD GPUs
  - Support for ndarray operations, universal functions, advanced indexing
  - Support for FFT, sparse matrices, linear algebra (through cuBLAS, cuSOLVER)
- Getting started is as simple as:

```
import numpy as np  
myarray = np.arange(10)
```



```
import cupy as cp  
myarray = cp.arange(10)
```

# Creation and transfer of arrays



- Arrays created on the CPU can be copied to the GPU ..
- .. and copied back to the CPU
- They can be created directly on the GPU avoiding copies if necessary
- If accessed (e.g. printed) they are copied to the CPU automatically

```
arr_cpu = np.arange(100,  
                   dtype=np.float32)  
on_gpu = cp.asarray(arr_cpu)  
to_cpu = cp.asnumpy(on_gpu)
```

```
arr1 = cp.arange(100,  
                dtype=cp.float32)  
arr2 = cp.ones(100)  
arr3 = cp.zeros_like(arr2)
```

NOTE: GPU operations in CuPy are asynchronous by default, meaning that when a function call returns immediately, while the operation may still be running on the GPU. To measure execution time accurately, you must synchronize the GPU before stopping your timer to ensure that all GPU computations have finished. This can be done with: `cp.cuda.Device().synchronize()`

# Conclusion



- Programming GPUs forces you to think parallel
  - CUDA is very well mapped to the properties of the hardware
- Portable code is key for long-term maintainability, testability and support for new accelerator devices
  - It improves the CPU performance as well if done properly, aiding automatic vectorization
  - Many possible solutions, not so many viable ones, even less production ready or compatible with existing infrastructure
- Starting from a CUDA code rather than sequential C++ made our life so much easier in our portability endeavour

# Where to go from here?



- References for learning C++
  - *Professional C++* by M. Gregoire
  - *High-performance C++* by B. Andrist and V. Sehr
  - *Template metaprogramming* by M. Bancila
- References for mastering GPU computing
  - *Programming massively parallel processors* by W-M. Hwu, D. B. Kirk
  - *Data-oriented design* by R. Fabian
  - *NVIDIA CUDA programming guide*