CUDA Course

Victor Eijkhout

2024 COE 379L / CSE 392





Justification

CUDA is the basic programming model for NVidia GPUs ... which TACC is getting a lot of.





Introduction





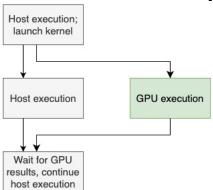
Why GPUs

- Massive parallelism (but that's really multicore and SIMD)
- Simpler processor: less power / more computing for the same power.
- Simpler processor: no general programs, only 'data parallel' operations
 - ... such as machine learning
- Latest GPUs have very high speed on low precision
 - ... such as for machine learning





Attached processor



- GPU is attached to regular CPU
- computation is 'offloaded'
- GPU has its own memory: you may need to synchronize





Data parallel programming

- Kernel: operation for a single thread; in effect: body of inner loop
- Grid: multi-dimensional structure of available 'threads' executed on cores (not a CPU core, more SIMD vector lane)





Processor structure

• Streaming Multiprocessor: like a core

• Warp: vector instruction

• Thread: like a simd-lane

 Single Instruction Multiple Thread (SIMT) 'massive multi-threading': you can have way more 'threads' than blocks × blocksize processor can switch very fast between threads





Loops vs CUDA





GPU memory

- Very much graphics inspired: 'constant' memory, 'texture' memory
- Main memory: has grown from 8G to 100G-ish
- Shared memory inside each SM: managed cache
- Actual caches and registers.





Comparison

- Dozens of cores
- Core has vector instructions
- · a vector instruction has 'lanes'
- Latency hiding by caches, prefetching

- Dozens of 'Streaming Multiprocessors'
- SM has 'warp's
- a warp has threads
- Latency hiding by massive multi-threading



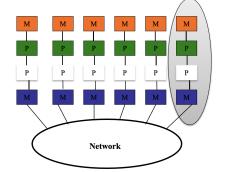


GPU cluster

Memory GPU

GPU

CPU
Memory CPU



Step 1: One node

Step 2: Multiple nodes

- Use GPU on a node (possible combined with OpenMP or other threading)
- Use MPI between nodes;
 NVlink also exists but not on large scale





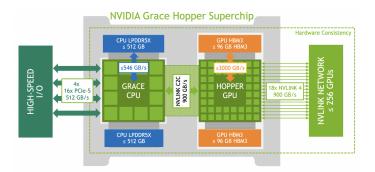
TACC's Vista

- 256 Grace-Grace two-socket CPU nodes 4 racks of 64 nodes each
- Grace-grace node has twice 72 cores: 144 total per node.
- 600 Grace-Hopper nodes
 19 racks of 32 nodes
- GH200 Grace-Hopper superchip combines each containing one 72-core Grace CPU (120 GiB LPDDR5X) and H200 GPU (96 GiB HBM3)
- GH node: 34 Tflops of FP64
 1979 Tflops of FP16 for Machine Learning.





Grace Hopper superchip







CPU vs GPU compared Hardware Comparison

(Frontera and Lonestar)

	Dual socket Intel Cascade Lake	Dual socket AMD Milan	Nvidia H100	Nvidia A100	Nvidia A30
Core/SM	56	128	144	108	56?
Speed (GHz)	2.7	2.45	1.6	1.41	1.24
SIMD / SIMT width (32 bit)	16	8	32	32	32
Vector lanes/ CUDA cores	896*2	1024*2	16896	6912	3584
	FMA	FMA			





CUDA and others

- CUDA is for C and Fortran; limited to NVidia
- OpenACC, OpenCL: portable but tough to program
- AMD GPUs: 'HIP'
- Intel GPUs (hm): Sycl
- OpenMP offload: universal but maybe not highest performance
- Kokkos: portable to OpenMP and all GPUs





Kernels





Translation code to kernel

```
1 for ( int i=0; i<n; ++i ) {
2   int threadIdx.x = i;
3   f(i);
4 }
5   6 __global__ computef( /* fargs */ ) {
7   int i = threadIdx.x;
8   f(i);
9 }
10   11 int gridsize,blocksize;
12 computef<<<gri>i1 computef
( fargs );
```

- Kernel is the body of the inner loop;
- Kernel is executed for each thread in the grid;
- Kernel consults treadIdx to find the 'iteration' number.





About threads

- CUDA creates thread blocks until all data has been processed
- You can tune the block size up to some limit;
- Make sure to create enough blocks.





Thread indexing

```
1 int global_thread_id = blockIdx.x * blockDim.x + threadIdx.x; // or 2d, 3d variant
2 if ( global_thread_id >= datasize ) return;
```

- Global thread index from block and thread coordinate;
- Global index often used as index in the data;
- You can have more threads than data; no problem but you do need to test.





Data on the device

```
1 float
2 *x, *y, *z;
3 size_t nbytes = N*sizeof(float);
4 CU_ASSERT( cudaMallocManaged( &x,nbytes ) );
5 CU_ASSERT( cudaMallocManaged( &y,nbytes ) );
6 CU_ASSERT( cudaMallocManaged( &z,nbytes ) );

1 // Run kernel on 1M elements on the CPU
2 int blocksize = 1024;
3 int gridsize = N/blocksize+1;
4 add1<<<gri>4 cyclisize, blocksize>>>(N,x,y,z);
5 // Wait for GPU to finish before accessing data on host
6 CU_ASSERT( cudaDeviceSynchronize() );
```

- Managed allocation is accessible on host and device
- Device is asynchronous, so synchronization needed





Utility macro

```
1 #define CU_ASSERT( code ) {
2    cudaError_t err = code ; \
3    if (err!=cudaSuccess) { \
4       printf("error <<%s>> on line %d\n",cudaGetErrorString(err),__LINE__); \
5       return 1; } }
```

- CUDA routines can return error code
- Check on fatal errors.





Exercise 1

Write a CUDA kernel that adds two arrays into a third.

Using the above snippets, finish the code, test for correctness.





Two-dimensional grid

```
1 dim3 blocksize(16,16);
2 dim3 gridsize( rowsize/blocksize.x+1,colsize/blocksize.y+1 );
3 printf(" .. grid is (%d,%d) blocks of (%d,%d)\n",
4 gridsize.x,gridsize.y,blocksize.x,blocksize.y);
```

- Two-dimensional data calls for two-dimensional thread organization (same for three of course)
- Specify grid/block size with dim3 objects.





Two-dimensional indexing

```
1 // index2d.cu
2 // Calculate the row and column indices for this thread
3 size_t row = blockIdx.y * blockDim.y + threadIdx.y;
4 size_t col = blockIdx.x * blockDim.x + threadIdx.x;
5
6 // Make sure we don't go out of bounds
7 if (row < height && col < width) {
8 // Calculate the linear index for the matrices
9 size_t idx = row * width + col;
10
11 // Perform the addition
12 d c[idx] = d a[idx] + d b[idx];</pre>
```

- Compute row/col from x/y components
- Adjacent threads need to access adjacent data



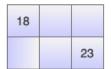


Thread numbering

0	1	2
3	4	5

6	7	
		11

12	
	17



- Lexicographic ordering of blocks;
- Lexicographic ordering of threads in blocks.





Exercise 2

Create a two-dimensional array and let each thread write its global number in the elemenent where it is active:

```
1 output[output_loc] = global_thread_num;
```

Your program should recreate the layout of figure ??.





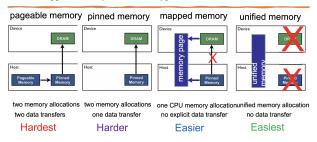
Memory





Host and device memory

Grace/Hopper Memory Allocation Types





72

- Host memory can be copied;
- can be visible on the device:
- can be unified with the device.





Explicit transfer

Allocate input data and copy host values:

```
1 // input allocated on device
2 int *dev_input;
3 cudaMalloc( &dev_input,bytesize );
4 // host values copied
5 cudaMemcpy( dev_input,input,bytesize,cudaMemcpyHostToDevice);
```

After the kernel execution copy output data back:

```
1 cudaDeviceSynchronize();
2 cudaMemcpy( output, dev_output, bytesize, cudaMemcpyDeviceToHost );
```

Free:

```
1 free(input); free(output);
2 cudaFree(dev_input); cudaFree(dev_output);
```

- Host holds pointer to device memory;
- used for copying in/out;
- device memory freed.





Implicit transfer

```
1 // addlm.cu
2 size_t nbytes = N*sizeof(float);
3 CU_ASSERT( cudaMallocManaged( &x,nbytes ) );
4 CU_ASSERT( cudaMallocManaged( &y,nbytes ) );
5 CU_ASSERT( cudaMallocManaged( &z,nbytes ) );
```

- Possibility one: allocate on the host with cudaHostAlloc using a flag of cudaHostAllocMapped.
 Pass cudaHostGetDevicePointer to the device.
- Easier: cudaMallocManaged; transfer is done through on-demand paging.
- (Impliciations for timing!)





Synchronization





Mechanisms

- Threads in a warp are always synchronized.
- Threads in a block can be synchronized with __syncthreads.
- Blocks in a grid can not be synchronized.
- Host and device can be synchronized with cudaDeviceSynchronize.





Warps and warp divergence

```
1 if ( threadIdx.%2=0 )
2 functionA();
3 else
4 functionB();
1 bool mask = (threadIdx.x%2=0)
2 if (mask) functionA();
3 if (not mask) functionB();
```

- Threads in a warp are synchronized
- Conditionals are serialized: warp divergence
- Try to rewrite your code...





Synchronization across blocks

- Threads blocks are not synchronized; can not be.
 Does this make sense from the original purpose of Graphics Processing Units (GPUs)?
- Solution 1: synchronize outside CUDA;
- Solution 2: use atomics;
- Solution 3: use a library that does it for you.





Example: sum reduction

• Recursive summation in each block:

```
1 // redumy.cu
2 __global__
3 void reduce0(size_t n, float *input, float *y);
```

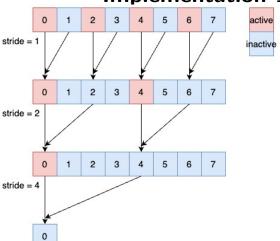
• Summing the blocks outside of CUDA:

```
1 for ( int b=0; b<gridsize; ++b)
2 value += y[b*blocksize];</pre>
```





Implementation 1



We let threads at distances $2, 4, 8, \ldots$ sum values from offsets $1, 2, 4, \ldots$





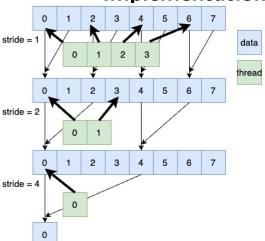
Code

```
1 __global__
2 void reduce0(size_t n, float *input, float *y) {
    size t
      thrd = threadIdx.x,
      global = blockDim.x*blockIdx.x + threadIdx.x;
    if (global>=n) return;
    y[global] = input[global];
    __syncthreads();
   for ( size_t offset=1; offset<=blockDim.x/2; offset*=2 ) {</pre>
      if (thrd\%(2*offset)==0)
10
        y[global] += y[global+offset];
11
      __syncthreads();
12
13
14 }
```





Implementation 2



ullet Non-obvious mapping thread \leftrightarrow data





Code

```
1 // pointer to this block
2 float *data = y*blockDim.x*blockIdx.x;
3 for ( size_t offset=1; offset<=blockDim.x/2; offset*=2 ) {
4    size_t index = 2 * offset * thrd;
5    if (index+offset<blockDim.x)
6    data[index] += data[index+offset];
7    __syncthreads();
8 }
9 // we need the following statement really only for thread 0:
10 y[global] = data[thrd];</pre>
```





Implementation 3

- Use 'Cuda UnBound'
- https://github.com/NVIDIA/cccl/tree/main/cub





Example: stencil



