

CUDA Programming Model (Part 2)

ECE 277

Cheolhong An

Heterogeneous Computing

```
#include <iostream>
#include <algorithm>

using namespace std;

#define N          1024
#define RADIUS     3
#define BLOCK_SIZE 16

__global__ void stencil_1d(int *in, int *out) {
    __shared__ int temp[BLOCK_SIZE + 2 * RADIUS];
    int gindex = threadIdx.x + blockDim.x * blockDim.x;
    int index = threadIdx.x + RADIUS;

    // Read input elements into shared memory
    temp[index] = in[gindex];
    if (threadIdx.x < RADIUS) {
        temp[index - RADIUS] = in[gindex - RADIUS];
        temp[index + BLOCK_SIZE] = in[gindex + BLOCK_SIZE];
    }

    // Synchronize (ensure all the data is available)
    __syncthreads();

    // Apply the stencil
    int result = 0;
    for (int offset = -RADIUS ; offset <= RADIUS ; offset++)
        result += temp[index + offset];

    // Store the result
    out[gindex] = result;
}

void fill_ints(int *x, int n) {
    fill_n(x, n, 1);
}

int main(void) {
    int *in, *out;           // host copies of a, b, c
    int *d_in, *d_out;       // device copies of a, b, c
    int size = (N + 2*RADIUS) * sizeof(int);

    // Alloc space for host copies and setup values
    in = (int *)malloc(size); fill_ints(in, N+2*RADIUS);
    out = (int *)malloc(size); fill_ints(out, N+2*RADIUS);

    // Alloc space for device copies
    cudaMalloc((void **) &d_in, size);
    cudaMalloc((void **) &d_out, size);

    // Copy to device
    cudaMemcpy(d_in, in, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_out, out, size, cudaMemcpyHostToDevice);

    // Launch stencil 1d kernel on GPU
    stencil_1d<<<N/BLOCK_SIZE,BLOCK_SIZE>>>(d_in + RADIUS,
d_out + RADIUS);

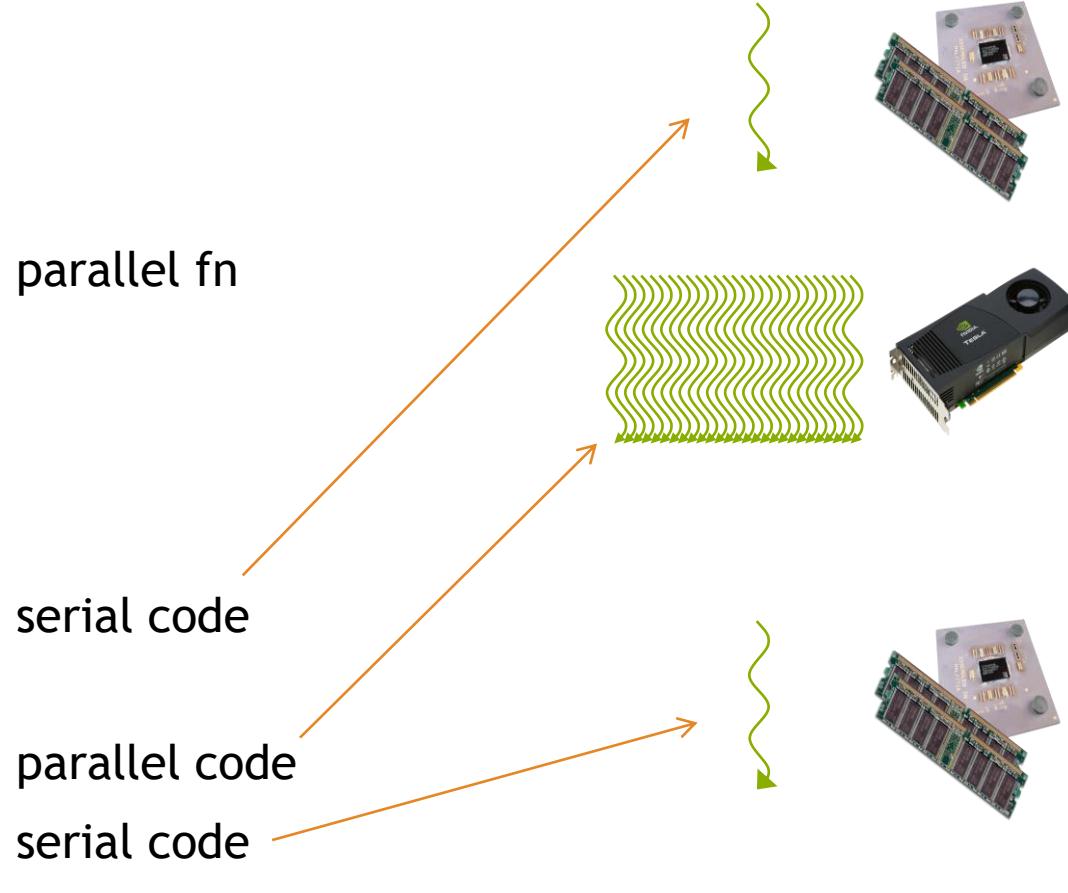
    // Copy result back to host
    cudaMemcpy(out, d_out, size, cudaMemcpyDeviceToHost);

    // Cleanup
    free(in); free(out);
    cudaFree(d_in); cudaFree(d_out);
    return 0;
}
```

parallel fn

serial code

parallel code
serial code



Common CUDA code structure

- Ex) checkThreadIndex

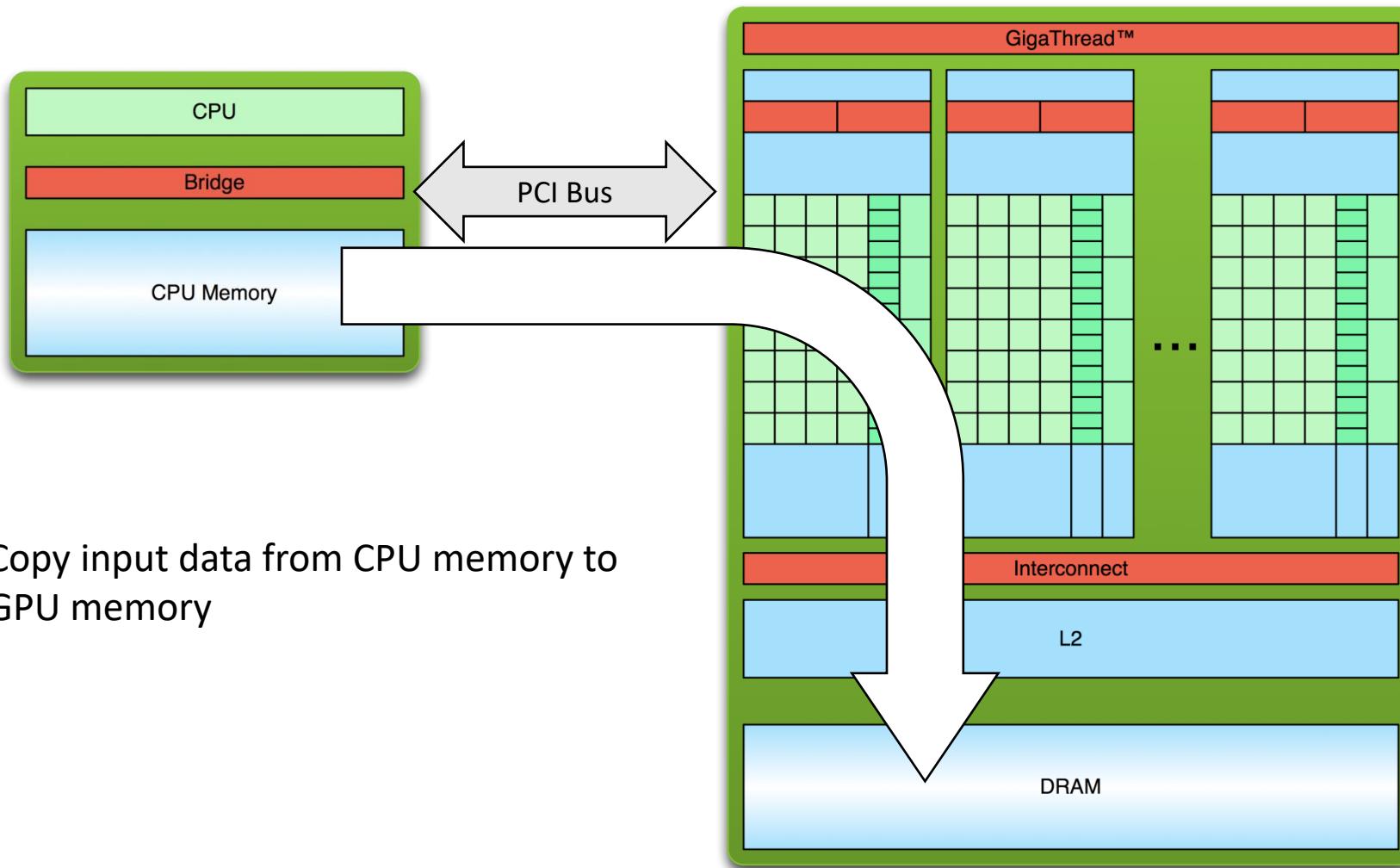
CPU memory allocation (malloc) <- Host
GPU memory allocation (cudaMalloc) <- Device

Transfer data from Host to Device
`cudaMemcpy(d_a, a, size, cudaMemcpyHostToDevice)`

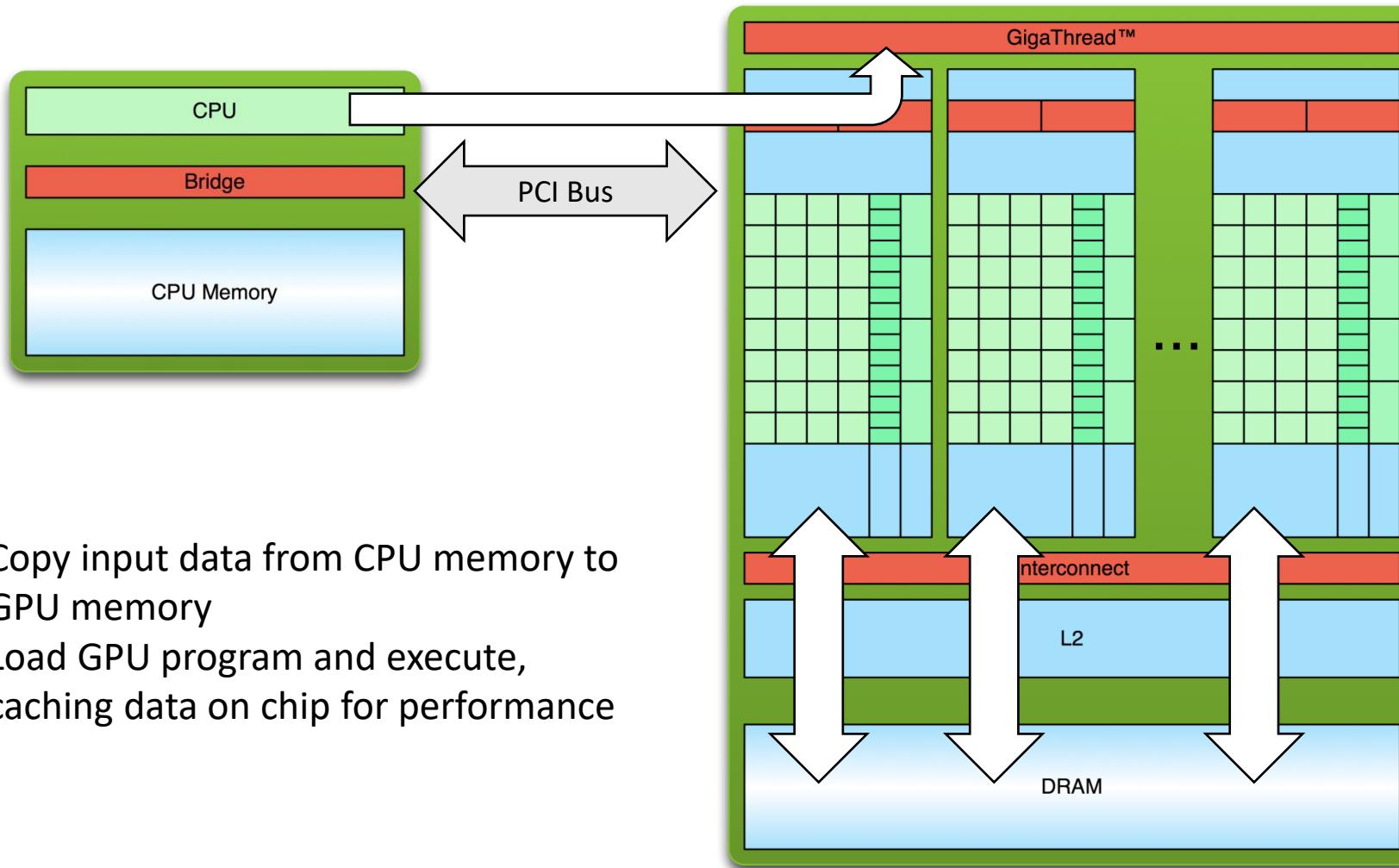
Run GPU codes (Call kernel functions)
`Kernel_fun<<<,>>> ()`

Transfer processed data from Device to Host
`cudaMemcpy(a, d_a, size, cudaMemcpyDeviceToHost)`

Transfer data from Host to Device

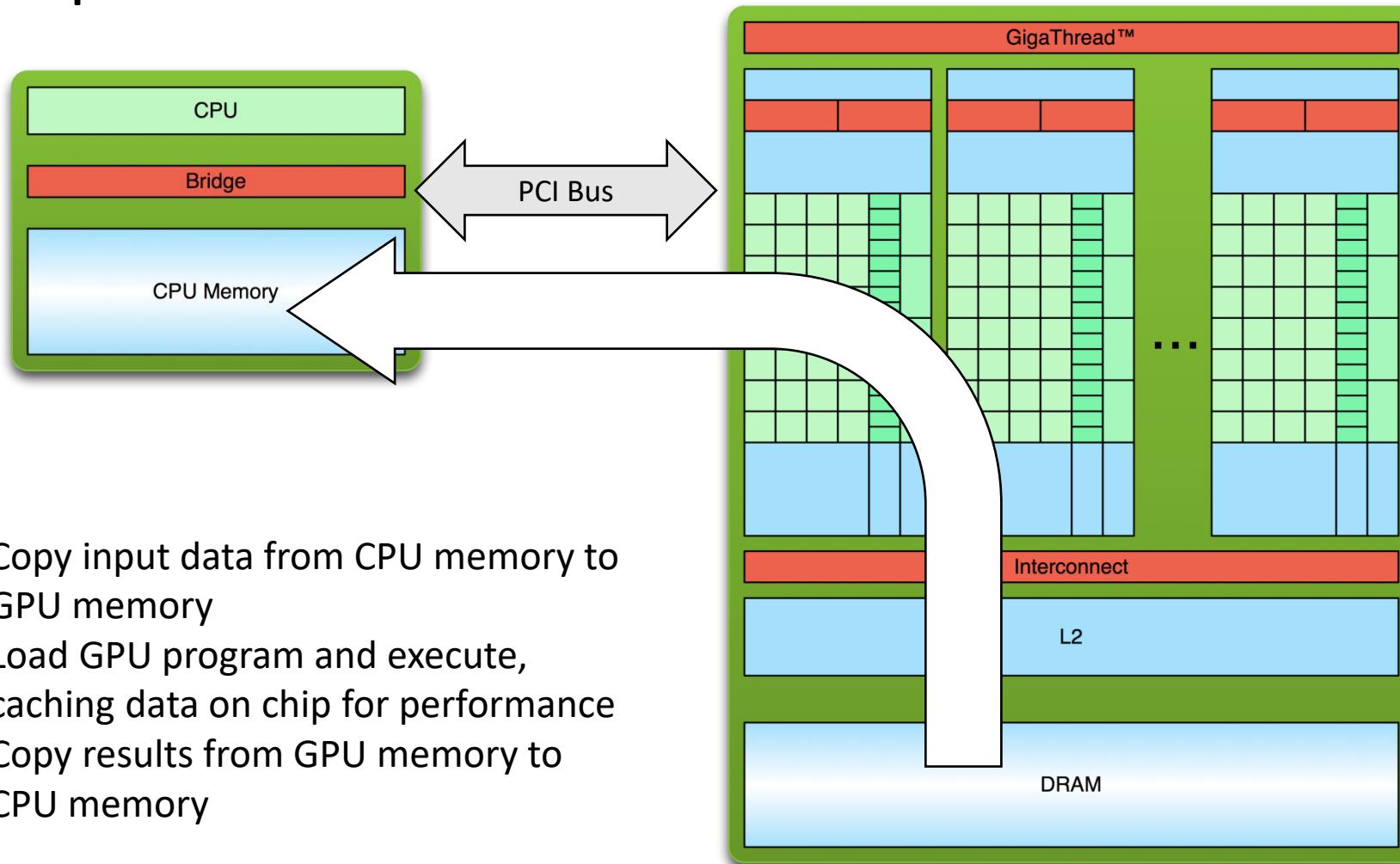


Run GPU codes



1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance

Transfer processed data from Device to Host



Multidimensional threads and threadblocks

- Why do we need 3 dimensional threads and threadblocks?
- What's advantages?
Application dependence
- We can think similar examples in C or C++.
 - 2D malloc allocation vs. 1D malloc allocation with 2D indexing

ex1) one-dimensional threads and threadblocks
dim3 nthreads;
nthreads.x = 4; nthreads.y = 1; nthreads.z = 1;
dim3 nblocks;
nblocks.x = 4; nblocks.y = 1; nblocks.z = 1;
kernel<<<nblocks, nthreads>>>()

ex2) two-dimensional threads and threadblocks
dim3 nthreads;
nthreads.x = 2; nthreads.y = 2; nthreads.z = 1;
dim3 nblocks;
nblocks.x = 2; nblocks.y = 2; nblocks.z = 1;
kernel<<<nblocks, nthreads>>>()

How to use blockIdx and threadIdx?

- **blockIdx and threadIdx are used to assign (map) data to threads**
- **blockIdx and threadIdx are used even for computation**

=> Domain knowledge is the key to know how to map effectively threads to data

For C++,

```
void copy_c(int *src_buf, *dst_buf) {  
    for (int w=0; w<10; w++) {  
        dst_buf[w] = src_buf[w] + w;  
    }  
}
```

source → destination

For CUDA,

```
--global__ void copy_cuda(int *src, *dst) {  
    dst[threadIdx.x] = src[threadIdx.x] + threadIdx.x;  
}
```

```
main(){  
    copy_c(src_buf, dst_buf);  
    copy_cuda<<<1,10>>>(src_buf, dst_buf);  
}
```



Mapping parallel threads to data (1/4)

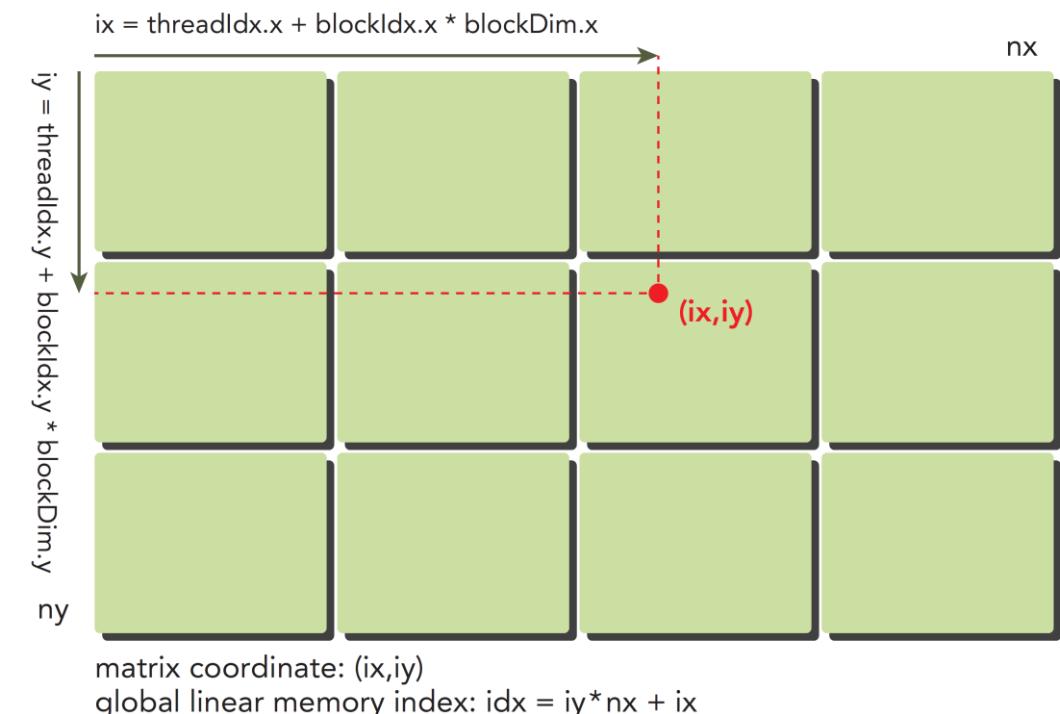
- For a given matrix ($n_x \times n_y$), we can map threads to elements of a matrix

$$\begin{aligned} ix &= \text{threadIdx.x} + \text{blockIdx.x} * \text{blockDim.x} \\ iy &= \text{threadIdx.y} + \text{blockIdx.y} * \text{blockDim.y} \end{aligned}$$

The linear index

$$\text{idx} = iy * n_x + ix$$

- Which mapping is used in this example?
1D grid and 1D blocks
1D grid and 2D blocks
2D grid and 1D blocks
2D grid and 2D blocks



Mapping parallel threads to data (2/4)

- 2D Grid, 2D Block

```
kernel<<<(2,3,1), (4,2,1) >>>(dst, src, nx)
```

```
--global__ void kernel ( int *dst , int *src , int nx ) {  
    idx_x = threadIdx.x + blockIdx.x*blockDim.x;  
    idx_y = threadIdx.y + blockIdx.y*blockDim.y;  
    linear_idx = idx_y*nx + idx_x;  
    dst [ linear_idx ] = src [ linear_idx ] + 2;  
}
```

If src and dst buffers are 1D array

If src and dst buffers are 2D array,
 $dst[idx_y][idx_x] = src[idx_y][idx_x] + 2;$

	Col 0	Col 1	Col 2	Col 3	Col 4	Col 5	Col 6	Col 7	nx
Row 0	0	1	2	3	4	5	6	7	
Row 1	8	9	10	11	12	13	14	15	
Row 2	16	17	18	19	20	21	22	23	
Row 3	24	25	26	27	28	29	30	31	
Row 4	32	33	34	35	36	37	38	39	
Row 5	40	41	42	43	44	45	46	47	

Mapping parallel threads to data (3/4)

- 2D Grid, 1D Block

kernel<<<(2,3,1), (8,1,1)>>>(dst, src, nx)

--global__ void kernel (int *dst, int *src, int nx) {
 idx_x = (threadIdx.x & 0x3) + blockIdx.x * blockDim.x;
 idx_y = (threadIdx.x >> 2) + blockIdx.y * blockDim.y;
 linear_idx = idx_y * nx + idx_x;
 dst[linear_idx] = src[linear_idx] + 2;
}

fusion operation
modular operation

0	1	2	3	4	5	6	7	nx
8	9	10	11	12	13	14	15	Row 0
16	17	18	19	20	21	22	23	Row 1
24	25	26	27	28	29	30	31	Row 3
32	33	34	35	36	37	38	39	Row 4
40	41	42	43	44	45	46	47	Row 5

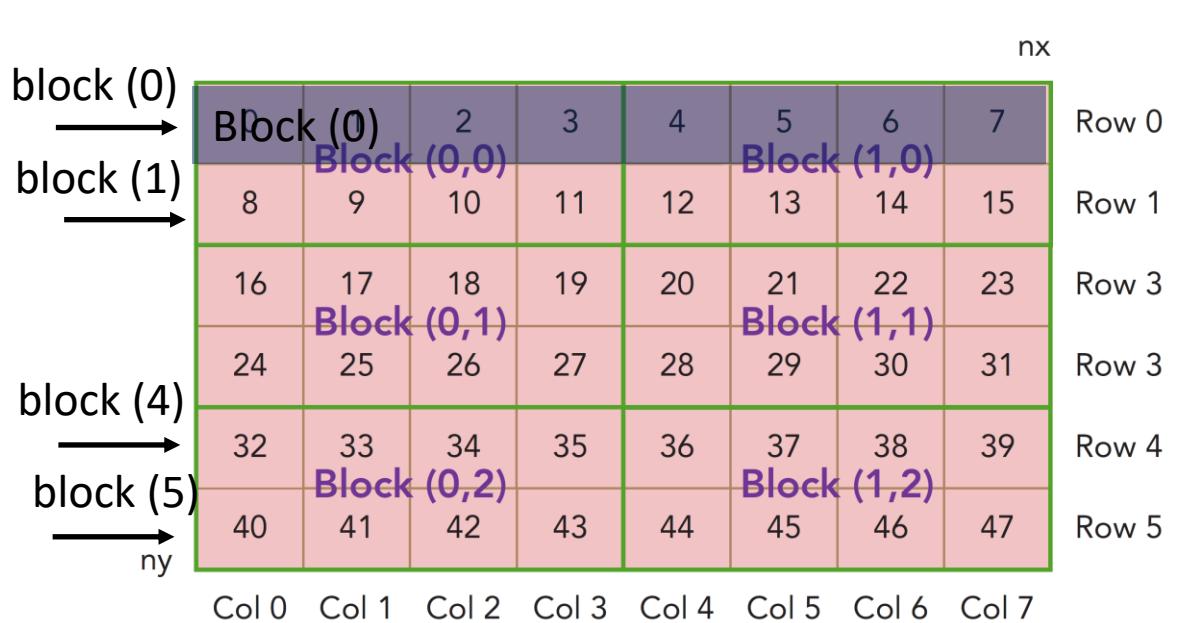
Mapping parallel threads to data (4/4)

- 1D Grid, 1D Block

```
kernel<<<(6,1,1), (8,1,1) >>>(dst, src, nx)
```

```
--global__ void kernel (int *dst, int *src, int nx) {
    idx_x = threadIdx.x;
    idx_y = blockIdx.x;
    linear_idx = idx_y*nx + idx_x;
    dst[linear_idx] = src[linear_idx] + 2;
}
```

Class_lab: c1_checkThreadIndex



Example: Grid and block partition (1/3)

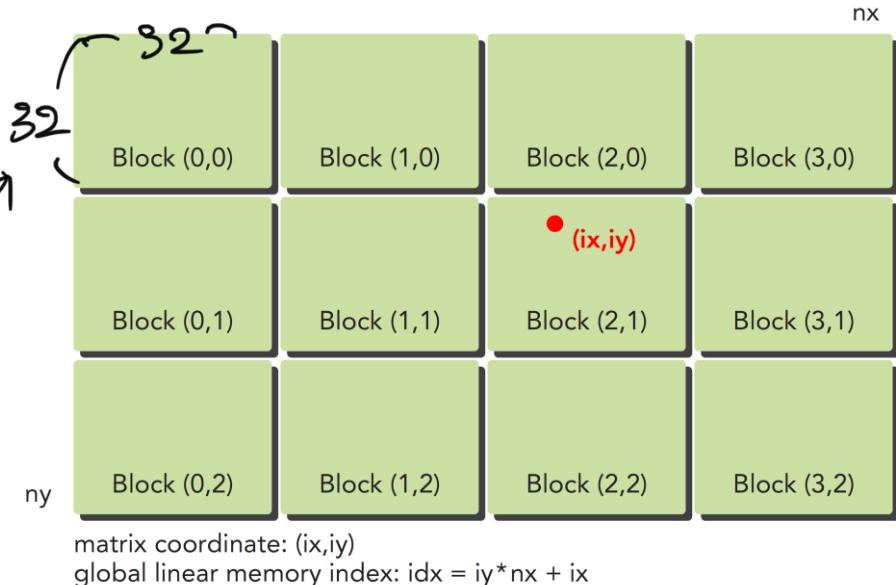
- 2D grid and 2D block partition

Map one thread to one matrix element

one threadblock processes a sub-block of a matrix (multiple columns and multiple rows)

```
--global__ void sumMatrixOnGPU2D( float *MatA, float *MatB, float *MatC, int nx, int ny )
{
    unsigned int ix = threadIdx.x + blockIdx.x * blockDim.x;
    unsigned int iy = threadIdx.y + blockIdx.y * blockDim.y;
    unsigned int idx = iy*nx + ix;
    if (ix < nx && iy < ny)
        MatC[ idx ] = MatA[ idx ] + MatB[ idx ];
}

int dimx = 32;
int dimy = 32;
dim3 block(dimx, dimy);
dim3 grid((nx+block.x-1)/block.x, (ny+block.y-1)/block.y);
sumMatrixOnGPU2D <<< grid, block >>>(d_MatA, d_MatB, d_MatC, nx, ny);
```



Example: Grid and block partition (2/3)

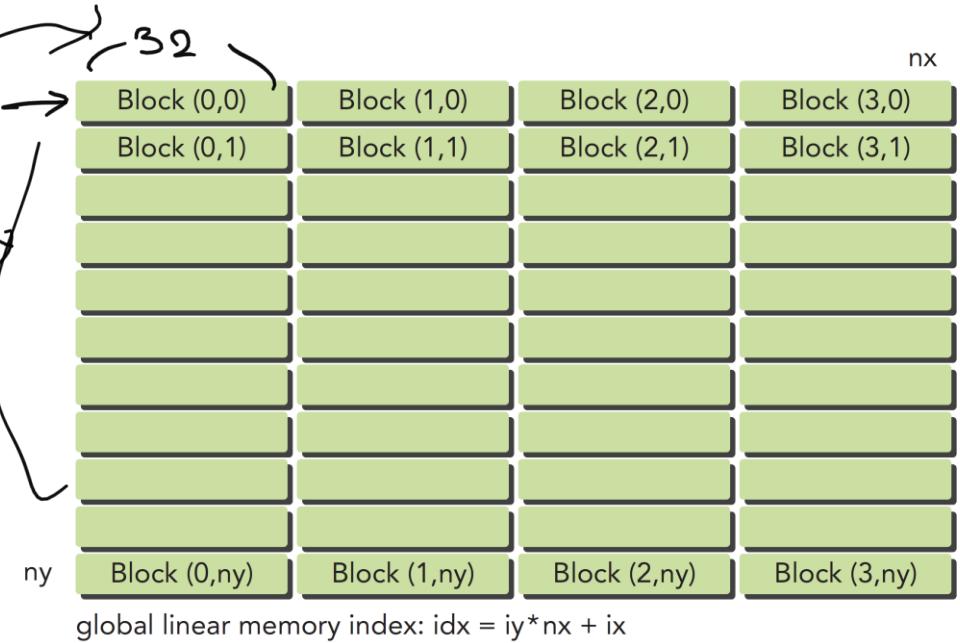
- 2D grid and 1D block partition

Map one thread to one matrix element

one threadblock processes multiple columns every row

```
--global__ void sumMatrixOnGPUMix( float *MatA, float *MatB, float *MatC, int nx, int ny )
{
    unsigned int ix = threadIdx.x + blockIdx.x * blockDim.x;
    unsigned int iy = blockIdx.y; one element/thread
    unsigned int idx = iy*nx + ix;
    if ( ix < nx && iy < ny )
        MatC[ idx ] = MatA[ idx ] + MatB[ idx ];
}

dim3 block ( 32 );
dim3 grid(( nx + block.x - 1 ) / block.x , ny );
sumMatrixOnGPUMIx <<< grid , block >>>(d_MatA , d_MatB , d_MatC , nx , ny );
```



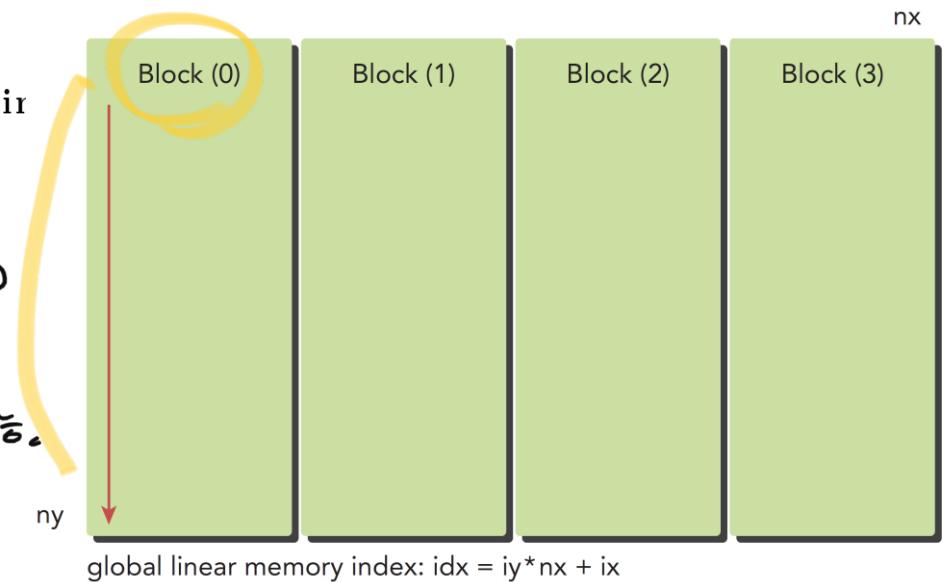
Example: Grid and block partition (3/3)

- 1D grid and 1D block partition

Map one thread to one matrix element

**However, one thread processes multiple elements
one threadblock processes a large matrix block**

```
--global__ void sumMatrixOnGPU1D( float *MatA, float *MatB, float *MatC, int nx, int ny ) {  
    unsigned int ix = threadIdx.x + blockIdx.x * blockDim.x;  
    if ( ix < nx ) {  
        for ( int iy=0; iy<ny; iy++ ) {  
            int idx = iy*nx + ix;  
            MatC[ idx ] = MatA[ idx ] + MatB[ idx ];  
        }  
    }  
}  
  
dim3 block ( 32 , 1 );  
dim3 grid (( nx+block.x-1)/block.x, 1 );  
sumMatrixOnGPU2D <<< grid , block >>>(d_MatA , d_MatB , d_MatC , nx , ny );
```



Performance vs. grid and block partitions

- The performance of GPU depends on grid and block partitions
- There is no simple rule to derive the best performance

It depends on GPU architecture, data size, access pattern and so on

However, **generally a single thread should process multiple elements in a threadblock with large data size**

Data size per operation per thread should be larger 4byte

(int, int2, int4, float, float2, float4)

(6 Byte/thread)

KERNEL	EXECUTION CONFIGURE	TIME ELAPSED
sumMatrixOnGPU2D	(512,1024), (32,16)	0.038041
sumMatrixOnGPU1D	(128,1), (128,1)	0.044701
sumMatrixOnGPUMix	(64,16384), (256,1)	0.030765

Example: CUDA Discrete Cosine Transform

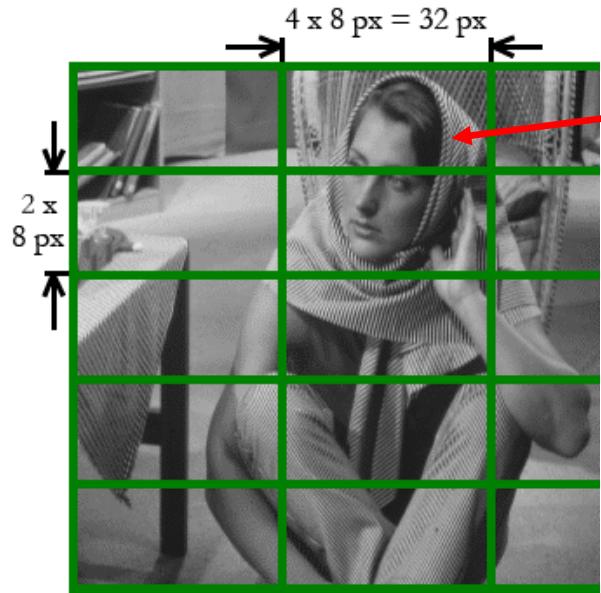
* lab
★ $C_1 - \text{block1d-grid1d}$, ⇒ which one is better?
(profiler → execution time)

8x4x2 threads to 8 8x8 blocks
samples (8 samples/thread)
2 warps/block

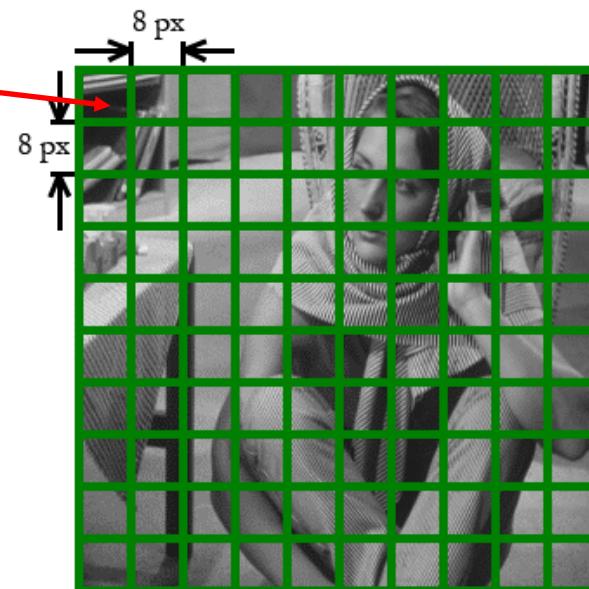
↙ better result.

vs.

8x8 threads to 8x8 block samples
(one sample/thread)
2 warps/block



thread block



Kernel <<< ((W+31)/32, (H+15)/16, 1) , (4,16,1) >>>

Kernel <<< ((W+7)/8, (H+7)/8, 1) , (8,8,1) >>>