

# CME213/ME339

## Lecture 3

Eric Darve  
Erich Elsen  
Austin Gibbons

Department of Mechanical Engineering  
Institute for Computational and Mathematical Engineering  
Stanford University

Spring 2012



# A simple example

---

```
1 // Host code
2 int main()
3 {
4     size_t size = ...;
5     float *h_A, *h_B, *h_C, *d_A, *d_B, *d_C;
6     // Allocate input vectors h_A and h_B in host memory
7     h_A = (float*)malloc(size);
8     h_B = (float*)malloc(size);
9     h_C = (float*)malloc(size);
10
11     // Initialize input vectors
12     ...
13
14     // Allocate vectors in device memory
15     cudaMalloc(&d_A, size);
16     cudaMalloc(&d_B, size);
17     cudaMalloc(&d_C, size);
18     ...
```



```
1      ...
2      // Copy vectors from host memory to device memory
3      cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
4      cudaMemcpy(d_B, h_B, size, cudaMemcpyHostToDevice);
5
6      // Invoke kernel
7      int nthreads = 256;
8      int nblocks = (N + nthreads - 1) / nthreads;
9      VecAdd<<<nblocks, nthreads>>>(d_A, d_B, d_C, N);
10
11     // Copy result from device memory to host memory
12     // h_C contains the result in host memory
13     cudaMemcpy(h_C, d_C, size, cudaMemcpyDeviceToHost);
14     ...
15 }
```



```
1  // Device code
2  __global__ void VecAdd(const float* A,
3                        const float* B,
4                        float* C,
5                        const int N)
6  {
7      int i = blockDim.x * blockIdx.x + threadIdx.x;
8      if (i < N)
9          C[i] = A[i] + B[i];
10 }
```



# A few keywords will get you going

---

- Kernel: function that executes on device (GPU) and can be called from host (CPU)
  - Can only access GPU memory
  - Not recursive
  - Must have void return type
  - No static variables
  - No variable number of arguments
- Functions must be declared with a qualifier:
  - `__global__`: GPU kernel function launched by CPU
  - `__device__`: GPU kernel called from GPU
  - `__host__`: CPU function called from CPU (default)
  - `__host__` and `__device__` qualifiers can be combined
- These qualifiers determine how functions are compiled.



# Memory allocation

---

```
1  cudaMalloc(void **pointer, size_t nbytes);  
2  cudaMemset(void *pointer, int value, size_t count);  
3  cudaFree(void *pointer);
```

Example:

```
1  int n = 1024;  
2  int nbytes = 1024*sizeof(int);  
3  int *d_a = 0;  
4  cudaMalloc(&d_a, nbytes);  
5  cudaMemset(d_a, 0, nbytes);  
6  cudaFree(d_a);
```



# Copying data from anywhere to anywhere

---

```
1 cudaMemcpy(void *dst, void *src, size_t nbytes,  
2           enum cudaMemcpyKind direction);
```

- direction specifies locations (host or device) of src and dst
- Blocks CPU thread: returns only after the copy is complete
- Doesn't start copying until previous CUDA calls complete
- enum cudaMemcpyKind:
  - cudaMemcpyHostToDevice
  - cudaMemcpyDeviceToHost
  - cudaMemcpyDeviceToDevice



# Three key abstractions

---

- ① hierarchy of thread groups
- ② shared memories
- ③ barrier synchronization





# CUDA threads

---

Kernel: function executed on the GPU as an array of parallel threads

```
1 // Device code
2 __global__ void VecAdd(...) { ... }
```

All threads execute the same kernel code, but can take different paths.

Each thread has an ID. It can be used to:

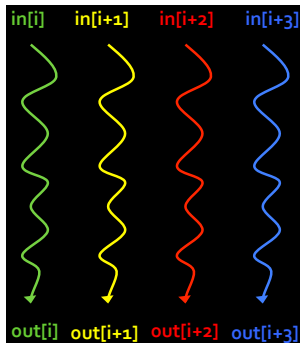
- Select input/output data
- Control decisions



# CUDA threads

ID is used to:

- Select input/output data
- Control decisions



```
1 // Device code
2 __global__ void VecAdd(const float* A, const float* B,
3                       float* C, const int N) {
4     int i = blockDim.x * blockIdx.x + threadIdx.x;
5     if (i < N) C[i] = A[i] + B[i];
6 }
```



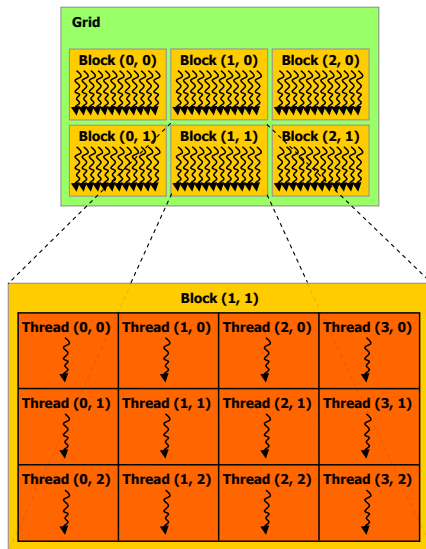
# Many many threads

---

- **CPU threads:** heavyweight entities.
- The operating system must swap threads on and off CPU execution channels to provide multithreading capability.
- Context switches (when two threads are swapped) are therefore slow and expensive.
- **GPUs threads:** numerous and lightweight.
- Critical to achieve peak performance: having many threads help hide memory (and instruction) latency.



# Blocks of threads and grid of blocks



# Thread index

---

threadIdx is a 3-component vector, so that threads can be identified using a one-dimensional, two-dimensional, or three-dimensional thread index, forming a 1D, 2D or 3D thread block.

```
1  __global__ void MatAdd(const float* A, const float* B,  
2                          float* C, const int N) {  
3      int i = threadIdx.x; int j = threadIdx.y;  
4      if (i < N && j < N) C[i*N+j] = A[i*N+j] + B[i*N+j];  
5  }
```



# Grid

---

Blocks are organized into a 1D, 2D or 3D grid of thread blocks.

blockIdx: 1D, 2D or 3D index to identify a thread block.

```
1  __global__ void MatAdd(const float* A, const float* B,  
2                          float* C, const int N) {  
3      int i = blockIdx.x * blockDim.x + threadIdx.x;  
4      int j = blockIdx.y * blockDim.y + threadIdx.y;  
5      if (i < N && j < N) C[i*N+j] = A[i*N+j] + B[i*N+j];  
6  }
```

blockDim: dimension of the thread block

gridDim: dimension of the grid



```
1  int main() {  
2  ...  
3  // Kernel invocation  
4  dim3 nthreads(16, 16);  
5  dim3 nblocks(N/nthreads.x, N/nthreads.y);  
6  MatAdd<<<nblocks, nthreads>>>(A, B, C);  
7  ...  
8  }
```

The number of threads per block and the number of blocks per grid are specified in the `<<<...>>>` syntax.

It can be of type `int` or `dim3`.

Note:

```
1  int nblocks = N / nthreads;  
2  int nblocks = (N + nthreads - 1) / nthreads;
```



# Example of indexing

```
1  __global__ void kernel( int *a ) {  
2      int idx = blockIdx.x*blockDim.x + threadIdx.x;  
3      a[idx] = 7;  
4  }
```

Output: 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7

```
1      int idx = blockIdx.x*blockDim.x + threadIdx.x;  
2      a[idx] = blockIdx.x;
```

Output: 0 0 0 0 1 1 1 1 2 2 2 2 3 3 3 3

```
1      int idx = blockIdx.x*blockDim.x + threadIdx.x;  
2      a[idx] = threadIdx.x;
```

Output: 0 1 2 3 0 1 2 3 0 1 2 3 0 1 2 3





# Limits

---

There are many limits on various quantities due to hardware limitations. We will provide a more complete list later on.

Max x-, y-, or z-dimension of a grid of thread blocks	65535
Max x- or y-dimension of a block	1024
Max z-dimension of a block	64
Max # threads per block	1024

