# CUDA: Going Parallel with Threads and Kernels

Ben Cumming, CSCS
July 24, 2016

# Going Parallel : Kernels and Threads

## Threads and kernels

- **threads** are streams of execution, run simultaneously on GPU (1000s)
- **kernel** is the task run by each thread
- CUDA provides language support for
    - writing kernels
    - launching many threads to execute a kernel in parallel
- CUDA hides the low-level details of launching threads

## The process for porting to CUDA

1. formulate algorithm in terms of parallel work items
2. write a kernel implementing a work item on one thread
3. launch the kernel with the required number of threads

CSCS

**ETH** *zürich*

## Scaled Vector Addition (`axpy`)

The exercise in the first section used CUBLAS to perform scaled vector addition

$$\mathbf{y} = \mathbf{y} + \alpha \mathbf{x}$$

- $\mathbf{x}$ and $\mathbf{y}$ are vectors of length $n$ $\qquad x, y \in \mathbb{R}^n$
- $\alpha$ is scalar $\qquad \alpha \in \mathbb{R}$

`axpy` can be performed as $n$ **independent** operations

$$y_i \leftarrow y_i + a * x_i, \quad i = 0, 1, \ldots, n-1$$

which can be performed independently and in any order

### `axpy` implemented with for loop

```
void axpy(double *y, double *x, double a, int n) {
  for(int i=0; i<n; ++i)
    y[i] = y[i] + a*x[i];
}
```

## What is a kernel?

- a kernel defines the work item for a single thread
- the work is performed by many threads executing the same kernel **simultaneously**
- Conceptually corresponds to the inner part of a loop for BLAS1 operations like `axpy`

| host : add two vectors | CUDA : add two vectors |
|---|---|

```
void add_cpu(int *a, int *b, int n){
  for(auto i=0; i<n; ++i)
    a[i] = a[i] + b[i];
}
```

```
__global__
void add_gpu(int *a, int *b, int n){
  auto i = threadIdx.x;
  a[i] = a[i] + b[i];
}
```

- `__global__` keyword indicates a kernel
- `threadIdx` used to find unique id of each thread

**CSCS**

**ETH**zürich

## launching a kernel

- host code launches a kernel on the GPU **asyncronously**
- CUDA provides special `<<<_,_>>>` syntax for launching a kernel
  - `add_gpu<<<1, num_threads>>>(args... )` will launch the kernel `add_gpu` with `num_threads` parallel threads.

### host : add two vectors
```
auto n = 1024;
auto a = host_malloc<int>(n);
auto b = host_malloc<int>(n);
add_cpu(a, b, n);
```

### CUDA : add two vectors
```
auto n = 1024;
auto a = device_malloc<int>(n);
auto b = device_malloc<int>(n);
add_gpu<<<1,n>>>(a, b, n);
```

# Exercise: My First Kernel

Open `topics/cuda/practicals/axpy/axpy_kernel.cu`

1. Write a kernel that implements `axpy` for `double`
   - `axpy_kernel(double *y, double *x, double a, int n)`
   - **extra**: can you write a C++ templated version for any type?
2. launch the kernel (look for `TODO`)
3. Compile the test and run
   - it will pass with no errors on success
   - first try with small vectors of size 8
   - try increasing launch size... what happens?
4. **extra**: can you extend the kernel to work for larger arrays?

**CSCS**

**ETH**zürich

# Scaling Up : Thread Blocks

In the `axpy` exercises we were limited to 1024 threads for a kernel launch

- but we need to scale beyond 1024 threads for the **massive parallelism** we were promised!

### Thread blocks and grids

kernels are executed in groups of threads called **thread blocks**
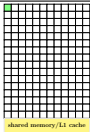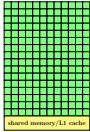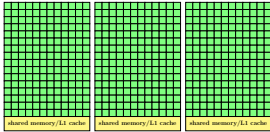
- the launch configuration `axpy<<<grid_dim, block_dim>>>(...)`
  - launch a **grid** of `grid_dim` **blocks**
  - each **block** has `block_dim` **threads**
  - for a total of `grid_dim` $\times$ `block_dim` threads
- previously we launched just one thread block

  `axpy<<<1, n>>>(...)`

CSCS

**ETH**zürich

## Why the additional complexity of grids+blocks+threads?

**Because coordination between threads doesn't scale**

- threads in a block can synchronize and share resources
- this does not scale past a certain number of cores/threads
- on the K20X GPU streaming multiprocessor (SMX) has 192 CUDA cores, and can run 2028 threads
- threads in a block run on the same SMX, with shared resources and thread cooperation
- work is broken into blocks, which are distributed over the 14 SMXs in the K20X GPU

| concept | hardware | |
|---------|----------|---|
| thread |  | ▪ each thread executed on one core |
| block |  | ▪ block executed on 1 SMX<br>▪ multiple blocks per SMX if sufficient resources<br>▪ threads in a block share SMX resources |
| grid |  | ▪ kernel is executed in grid of blocks<br>▪ blocks distributed over SMXs<br>▪ multiple kernels can run at same time |

CSCS

**ETH**zürich

## Calculating thread indexes

A kernel has to calculate the index of its work item

- in `axpy` we used `threadIdx.x` for the index
- when using multiple blocks, we need more information, which is available in the following **magic variables**:

| | |
|---|---|
| `gridDim` | : total number of blocks in the grid |
| `blockDim` | : number of threads in a thread block |
| `blockIdx` | : index of block `[0, gridDim-1]` |
| `threadIdx` | : index of thread in thread block `[0, blockDim-1]` |

**ETH** zürich

## Calculating thread indexes

Consider accessing an array of length 24 with 8 threads per block. The **dimensions** of the kernel launch are:

- `blockDim.x == 8` (8 threads/block)
- `gridDim.x == 3` (3 blocks)

We calculate the index for our thread using the formula

```
auto index = threadIdx.x + blockIdx.x*blockDim.x
```

```
index = threadIdx.x + blockDim.x*blockIdx.x
      = 5 + 8 * 1
      = 13
```

threadIdx.x    | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

blockIdx.x = 0          blockIdx.x = 1          blockIdx.x = 2

## Calculating grid dimensions

The number of thread blocks and the number of threads per block are parameters for the kernel launch:

```
kernel<<<blocks, threads_per_block>>>(...)
```

Remember to guard against overflow when the number of work items is not divisible by the thread block size

## vector addition with multiple blocks

```
__global__
void add_gpu(int *a, int *b, int n){
  auto i = threadIdx.x + blockIdx.x*blockDim.x;
  if(i<n) { // guard against access off end of arrays
    a[i] += b[i];
  }
}

// in main()
auto block_size = 512;
auto num_blocks = (n + (block_size-1)) / block_size;
add_gpu<<<num_blocks, block_size>>>(a, b, n);
```

> ### Calculating grid dimensions
>
> We have to take care when calculating the number of blocks in the grid, i.e. `blocks`:
>
> ```
> kernel<<<blocks, threads_per_block>>>(...)
> ```
>
> Most likely, the number of work items `n` is not a multiple of `threads_per_block`.
>
> - in which case some threads in the last thread block will do any work

> ### Calculating grid dimensions
>
> ```
> // in main()
> auto block_size = 512;
> auto num_blocks = (n + (block_size -1)) / block_size;
> add_gpu<<<num_blocks, block_size>>>(a, b, n);
> ```

CSCS

**ETH**zürich

> **The number of threads per block impacts performance**
>
> - the optimal number depends on the resources (registers, shared memory, etc) that a kernel requires

> **Choosing block size automatically (CUDA 6.5 and later)**
>
> ```
> int block_size, min_grid_dim;
>
> cudaOccupancyMaxPotentialBlockSize(&min_grid_size, &block_size,
>                                    add_gpu, 0, n);
> auto num_blocks = (n + (block_size-1)) / block_size;
>
> add_gpu<<<num_blocks, block_size>>>(a, b, n);
> ```

> The variable `min_grid_size` is set to the minimum number of blocks required to **saturate** the GPU, i.e. provide the GPU with enough work to utilize all of the SMXs.
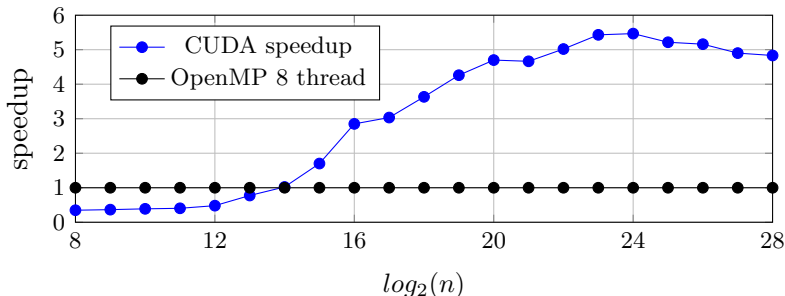
# Exercise: Blocks

Open `cuda/exercises/axpy/axpy.cu` from the last exercise

1. Extend the `axpy` kernel for arbitrarily large input arrays (any `n`)
2. Update the call site to calculate the grid configuration
3. Compile the test and run
   - it will pass with no errors on success
4. Experiment with varying the size of the arrays (scaling)
   - start small and increase
5. finish the `newton.cu` example
   - how do the h2d, d2h and kernel timings compare?
6. **extra**: Compare scaling with the `axpy_omp` benchmark
7. **extra**: Experiment with varying the block size
   - try `cudaOccupancyMaxPotentialBlockSize`.

CSCS **ETH**zürich

# Exercise: Results



The GPU is a throughput device:

- the CUDA implementation is faster for $2^{15} \approx 32,000$
- requires $2^{20} \approx 1,000,000$ to get "advertised" $5\times$ speedup

You have to provide enough parallelism to exploit many cores

CSCS

**ETH** *zürich*