

General-purpose programming on GPU

Launching kernels

Eugenio Rustico
rustico@dmi.unict.it

D.M.I. - Università di Catania

Updated: April 9, 2011



Overview:

- Recall: launching kernels
- Further steps: 2D grids, shared memory
- Bottom-up approach
- Interact, any time!

Requisites:

- Mastery of C/C++
- Coffee!

Web: <http://www.dmi.unict.it/~bilotta/gpgpu/>

Forum: <http://forum.sdai.unict.it/index.php?board=116.0>

Overview

1 Introduction

2 Dive in: array initialization

- Constant init
- Non-constant init
- Function init
- Threaded init
- CUDA init

3 Launch granularity: blocks

- 1D blocks
- 2D blocks

4 Hands on code



oooooooooo

Let's start with a trivial example: initializing an array.



oooooooooo

Let's start with a trivial example: initializing an array.

Allocate space:

```
1 unsigned int vector_size = 10;  
2 hVector = (int*) malloc(vector_size*sizeof(int));
```

Let's start with a trivial example: initializing an array.

Allocate space:

```
1 unsigned int vector_size = 10;
2 hVector = (int*) malloc(vector_size*sizeof(int));
```

Iterate on elements:

```
3 for (int i=0; i<vector_size; i++)
4     hVector[i] = 0;
```

Let's start with a trivial example: initializing an array.

Allocate space:

```
1 unsigned int vector_size = 10;
2 hVector = (int*) malloc(vector_size*sizeof(int));
```

Use memset:

```
3 // must include <string.h>
4 memset(hVector, 0, vector_size*sizeof(int));
```

Initialization with non-constant value

Allocate space:

```
1 unsigned int vector_size = 10;
2 hVector = (int*) malloc(vector_size*sizeof(int));
```

Iterate on elements (cannot use memset):

```
3 for (int i=0; i<vector_size; i++)
4     hVector[i] = i;
```

Same example, with body of for cycle replaced by a function

```
1 void set_array_element( int* array , \
2     unsigned int element )
3 {
4     array [element] = element ;
5 }
```

Same example, with body of for cycle replaced by a function

```
1 void set_array_element( int* array , \
2     unsigned int element)
3 {
4     array[element] = element;
5 }

6 hVector = (int*) malloc(vector_size*sizeof(int));
7 for (int i=0; i<vector_size; i++)
8     set_array_element(hVector, i);
9 //hVector[i] = i;
```

With threads: iterations of for cycle are now *parallel* (in theory)

```
1 static void set_array_element( ... )
2 {
3     array[element] = element;
4 }

5 hVector = (int*) malloc(vector_size*sizeof(int));
6 for (int i=0; i<vector_size; i++)
7     pthread_create(&threadID, NULL, \
8         set_array_element, (...));
9 //hVector[i] = i;
```

With threads: iterations of for cycle are now *parallel* (in theory)

```
1 static void set_array_element( ... )
2 {
3     array[element] = element;
4 }
5
5 hVector = (int*) malloc(vector_size*sizeof(int));
6 for (int i=0; i<vector_size; i++)
7     pthread_create(&threadID, NULL, \
8         set_array_element, (...));
9     //hVector[i] = i;
```

Note: actually, pthreads can take only one parameter (e.g. pointer to class or struct).

Same initialization with CUDA:

```
1 __global__ void set_array_element( int* array)
2 {
3     array[threadIdx.x] = threadIdx.x;
4 }
5
6 cudaMalloc(&dVector, vector_size*sizeof(int));
7 //for (int i=0; i<vector_size; i++)
8 set_array_element<<<1, vector_size>>>(dVector);
9     //hVector[i] = i;
```

Same initialization with CUDA:

```
1 __global__ void set_array_element( int* array)
2 {
3     array[threadIdx.x] = threadIdx.x;
4 }
5
6 cudaMalloc(&dVector, vector_size*sizeof(int));
7 //for (int i=0; i<vector_size; i++)
8 set_array_element<<<1, vector_size>>>(dVector);
9     //hVector[i] = i;
10
11 // download array to access data
12 cudaMemcpy(hVector, dVector, \
13             vector_size*sizeof(int), \
14             cudaMemcpyDeviceToHost);
```

CUDA init

```
1 void set_array_element( int* array , uint element)
2 {
3     array[element] = element;
4 }
5 ...
6 hVector = (int*) malloc(vector_size*sizeof(int));
7 for (int i=0; i<vector_size; i++)
8     set_array_element(hVector, i);
```

```
1 __global__ void set_array_element( int* array )
2 {
3     array[threadIdx.x] = threadIdx.x;
4 }
5 ...
6 cudaMalloc(&dVector, vector_size*sizeof(int));
7 //for (int i=0; i<vector_size; i++)
8 set_array_element<<<1, vector_size>>>(dVector);
```

Find the differences



```
1 __global__ void set_array_element( int* array)
2 {
3     array[threadIdx.x] = threadIdx.x;
4 }
5 ...
6 cudaMalloc(&dVector, vector_size*sizeof(int));
7 set_array_element<<<1, vector_size>>>(dVector);
```

```
1 __global__ void set_array_element( int* array)
2 {
3     array[threadIdx.x] = threadIdx.x;
4 }
5 ...
6 cudaMalloc(&dVector, vector_size*sizeof(int));
7 set_array_element<<<1, vector_size>>>(dVector);
```

- `threadIdx` is *implicit* (it is a `dim3`: has x, y, and z fields)

```
1 __global__ void set_array_element( int* array)
2 {
3     array[threadIdx.x] = threadIdx.x;
4 }
5 ...
6 cudaMalloc(&dVector, vector_size*sizeof(int));
7 set_array_element<<<1, vector_size>>>(dVector);
```

- `threadIdx` is *implicit* (it is a `dim3`: has x, y, and z fields)
- `--global__` and `void` are **mandatory** for kernels

```
1 __global__ void set_array_element( int* array)
2 {
3     array[threadIdx.x] = threadIdx.x;
4 }
5 ...
6 cudaMalloc(&dVector, vector_size*sizeof(int));
7 set_array_element<<<1, vector_size>>>(dVector);
```

- `threadIdx` is *implicit* (it is a `dim3`: has x, y, and z fields)
- `--global__` and `void` are **mandatory** for kernels
- Non-standard syntax: compile with `nvcc`

```
1 __global__ void set_array_element( int* array)
2 {
3     array[threadIdx.x] = threadIdx.x;
4 }
5 ...
6 cudaMalloc(&dVector, vector_size*sizeof(int));
7 set_array_element<<<1, vector_size>>>(dVector);
```

- `threadIdx` is *implicit* (it is a `dim3`: has x, y, and z fields)
- `--global__` and `void` are **mandatory** for kernels
- Non-standard syntax: compile with `nvcc`
- `hVector` and `dVector` belong to separate *address spaces*:

```
1 __global__ void set_array_element( int* array)
2 {
3     array[threadIdx.x] = threadIdx.x;
4 }
5 ...
6 cudaMalloc(&dVector, vector_size*sizeof(int));
7 set_array_element<<<1, vector_size>>>(dVector);
```

- `threadIdx` is *implicit* (it is a `dim3`: has x, y, and z fields)
- `--global__` and `void` are **mandatory** for kernels
- Non-standard syntax: compile with `nvcc`
- `hVector` and `dVector` belong to separate *address spaces*:
 - `hVector` is allocated with `malloc()`, `dVector` with `cudaMalloc()`

```
1 __global__ void set_array_element( int* array)
2 {
3     array[threadIdx.x] = threadIdx.x;
4 }
5 ...
6 cudaMalloc(&dVector, vector_size*sizeof(int));
7 set_array_element<<<1, vector_size>>>(dVector);
```

- `threadIdx` is *implicit* (it is a `dim3`: has x, y, and z fields)
- `--global__` and `void` are **mandatory** for kernels
- Non-standard syntax: compile with `nvcc`
- `hVector` and `dVector` belong to separate *address spaces*:
 - `hVector` is allocated with `malloc()`, `dVector` with `cudaMalloc()`
 - we need to copy back the array as soon as we need to access it ("download")

Blocks

- Kernel launches are grouped in *blocks*, reflecting the hardware architecture

Blocks

- Kernel launches are grouped in *blocks*, reflecting the hardware architecture
- Each block is assigned to a multiprocessor

Blocks

- Kernel launches are grouped in *blocks*, reflecting the hardware architecture
- Each block is assigned to a multiprocessor
- The block size is up to the programmer (we will discuss it later)

Blocks

- Kernel launches are grouped in *blocks*, reflecting the hardware architecture
- Each block is assigned to a multiprocessor
- The block size is up to the programmer (we will discuss it later)
- Thread indexing becomes two-level

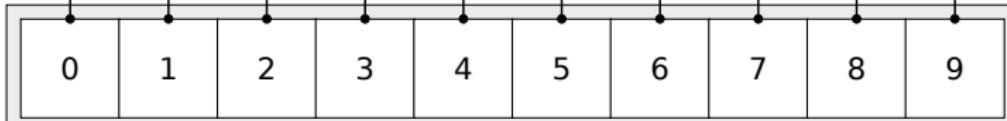
1D blocks

```
blockDim.x = 10    gridDim.x = 1
blockDim.y = 1    gridDim.y = 1
blockDim.z = 1    gridDim.z = 1
```

kernel grid

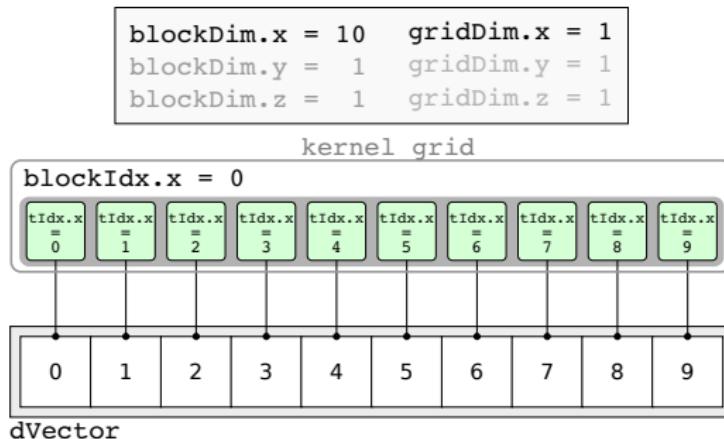
blockIdx.x = 0

tIdx.x = 0	tIdx.x = 1	tIdx.x = 2	tIdx.x = 3	tIdx.x = 4	tIdx.x = 5	tIdx.x = 6	tIdx.x = 7	tIdx.x = 8	tIdx.x = 9
------------	------------	------------	------------	------------	------------	------------	------------	------------	------------

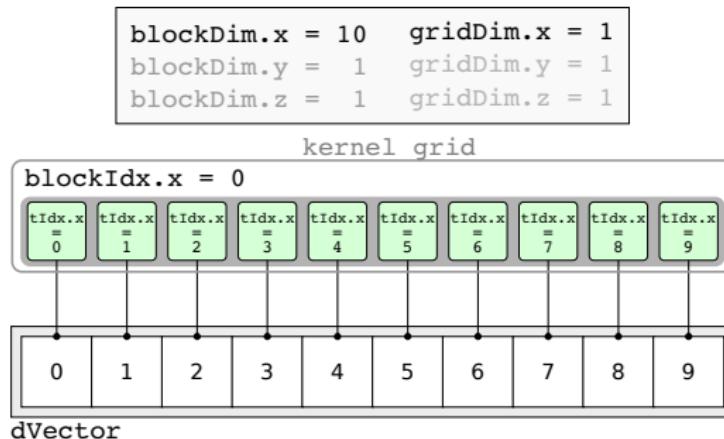


dVector

1D blocks



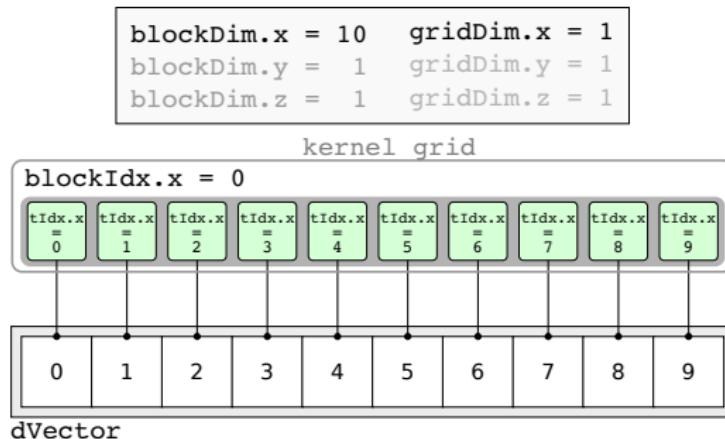
1D blocks



Kernel call:

```
1 kernel_name<<<1, vector_size>>>(dVector);
```

1D blocks



Kernel call:

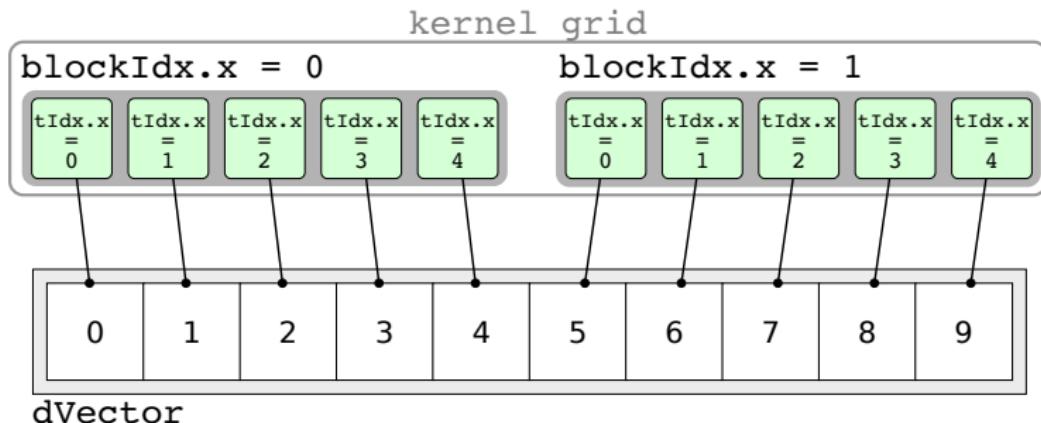
```
1 kernel_name<<<1, vector_size>>>(dVector);
```

Inside kernel:

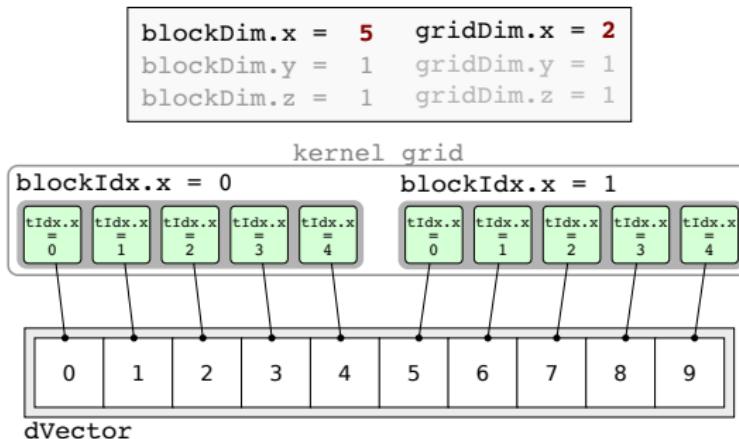
```
1 int myidx = threadIdx.x;
2 dVector[myidx] = ...;
```

1D blocks

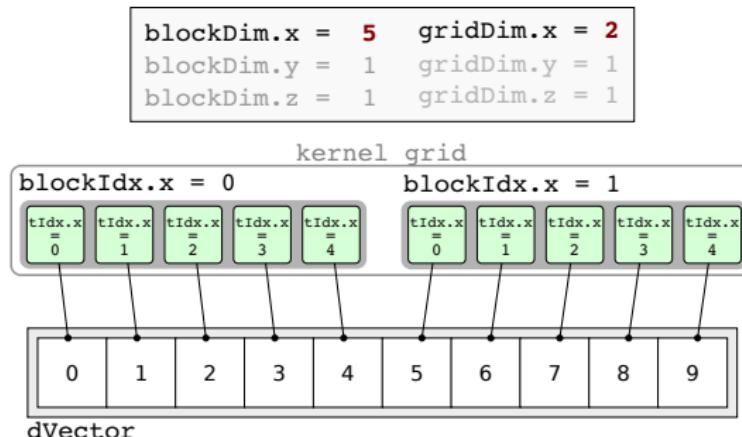
```
blockDim.x = 5    gridDim.x = 2
blockDim.y = 1    gridDim.y = 1
blockDim.z = 1    gridDim.z = 1
```



1D blocks



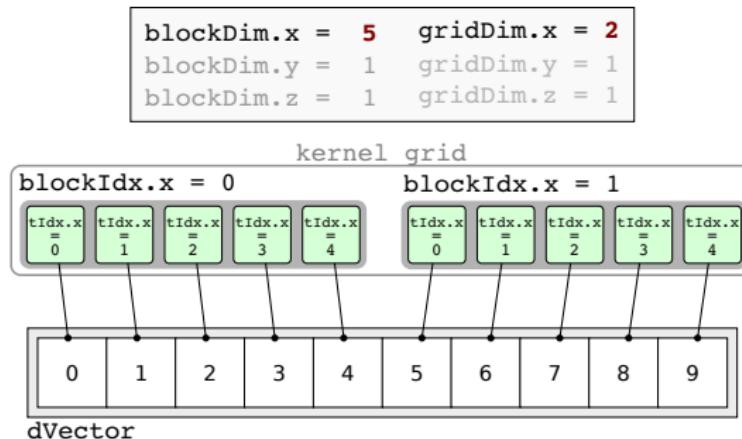
1D blocks



Kernel call:

```
1 kernel_name<<<2, vector_size/2>>>(dVector);
```

1D blocks



Kernel call:

```
1 kernel_name<<<2, vector_size/2>>>(dVector);
```

Inside kernel:

```
1 int myidx = blockDim.x*blockIdx.x + threadIdx.x;
2 dVector[myidx] = ...;
```

Kernel launch syntax:

```
kernel_name<<<dimGrid , dimBlock>>>(...);
```

Where

Kernel launch syntax:

```
kernel_name<<<dimGrid , dimBlock>>>(...);
```

Where

- `dimBlock` and `dimGrid` are `dim3`

Kernel launch syntax:

```
kernel_name<<<dimGrid , dimBlock>>>(...);
```

Where

- `dimBlock` and `dimGrid` are `dim3`
- `dimBlock` is 1D, 2D or 3D, `dimGrid` is 1D or 2D (up to CUDA 3.2)

Kernel launch syntax:

```
kernel_name<<<dimGrid , dimBlock>>>(...);
```

Where

- `dimBlock` and `dimGrid` are `dim3`
- `dimBlock` is 1D, 2D or 3D, `dimGrid` is 1D or 2D (up to CUDA 3.2)
- specifying an `int` is equivalent to a 1D `dim3`

Kernel launch syntax:

```
kernel_name<<<dimGrid , dimBlock>>>(...);
```

Where

- `dimBlock` and `dimGrid` are `dim3`
- `dimBlock` is 1D, 2D or 3D, `dimGrid` is 1D or 2D (up to CUDA 3.2)
- specifying an `int` is equivalent to a 1D `dim3`
- maximum `dimBlock` is smaller than maximum `gridDim`, see `deviceQuery`

```
kernel_name<<<dimGrid , dimBlock>>>(....);
```

What if dimBlock does not multiply dimGrid?

1D blocks

```
kernel_name<<<dimGrid , dimBlock>>>(....);
```

What if dimBlock does not multiply dimGrid?

- Launch sufficient number of blocks to cover all data (round up)

1D blocks

```
kernel_name<<<dimGrid , dimBlock>>>(....);
```

What if dimBlock does not multiply dimGrid?

- Launch sufficient number of blocks to cover all data (round up)
- Inside kernel, check address before accessing

```
1 int vector_size = 101;
2 int dimBlock = 10;
3 int dimGrid = \
4     (int) ceil(vector_size/(float)dimBlock);
5 kernel<<<dimGrid , dimBlock>>>(uint vec_size,....);
```

1D blocks

```
kernel_name<<<dimGrid , dimBlock>>>(....);
```

What if dimBlock does not multiply dimGrid?

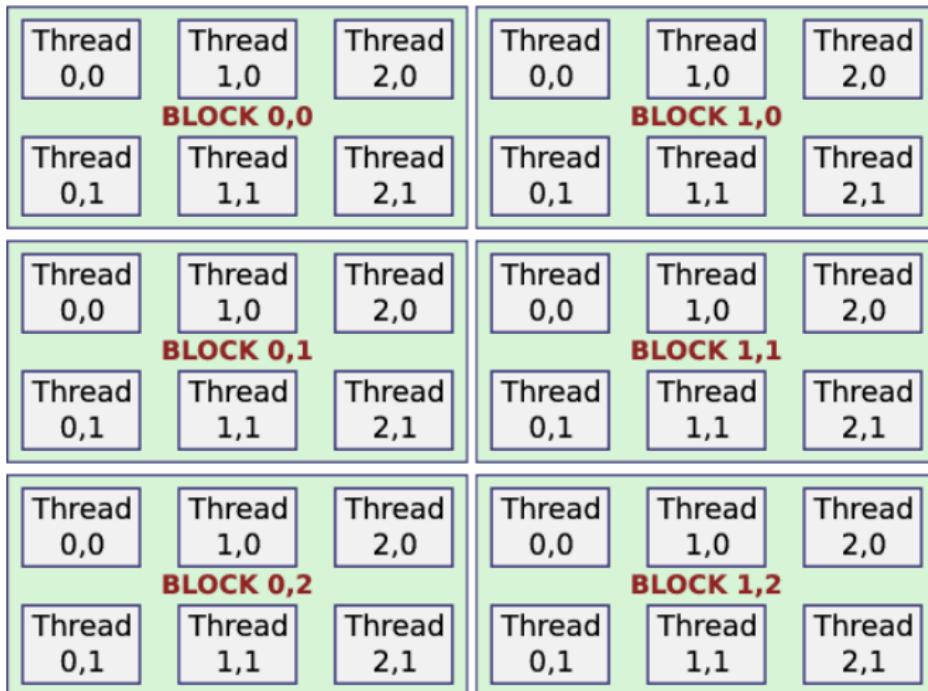
- Launch sufficient number of blocks to cover all data (round up)
- Inside kernel, check address before accessing

```
1 int vector_size = 101;
2 int dimBlock = 10;
3 int dimGrid = \
4     (int) ceil(vector_size/(float)dimBlock);
5 kernel<<<dimGrid , dimBlock>>>(uint vec_size,....);
```

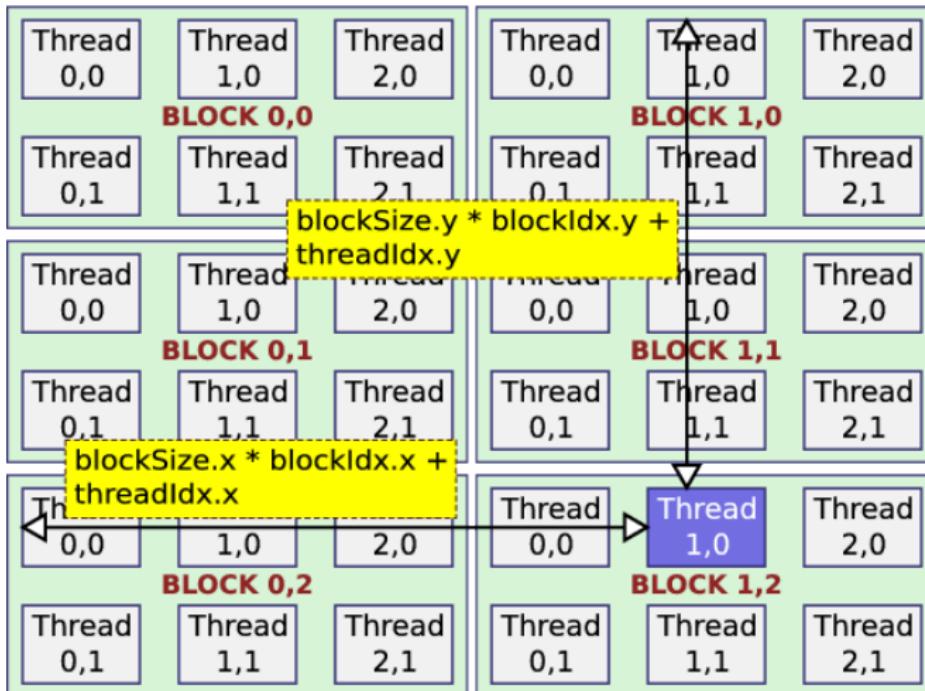
Inside kernel:

```
6 int myidx = blockDim.x*blockIdx.x + threadIdx.x;
7 if (myidx >= vec_size) return;
8 dVector[myidx] = ... ;
```

2D blocks



2D blocks



2D blocks: same addressing on both X and Y

```
1 dim3 vector_size(15,10,1);
2 dim3 dimBlock(10,10,1);
3 dim3 dimGrid( \
4     vector_size.x / dimBlock.x),
5     vector_size.y / dimBlock.y),
6     1);
7 kernel<<<dimGrid, dimBlock>>>(dim3 vec_size,...);
```

2D blocks

2D blocks: same addressing on both X and Y

```
1 dim3 vector_size(15,10,1);
2 dim3 dimBlock(10,10,1);
3 dim3 dimGrid( \
4     vector_size.x / dimBlock.x),
5     vector_size.y / dimBlock.y),
6     1);
7 kernel<<<dimGrid, dimBlock>>>(dim3 vec_size,...);
```

Inside kernel:

```
8 int x_coord = blockDim.x*blockIdx.x + threadIdx.x;
9 int y_coord = blockDim.y*blockIdx.y + threadIdx.y;
10 dVector[x_coord][y_coord] = ...;
```

...and same checks: what if does not multiply?

...and same checks: what if does not multiply?

```
1 dim3 vector_size(150,100,1);
2 dim3 dimBlock(16,16,1);
3 dim3 dimGrid( \
4     (int) ceil(vector_size.x/(float)dimBlock.x),
5     (int) ceil(vector_size.y/(float)dimBlock.y),
6     1);
7 kernel<<<dimGrid, dimBlock>>>(dim3 vec_size,...);
```

...and same checks: what if does not multiply?

```
1 dim3 vector_size(150,100,1);
2 dim3 dimBlock(16,16,1);
3 dim3 dimGrid( \
4     (int) ceil(vector_size.x/(float)dimBlock.x),
5     (int) ceil(vector_size.y/(float)dimBlock.y),
6     1);
7 kernel<<<dimGrid, dimBlock>>>(dim3 vec_size,...);
```

Inside kernel:

```
8 int x_coord = blockDim.x*blockIdx.x + threadIdx.x;
9 if (x_coord >= vec_size.x) return;
10 int y_coord = blockDim.y*blockIdx.y + threadIdx.y;
11 if (y_coord >= vec_size.y) return;
12 dVector[x_coord][y_coord] = ... ;
```

Hands on code:

- GPU vector initialization
- GPU vector initialization, k elements per thread
- GPU vector average
- GPU image smoothing, variable K
- Minimum scan