# That was Blocks in Parallel. What about Threads in Parallel

GPU

C

HOST

Kernel

Cuda

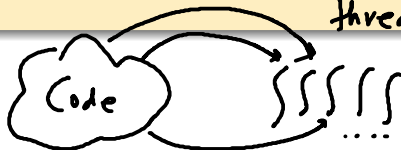## Function Call Change

```
// add<<<1, 1>>>(da, db, dc); // single thread GPU
// add<<<N, 1>>>(da, db, dc); // N blocks on GPU
   add<<<1, N>>>(da, db, dc); // N threads on GPU
```

## Changes in Kernel Code

```
__global__ void add(int *a, int *b, int *c)
{
    c[threadIdx.x] = a[threadIdx.x] + b[threadIdx.x];
}
```
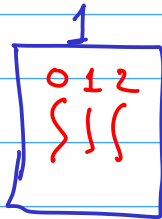
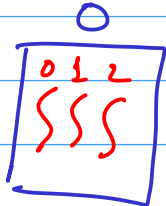Single instruction
multiple
threads

- Rest of Host code would be the same

id scheme .

body

Code

$3 \times 2 = 6$ threads.

ftnName <<< 2 , 3 >>>( ———— );

blocks          threads.

0               1

| 0 1 2 | | 0 1 2 |
| } } } | | } } } |

☆

| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

A[8]

+

| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

B[8]

GPU RAM

=

| 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |

C[8]

8 cores

$\frac{1}{8} = 12.5\%$

add <<< 8 , 1 >>> ( A, B, c );

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
 0   1   2   3   4   5   6   7

blockId

thread Id

$C[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x]$

| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

A[8]

+

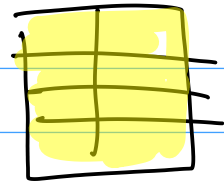| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

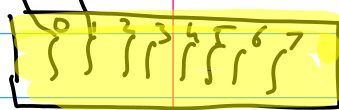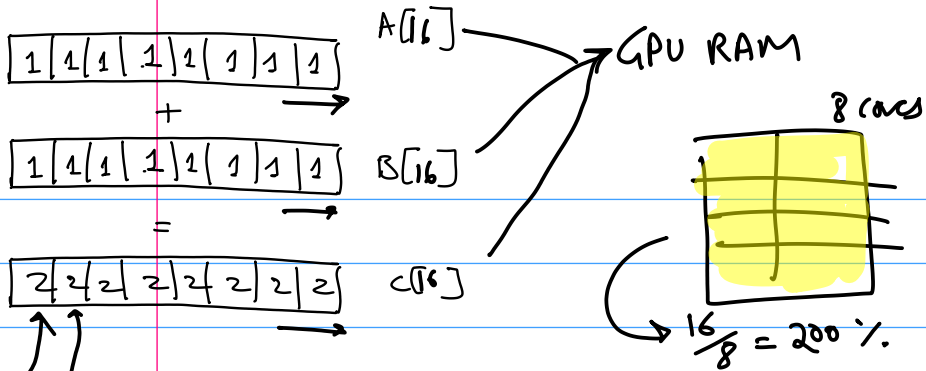B[8]

=

| 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |

C[8]

GPU RAM

8 cores

$$\frac{8}{8} = 100\%.$$

add $\lll \overline{1}, 8 \ggg (A, B, C);$

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

0

blockId

threadId

$$C[threadIdx.x] = a[threadIdx.x] + b[threadIdx.x]$$

| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

A[16]

GPU RAM

+

| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

B[16]

8 cores

=

| 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |

C[16]

$\frac{16}{8} = 200\%$

add <<< $\overline{1}$ , 16 >>> ( A, B, C );

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

0

blockId

thread Id

$$C[threadIdx.x] = a[threadIdx.x] + b[threadIdx.x]$$

0 . . . . . . . . . . . . . . 15          A[16] ───────→ GPU RAM

| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
+
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |          B[16]
=
| 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |          C[16]

8/8 = 100%          8/8 = 100%

add <<< 2 , 8 >>> (A, B, C);

0          1

01234567          01234567

0....1          0....1

blockId

thread Id

C[ Index ] = a[ Index ] + b[ Index ]
   0...7          0....7          0....7

8    8  → blockDim

0    1  → block Id



0 1 2 3 4 5 6 7    0 1 2 3 4 5 6 7  → thread Id.

$0 \times 8 + 0$

⓪ 1 2 3 4 5 6 ⑦ ⑧ 9 10 11 12 13 14 ⑮

$0 \times 8 + 7$    $1 \times 8 + 0$    $1 \times 8 + 7$

blockId
thread Id    ⟹    Inde

$$index = blockIdx.x * blockDim.x + threadIdx.x$$

blockDim    4        4         4.         4

blockIdx.x   0       1         2          3



threadId    0123    0123      0123       0123
x.x

$$index = blockIdx.x * blockDim.x + threadIdx.x$$

⟪ 2, 8 ⟫

⟪ 4, 4 ⟫

X, Y

0×4+0
0×8+0

  0  1  2  3  4  5  6  (7) (8) 9  10 11 12  13 14 15

0×8+7    1×8+0                    1×8+7

1×4+3   2×4+0                   3×4+3

for(j = 0; j < 4; j++)
   for(i = 0; i < 4; i++)

index = j*4 + i

block          thread

| j | | | | |
|---|---|---|---|---|
| 3 | 12 | 13 | 14 | 15 |
| 2 | 8 | 9 | 10 | 11 |
| 1 | 4 | 5 | 6 | 7 |
| 0 | 0 | 1 | 2 | 3 |
| | 0 | 1 | 2 | 3 |

i

$R \times 4 \times 2$
$+ j \times 4$
$+ i$

— slice 0 —

Slice 1

storage

| 4 | 5 | 6 | 7 |
|---|---|---|---|
| 0 | 1 | 2 | 3 |

| 12 | 13 | 14 | 15 |
|----|----|----|----|
| 8 | 9 | 10 | 11 |

1D 0-16

threads
1D

16x1
1x16

threads
2D

4x4

threads
3D

2x4x2



$j$

1

0

0   1   2   3

1

0

$k \times 4 \times 2$
$+ j \times 4$
$+ i$

$\longrightarrow$ Cuda

Index

blockIdx.z * blockDim.x * blockDim.y  ┐ 3D

+ blockIdx.y * blockDim.x  ┐ 2D

+ threadIdx.x  ┐ 1D

Data

Thread

(1D) → 1D

→ 2D

→ 3D

(2D)

What if

$8 \times 1 = 12\%$

$1 \times 8 = 50\%$

$4 \times 2 = 100\%$

$2 \times 2 \times 2 = 200\%$

Thread Dimensions

load balance

# Blocks or Threads. Whatever. Code is running in Parallel anyways

- Let's look at sample specs for NVIDIA Kepler K40
  - Cores = 2,880
  - Multiprocessors = 15
  - Cores / multi-processor = 192

- If we make parallel blocks, 1 block will be assigned to 1 multi-processor. This means all multi-processors will be busy, but within the multi-processor, 1 core has work to do, the remaining others are free.

- If we make parallel threads and 1 block, only 1 multi-processor will be busy, the remaining will be free.

- In either case, the GPU is **under-utilized**. For maximum utilization, use both (blocks + threads)

- **Note: Above is programmer's interpretation. The actual execution model loosely follows this interpretation but has other restrictions also.**
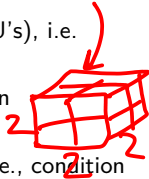
## Case: Higher Dimensions

$1D = <<< N, 1 >>>$
$<<< 1, N >>>$

**3D**

- So far, we have passed single values to kernel function call.
- For higher dimensions, use:
  - dimGrid($g_x$, $g_y$, $g_z$) for dimensions of the grid
  - dimBlock($t_x$, $t_y$, $t_z$) for dimensions of the thread block

dimGrid $(1,1,1)$
dimBlock $(2,2,2)$

- Total threads in thread-block cannot exceed 512 (or 1024 for some GPU's), i.e. condition $t_x \times t_y \times t_z \leq 512$ must be satisfied.
- Grid dimensions cannot exceed 32,768 in either dimension, i.e., condition $\max(g_x, g_y) \leq 32,768$
- dimGrid component must be evenly divisible by dimBlock component, i.e., condition $mod(g_n, t_n) = 0$ must be satisfied.
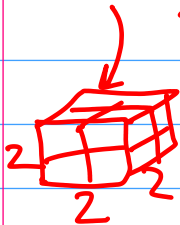
```
CL_DEVICE_MAX_WORK_ITEM_DIMENSIONS:3
CL_DEVICE_MAX_WORK_ITEM_SIZES:       1024 / 1024 / 64 (512/512/32)
CL_DEVICE_MAX_WORK_GROUP_SIZE:       1024 (512)
```

$2D = <<< M, N >>>$

dimGrid (2,2,2)
dimBlock (2,2,2) ≠



2 2
   2

3D threads



2   2   2   2
   2       2

2       2       2
   2       2       2

2   2       2
   2       2

6D

# Case: Higher Dimensions (cont.)
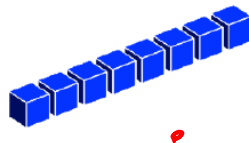
```
dim3 dimGrid(3,1,1);
dim3 dimBlock(8,1,1);

kernelfn <<< dimGrid, dimBlock >>> (arg1, arg2, ...);
```
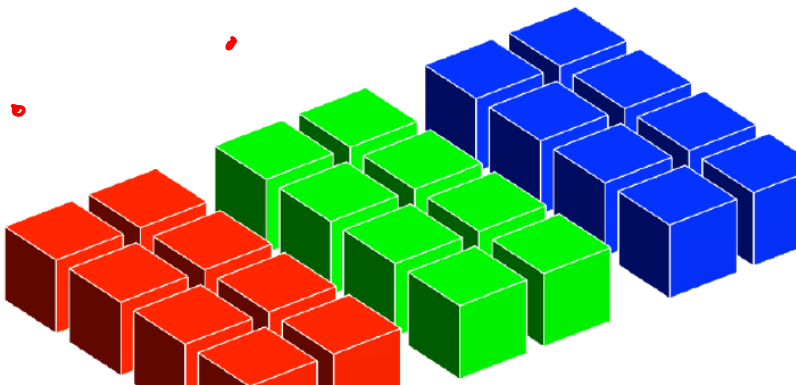
3          8

2D

→x

dim2   dim Grid (3,1)

dim2   dim Block ( 8, 1)

0 1 2 3 4 5 6 7

# Case: Higher Dimensions (cont.)

```
dim3 dimGrid(3,1,1);
dim3 dimBlock(2,4,1);

kernelfn <<< dimGrid, dimBlock >>> (arg1, arg2, ...);
```

# Case: Higher Dimensions (cont.)

$$I = I_x * gridDim.x * blockDim.x + I_y$$
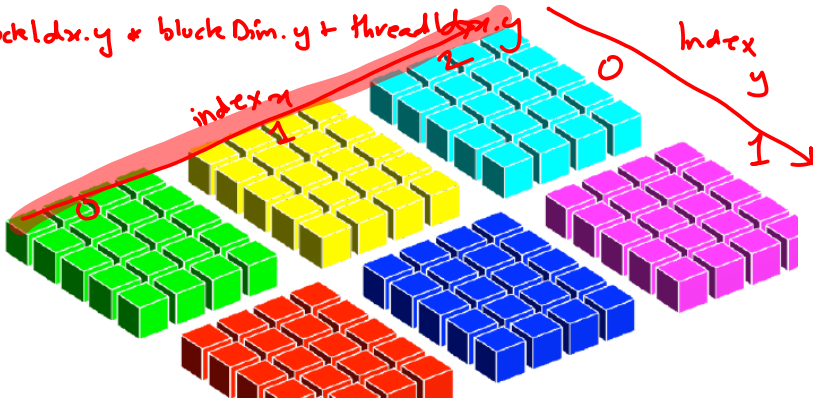
```
dim3 dimGrid(3,2,1);
dim3 dimBlock(4,5,1);

kernelfn <<< dimGrid, dimBlock >>> (arg1, arg2, ...);
```
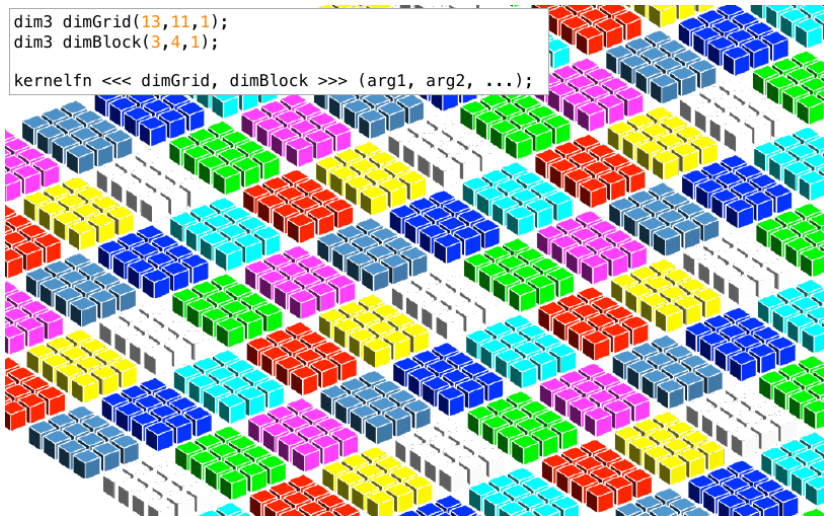
$I_x = blockIdx.x * blockDim.x + threadIdx.x$

$I_y = blockIdx.y * blockDim.y + threadIdx.y$

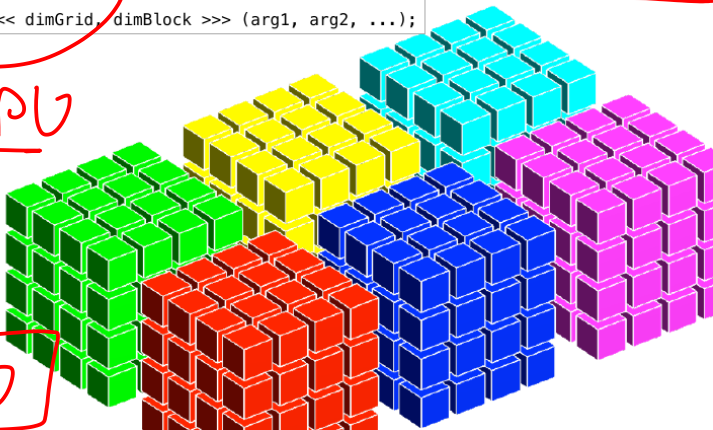# Case: Higher Dimensions (cont.)

```
dim3 dimGrid(13,11,1);
dim3 dimBlock(3,4,1);

kernelfn <<< dimGrid, dimBlock >>> (arg1, arg2, ...);
```

# Case: Higher Dimensions (cont.)

```
dim3 dimGrid(3,2,1);
dim3 dimBlock(4,4,4);

kernelfn <<< dimGrid, dimBlock >>> (arg1, arg2, ...);
```

threads

1 GPU

M GPU

② Host Code ⟶ Kernel

① Cuda Drivers.
    ↳ NVIDIA
        ↳ Propietary.

Open Source

Black list. ✗ Nouveau

White list
NVIDIA

# Combining both Concepts of Blocks and Threads

### Formula to identify a point

```
int index = blockIdx.x * blockDim.x + threadIdx.x;
```
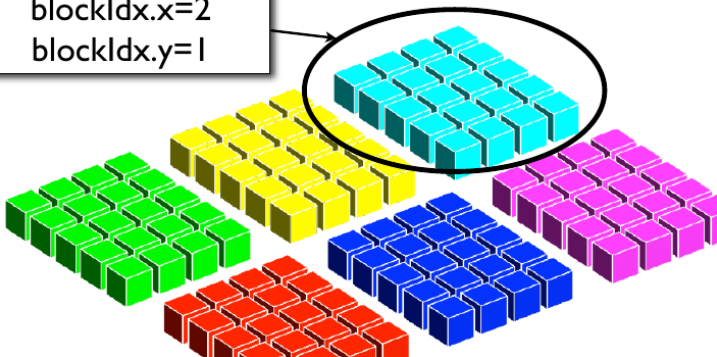
- blockIdx.x = Index of thread block in grid
- blockDim.x = Number of threads in 1D thread-block
- threadIdx.x = Index of thread in 1D thread block

- Each thread execute kernel code is automatically provided the following variables
  - threadIdx.x, threadIdx.y, threadIdx.z
  - blockDim.x, blockDim.y, blockDim.z
  - blockIdx.x, blockIdx.y

* gridDim.x , gridDim.y , gridDim.z

# Combining both Concepts of Blocks and Threads (cont.)

```
dim3 dimGrid(3,2,1);
dim3 dimBlock(4,5,1);

kernelfn <<< dimGrid, dimBlock >>> (arg1, arg2, ...);
```
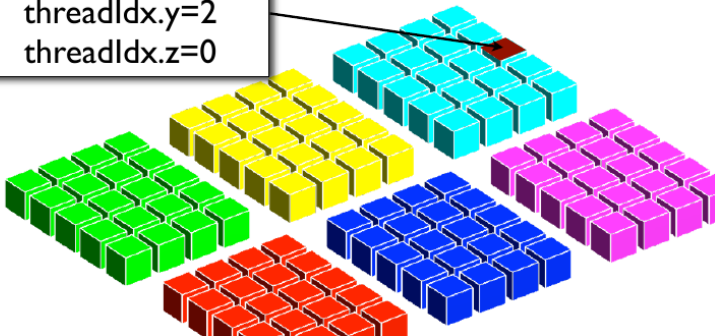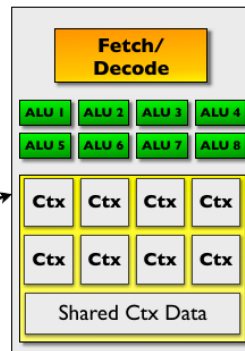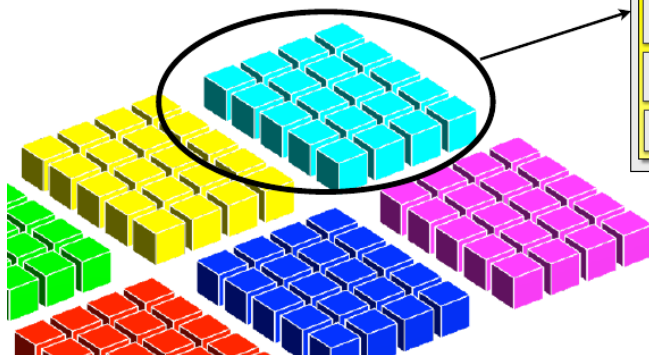
blockIdx.x=2
blockIdx.y=1

# Higher Dimensions: Execution Model Revisited

All the threads in an individual thread-block are handled by the same streaming multiprocessor.

# Higher Dimensions: Execution Model Revisited (cont.)



In this example batches of 8 threads will be processed concurrently