

# Introduction to Parallel and Distributed Processing

## CUDA Programming

Jaroslav 'Jaric' Zola

<http://www.jzola.org/>

## Suggested Reading

- Blog of Mark Harris (NVIDIA Chief Technologist for Software)  
<https://devblogs.nvidia.com/parallelforall/author/mharris/>

# NVIDIA CUDA

- SDK and platform for programming NVIDIA's GPGPUs
- Designed for CUDA-enabled GPGPUs
- Compiler, runtime and supporting tools (e.g. profiler)
- Support for C, C++11 and Fortran

# Basic CUDA Workflow

- Allocate memory on host/device
- Transfer data from host to device
- Configure and run kernel(s)
- Transfer data from device to host
- Free memory

# Basic CUDA Workflow

- Allocate memory on host/device
- Transfer data from host to device
- Configure and run kernel(s)
- Transfer data from device to host
- Free memory
- Since CUDA 6 memory can be managed (i.e. no need to copy)

# Simple CUDA Code

```
1 // ex01.cpp
2 #include <vector>
3 #include <cuda_runtime_api.h>
4
5 int main(int argc, char* argv[]) {
6     int n = 4 * 1024 * 1024;
7     int size = n * sizeof(float);
8
9     std::vector<float> x(n);
10    std::vector<float> y(n);
11
12    float* d_x;
13    float* d_y;
14
15    cudaMalloc(&d_x, size);
16    cudaMalloc(&d_y, size);
17
18    cudaMemcpy(d_x, x.data(), size, cudaMemcpyHostToDevice);
19    cudaMemcpy(d_y, y.data(), size, cudaMemcpyHostToDevice);
20    // ...
21    cudaMemcpy(d_y, y.data(), size, cudaMemcpyDeviceToHost);
22
23    cudaFree(d_x);
24
25    return 0;
26 }
```

# Simpler CUDA Code

```
1 | // ex01.cpp
2 | #include <vector>
3 | #include <cuda_runtime_api.h>
4 |
5 | int main(int argc, char* argv[]) {
6 |     int n = 4 * 1024 * 1024;
7 |     int size = n * sizeof(float);
8 |
9 |     std::vector<float> x(n);
10 |    std::vector<float> y(n);
11 |
12 |    float* d_x;
13 |    float* d_y;
14 |
15 |    cudaMallocManaged(&d_x, size);
16 |    cudaMallocManaged(&d_y, size);
17 |    // ...
18 |    // may need sync: cudaDeviceSynchronize();
19 |    // ...
20 |    cudaFree(d_x);
21 |
22 |    return 0;
23 | }
```

```
1 | nvcc -O3 ex01.cpp -o ex01
2 | nvprof ./ex01
```

# Let's SAXPY

- Just a reminder, the most basic BLAS1 routine :-)

$$z = \alpha \cdot x + y$$



# Let's SAXPY

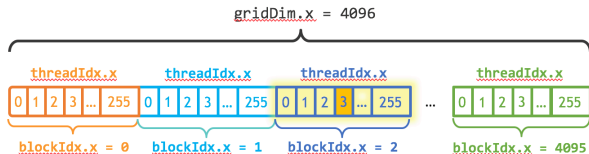
- Idea: we want each CUDA thread to run just one operation:

$$y_i = a \cdot x_i + y_i$$

```

1 | void saxpy(int n, float a, const float* x, float* y) {
2 |     // ...
3 |     y[i] = a * x[i] + y[i];
4 | }
    
```

- $x$  and  $y$  are 1D so it makes sense to organize threads into 1D grid/blocks:



$$\text{index} = \text{blockIdx.x} * \text{blockDim.x} + \text{threadIdx.x}$$

$$\text{index} = (2) * (256) + (3) = 515$$

# Let's SAXPY

```
1 // ex01.cu
2 __global__
3 void saxpy(int n, float a, const float* x, float* y) {
4     int i = blockIdx.x * blockDim.x + threadIdx.x;
5     if (i < n) y[i] = a*x[i] + y[i];
6 }
7
8 __host__
9 void run_saxpy(int n, float a, const float* x, float* y) {
10     const int block_size = 1024;
11     int num_blocks = (n + block_size - 1) / block_size;
12     saxpy<<<num_blocks, block_size>>>(n, 0.1, x, y);
13 }
```

# CUDA Syntax Basics

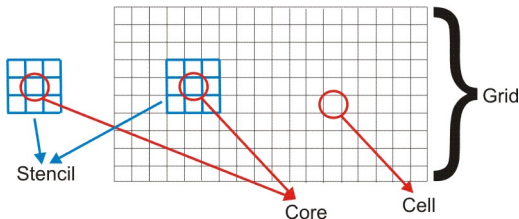
- `__global__` function runs on device but callable from host
- `__device__` runs and callable from device only
- `__host__` function runs on host only
- `fun<<< >>>(...)` calls kernel with grid/block configuration
- `threadIdx`, `blockIdx`, `blockDim` are built-in variables
- Inside block threads indexed always in row-wise manner
- CUDA code must be in `.cu` file

# SAXPY with Thrust

```
1 // ex02.cpp
2 #include <thrust/device_vector.h>
3 #include <thrust/functional.h>
4 #include <thrust/transform.h>
5
6 using thrust::placeholders;
7
8 int main(int argc, char* argv[]) {
9     int n = 4 * 1024 * 1024;
10
11     thrust::host_vector<float> x(n);
12     thrust::host_vector<float> y(n);
13
14     thrust::device_vector<float> d_x = x;
15     thrust::device_vector<float> d_y = y;
16
17     float a = 0.1;
18     thrust::transform(d_x.begin(), d_x.end(),
19                     d_y.begin(), d_y.begin(), a * _1 + _2);
20
21     y = d_y;
22
23     return 0;
24 }
```

# Stencil Computation

- Stencil: pattern in which vector/array is updated

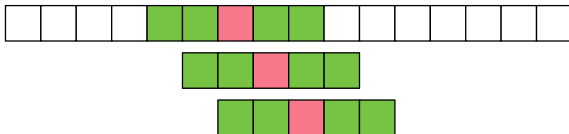


- One of the most basic (parallel) operations  
signal processing, statistics, ML, etc.
- Example 1D stencil:

$$y_i = \frac{1}{2r+1} \sum_{j=i-r}^{i+r} x_j$$

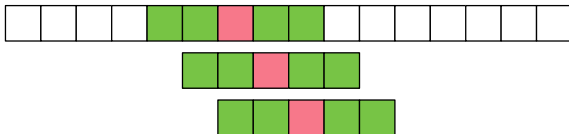
# CUDA Stencil

- Again 1D pattern: each thread computes  $y_i$
- To compute stencil around  $x_i$  we need  $2r$  flanking elements
- Hence, each element will be accessed  $2r + 1$  times



# CUDA Stencil

- Again 1D pattern: each thread computes  $y_i$
- To compute stencil around  $x_i$  we need  $2r$  flanking elements
- Hence, each element will be accessed  $2r + 1$  times



- Idea: make threads to cooperate within their block to “cache” and share/reuse  $x_i$

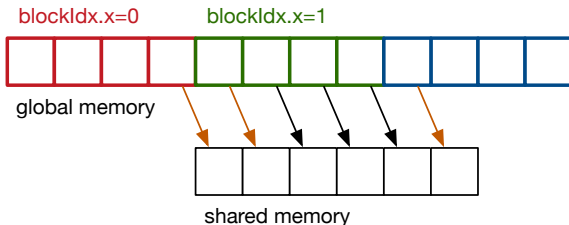
# Shared Memory

- Shared memory: extremely fast on-chip local memory access in 4 cycles
- Shared by threads within block, hence can be used for communication
- Declared with `__shared__`



# Shared Memory

- Shared memory: extremely fast on-chip local memory access in 4 cycles
- Shared by threads within block, hence can be used for communication
- Declared with `__shared__`
- Shared memory in stencil ( $r = 1$ ):



# CUDA Stencil

```
1  template <int block_size, int radius>
2  __global__ void stencil1D(int n, const float* x, float* y) {
3      __shared__ float buf[block_size + 2 * radius];
4
5      int gidx = blockIdx.x * block_size + threadIdx.x;
6      int lidx = threadIdx.x + radius;
7
8      buf[lidx] = x[gidx]; // yey, let's load
9
10     if (threadIdx.x < radius) { // ghost region
11         buf[lidx - radius] = x[gidx - radius];
12         buf[lidx + block_size] = x[gidx + block_size];
13     }
14
15     __syncthreads();
16
17     float S = 0.0;
18     for (int i = -radius; i <= radius; ++i) S += buf[lidx + i];
19     y[gidx] = S / (2 * radius + 1);
20 }
```

# For Fun

- Write host code to invoke stencil and then profile it!