

# GPGPU programming

Presented by Lukas Gamper and Alex Kosenkov

# Agenda

- CUDA
- (OpenCL)
- thrust
- OpenACC

# CUDA

- Compute Device Unified Architecture
- General-Purpose computing on Graphics Processing Units

# CUDA: \$ module load cuda

## Makefile

```
1 objects = main.o vectoradd.o
2
3 all: $(objects)
4     nvcc -arch=sm_20 $(objects) -o app
5
6 %.o: %.cpp
7     nvcc -x cu -arch=sm_20 -I. -dc $< -o $@
8
9 clean:
10    rm -f *.o app
```

```
$ make
$ nvcc -x cu -arch=sm_20 -I. -dc main.cpp -o main.o
$ nvcc -x cu -arch=sm_20 -I. -dc vectoradd.cpp -o vectoradd.o
$ nvcc -arch=sm_20 main.o vectoradd.o -o app
```

## vectoradd.h

```
1 #ifndef VECTORADD_H
2 #define VECTORADD_H
3
4 #include <cuda_runtime.h>
5
6 __global__ void vectorAdd(const float *A, const float *B, float *C, int n);
7 __device__ void vectorAdd_impl(const float& A, const float& B, float& C);
8
9 #endif
```

## vectoradd.cpp

```
1 #include "vectoradd.h"
2
3 __global__ void vectorAdd(const float *A, const float *B, float *C, int n){
4     int i = blockDim.x * blockIdx.x + threadIdx.x;
5     if(i < n) vectorAdd_impl(A[i], B[i], C[i]);
6 }
7
8 __device__ void vectorAdd_impl(const float& A, const float& B, float& C){
9     C = A + B;
10 }
```

# CUDA: \$ module load cuda

main.cpp

```
1 #include <stdio.h>
2 #include "vectoradd.h"
3 #define N 50000
4
5 int main(){
6     size_t size = N*sizeof(float);
7     float *h_A = (float *)malloc(size);
8     float *h_B = (float *)malloc(size);
9     float *h_C = (float *)malloc(size);
10    float *d_A = NULL; cudaMalloc((void **)&d_A, size);
11    float *d_B = NULL; cudaMalloc((void **)&d_B, size);
12    float *d_C = NULL; cudaMalloc((void **)&d_C, size);
13
14    for(int i = 0; i < N; ++i){
15        h_A[i] = rand()/(float)RAND_MAX;
16        h_B[i] = rand()/(float)RAND_MAX;
17    }
18    cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
19    cudaMemcpy(d_B, h_B, size, cudaMemcpyHostToDevice);
20
21    int threadsPerBlock = 256;
22    int blocksPerGrid = (N + threadsPerBlock - 1) / threadsPerBlock;
23    printf("CUDA kernel launch with %d blocks of %d threads\n", blocksPerGrid, threadsPerBlock);
24    vectorAdd <<< blocksPerGrid, threadsPerBlock >>> (d_A, d_B, d_C, N);
25
26    cudaMemcpy(h_C, d_C, size, cudaMemcpyDeviceToHost);
27
28    free(h_A); free(h_B); free(h_C);
29    cudaFree(d_A); cudaFree(d_B); cudaFree(d_C);
30    cudaDeviceReset();
31    return 0;
32 }
```

# CUDA: \$ module load cuda

Launch:

```
#define N 50000

int threadsPerBlock = 256;
int blocksPerGrid = (N + threadsPerBlock - 1) / threadsPerBlock;

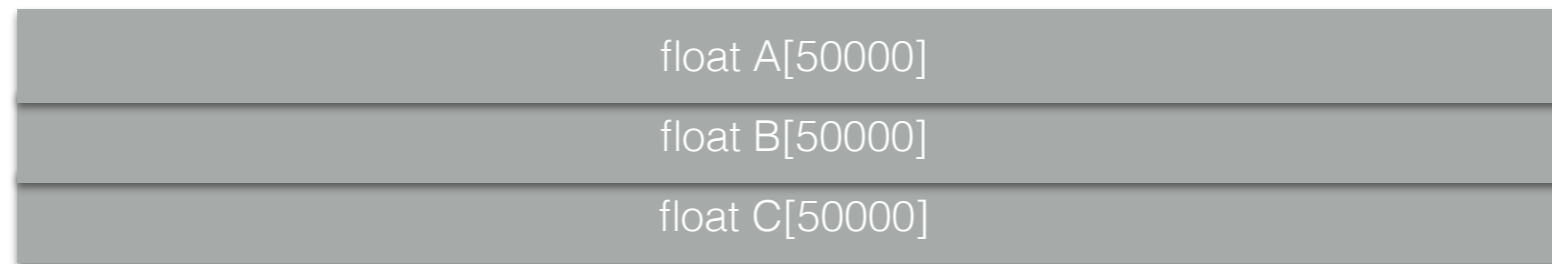
vectorAdd <<< blocksPerGrid, threadsPerBlock >>> (d_A, d_B, d_C, N);
```

Kernel:

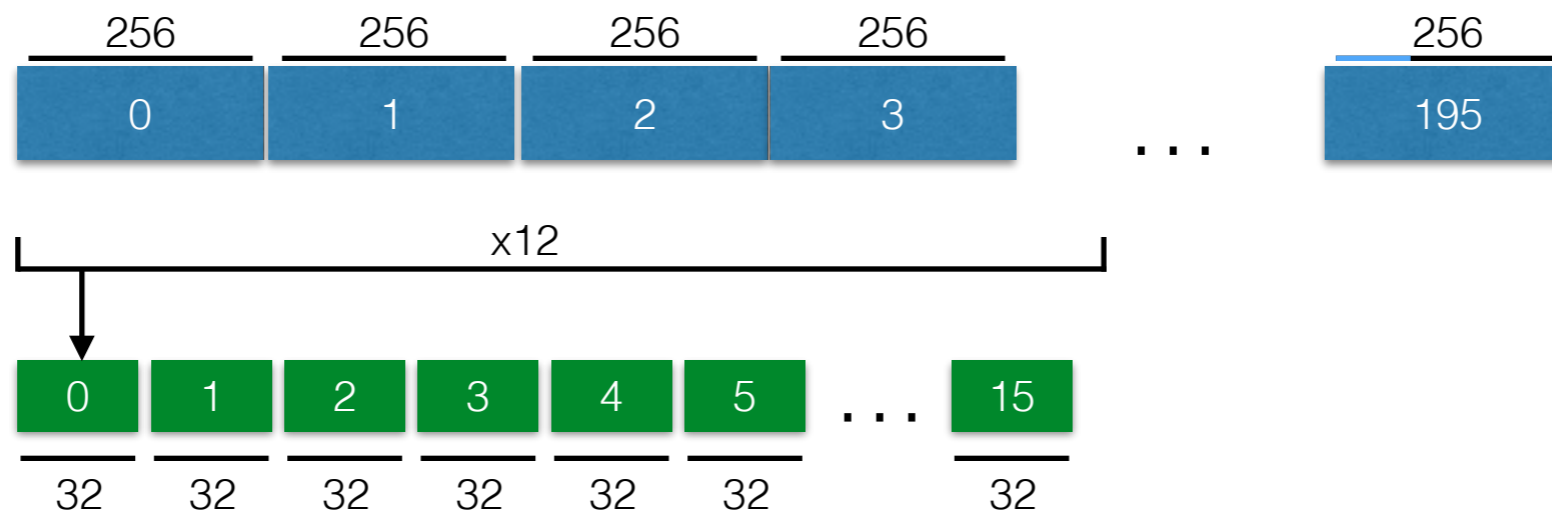
```
__global__ void vectorAdd(const float *A, const float *B, float *C, int n){
    int i = blockDim.x * blockIdx.x + threadIdx.x;
    if(i < n) C[i] = A[i] + B[i];
}
```

CUDA kernel launch with 196 blocks of 256 threads

# CUDA: \$ module load cuda



```
int i = blockDim.x * blockIdx.x + threadIdx.x; // 256*blockId + tid
```



12 blocks x 8 warps per MP

# CUDA: limitations

- Data transfer time
- Memory coalescing
- Execution divergence



# CUDA: parallel reduction (nvidia example)

main.cpp (kernel invocation)

```
6 float hostReduce(float* input, unsigned int n){
7     for(unsigned int stride = 1; stride < n; stride *= 2)
8     for(unsigned int k = stride; k < n; k += stride*2)
9     input[k-stride] += input[k];
10    return input[0];
11 }
...
29 void deviceReduce(int blocks, int block_size, float* input, float* output, unsigned int n){
30     int shmemSize = block_size*sizeof(float);
32     reduce<<< blocks, block_size, shmemSize >>>(input, output, n);
33 }
...
45 int main(){
...
61     int threadsPerBlock = 256;
62     int blocksPerGrid = (N + threadsPerBlock - 1) / threadsPerBlock;
...
70     deviceReduce(blocksPerGrid, threadsPerBlock, d_C, d_A, N);
71     deviceReduce(1, blocksPerGrid, d_A, d_B, blocksPerGrid);
72     cudaMemcpy(h_B, d_B, sizeof(float), cudaMemcpyDeviceToHost);
73
74     std::cout << "Host sum is " << hostReduce(h_C, N) << "\n";
75     std::cout << "Device first sum is " << h_B[0] << "\n";
...
81 }
```

# CUDA: parallel reduction #0

```
3 __global__ void reduce0(float *g_idata, float *g_odata, unsigned int n) {
4     extern __shared__ float sdata[];
5     unsigned int tid = threadIdx.x;
6     unsigned int i = blockIdx.x*blockDim.x + threadIdx.x;
7     sdata[tid] = g_idata[i];
8     __syncthreads();
9     for(unsigned int s=1; s < blockDim.x; s *= 2) {
10         if (tid % (2*s) == 0) {
11             sdata[tid] += sdata[tid + s];
12         }
13         __syncthreads();
14     }
15     if (tid == 0) g_odata[blockIdx.x] = sdata[0];
16 }
```

Tip: execution divergence

# CUDA: parallel reduction #1

```
3 __global__ void reduce1(float *g_idata, float *g_odata, unsigned int n) {
4     extern __shared__ float sdata[];
5     unsigned int tid = threadIdx.x;
6     unsigned int i = blockIdx.x*blockDim.x + threadIdx.x;
7     sdata[tid] = g_idata[i];
8     __syncthreads();
9     for (unsigned int s=1; s < blockDim.x; s *= 2) {
10         int index = 2 * s * tid;
11         if (index < blockDim.x) {
12             sdata[index] += sdata[index + s];
13         }
14         __syncthreads();
15     }
16     if (tid == 0) g_odata[blockIdx.x] = sdata[0];
17 }
```

Tip: memory banks conflicts

# CUDA: parallel reduction #2

```
3 __global__ void reduce2(float *g_idata, float *g_odata, unsigned int n) {
4     extern __shared__ float sdata[];
5     unsigned int tid = threadIdx.x;
6     unsigned int i = blockIdx.x*blockDim.x + threadIdx.x;
7     sdata[tid] = g_idata[i];
8     __syncthreads();
9     for (unsigned int s=blockDim.x/2; s>0; s>>=1) {
10         if (tid < s) {
11             sdata[tid] += sdata[tid + s];
12         }
13         __syncthreads();
14     }
15     if (tid == 0) g_odata[blockIdx.x] = sdata[0];
16 }
```

Tip: idle threads

# CUDA: parallel reduction #3

```
3 __global__ void reduce3(float *g_idata, float *g_odata, unsigned int n) {
4     extern __shared__ float sdata[];
5     unsigned int tid = threadIdx.x;
6     unsigned int i = blockIdx.x*(blockDim.x*2) + threadIdx.x;
7     sdata[tid] = g_idata[i] + g_idata[i+blockDim.x];
8     __syncthreads();
9     for (unsigned int s=blockDim.x/2; s>0; s>>=1) {
10         if (tid < s) {
11             sdata[tid] += sdata[tid + s];
12         }
13         __syncthreads();
14     }
15     if (tid == 0) g_odata[blockIdx.x] = sdata[0];
16 }
```

Tip: useless instructions

# CUDA: parallel reduction #4

```
1
2 __device__ void warpReduce(volatile float* sdata, int tid) {
3     sdata[tid] += sdata[tid + 16];
4     sdata[tid] += sdata[tid + 8];
5     sdata[tid] += sdata[tid + 4];
6     sdata[tid] += sdata[tid + 2];
7     sdata[tid] += sdata[tid + 1];
8 }
9
10 __global__ void reduce4(float *g_idata, float *g_odata, unsigned int n) {
11     extern __shared__ float sdata[];
12     unsigned int tid = threadIdx.x;
13     unsigned int i = blockIdx.x*(blockDim.x*2) + tid;
14     sdata[tid] = g_idata[i] + g_idata[i+blockDim.x];
15     __syncthreads();
16     for (unsigned int s=blockDim.x/2; s>16; s>>=1) {
17         if (tid < s) {
18             sdata[tid] += sdata[tid + s];
19         }
20         __syncthreads();
21     }
22     if (tid < 16) warpReduce(sdata, tid);
23     if (tid == 0) g_odata[blockIdx.x] = sdata[0];
24 }
```

# CUDA: parallel reduction #5

```
1
2 template <unsigned int blockSize>
3 __device__ void warpReduce (volatile float *sdata, unsigned int tid) {
4     if(blockSize >= 64) sdata[tid] += sdata[tid + 32];
5     if(blockSize >= 32) sdata[tid] += sdata[tid + 16];
6     if(blockSize >= 16) sdata[tid] += sdata[tid + 8];
7     if(blockSize >= 8)  sdata[tid] += sdata[tid + 4];
8     if(blockSize >= 4)  sdata[tid] += sdata[tid + 2];
9     if(blockSize >= 2)  sdata[tid] += sdata[tid + 1];
10 }
11 template <unsigned int blockSize>
12 __global__ void reduce5(float *g_idata, float *g_odata, unsigned int n) {
13     extern __shared__ float sdata[];
14     unsigned int tid = threadIdx.x;
15     unsigned int i = blockIdx.x*(blockSize*2) + tid;
16     sdata[tid] = g_idata[i] + g_idata[i+blockDim.x];
17     __syncthreads();
18     if(blockSize >= 512){ if(tid < 256) sdata[tid] += sdata[tid + 256]; __syncthreads(); }
19     if(blockSize >= 256){ if(tid < 128) sdata[tid] += sdata[tid + 128]; __syncthreads(); }
20     if(blockSize >= 128){ if(tid < 64)  sdata[tid] += sdata[tid + 64];  __syncthreads(); }
21     if(tid < 32) warpReduce<blockSize>(sdata, tid);
22     if(tid == 0) g_odata[blockIdx.x] = sdata[0];
23 }
```

# CUDA: parallel reduction #6

```
1
2 template <unsigned int blockSize>
3 __global__ void reduce(float *g_idata, float *g_odata, unsigned int n) {
4     extern __shared__ float sdata[];
5     unsigned int tid = threadIdx.x;
6     unsigned int i = blockIdx.x*(blockSize*2) + tid;
7     unsigned int gridSize = blockSize*2*gridDim.x;
8     sdata[tid] = 0.;
9     while(i < n){ sdata[tid] += g_idata[i] + g_idata[i+blockSize]; i += gridSize; }
10    __syncthreads();
11    if(blockSize >= 512){ if(tid < 256) sdata[tid] += sdata[tid + 256]; __syncthreads(); }
12    if(blockSize >= 256){ if(tid < 128) sdata[tid] += sdata[tid + 128]; __syncthreads(); }
13    if(blockSize >= 128){ if(tid < 64) sdata[tid] += sdata[tid + 64]; __syncthreads(); }
14    if(tid < 32) warpReduce<blockSize>(sdata, tid);
15    if(tid == 0) g_odata[blockIdx.x] = sdata[0];
16 }
```



# CUDA: thrust

## Containers:

`thrust::device_vector< T, Alloc >`

`thrust::host_vector< T, Alloc >`

## Examples of supported algorithms:

`thrust::for_each`

`thrust::fill`, `thrust::generate`, `thrust::transform`

`thrust::sort`, `thrust::find`, `thrust::find_if`

`thrust::copy`, `thrust::copy_if`

`thrust::merge`, `thrust::reduce`, `thrust::inner_product`

# CUDA: thrust

## Example: basic usage

```
1 #include <thrust/host_vector.h>
2 #include <thrust/device_vector.h>
3 #include <thrust/generate.h>
4 #include <thrust/sort.h>
5 #include <thrust/copy.h>
6 #include <algorithm>
7 #include <cstdlib>
8
9 int main(void){
10     // generate 32M random numbers serially
11     thrust::host_vector<int> h_vec(32 << 20);
12     std::generate(h_vec.begin(), h_vec.end(), rand);
13
14     // transfer data to the device
15     thrust::device_vector<int> d_vec = h_vec;
16
17     // sort data on the device
18     thrust::sort(d_vec.begin(), d_vec.end());
19
20     // transfer data back to host
21     thrust::copy(d_vec.begin(), d_vec.end(), h_vec.begin());
22
23     return 0;
24 }
```

# CUDA: thrust

Example: saxpy ( $y[i] += a \cdot x[i]$ )

```
1 #include <thrust/transform.h>
2 #include <thrust/device_vector.h>
3 #include <thrust/host_vector.h>
4 #include <thrust/functional.h>
5 #include <iostream>
6 #include <iterator>
7 #include <algorithm>
8
9 struct saxpy_functor : public thrust::binary_function<float,float,float> {
10     saxpy_functor(float _a) : a(_a) {}
11
12     __host__ __device__ float operator()(const float& x, const float& y) const {
13         return a * x + y;
14     }
15     const float a;
16 };
17
18 void saxpy_fast(float A, thrust::device_vector<float>& X, thrust::device_vector<float>& Y){
19     thrust::transform(X.begin(), X.end(), Y.begin(), Y.begin(), saxpy_functor(A)); // Y <- A * X + Y
20 }
21
22 int main(void){
23     float x[4] = {1.0, 1.0, 1.0, 1.0};
24     float y[4] = {1.0, 2.0, 3.0, 4.0};
25
26     thrust::device_vector<float> X(x, x + 4);
27     thrust::device_vector<float> Y(y, y + 4);
28
29     saxpy_fast(2.0, X, Y);
30     return 0;
31 }
```

# OpenACC

## Example: stencil computation

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <iostream>
4
5 #define M 64
6 #define N 128
7
8 void stencil(float A[M][N], float B[M][N]){
9
10     #pragma acc kernels pcopyin(A[0:M]) pcopy(B[0:M])
11     {
12         float c11, c12, c13, c21, c22, c23, c31, c32, c33;
13         c11 = +2.0f; c21 = +5.0f; c31 = -8.0f;
14         c12 = -3.0f; c22 = +6.0f; c32 = -9.0f;
15         c13 = +4.0f; c23 = +7.0f; c33 = +10.0f;
16
17         #pragma acc loop gang(64)
18         for (int i = 1; i < M - 1; ++i){
19             #pragma acc loop worker(128)
20             for (int j = 1; j < N - 1; ++j){
21                 B[i][j] = c11 * A[i - 1][j - 1] + c12 * A[i + 0][j - 1] + c13 * A[i + 1][j - 1]
22                     + c21 * A[i - 1][j + 0] + c22 * A[i + 0][j + 0] + c23 * A[i + 1][j + 0]
23                     + c31 * A[i - 1][j + 1] + c32 * A[i + 0][j + 1] + c33 * A[i + 1][j + 1];
24             }
25         }
26     }
27 }
28
29 int main(){
30     float A[M][N];
31     float B[M][N];
32
33     for(int i = 0; i < M; i++)
34     for(int j = 0; j < N; j++)
35     A[i][j] = std::rand();
36
37     stencil(A,B);
38     return 0;
39 }
```

# Questions

End of the fourth part