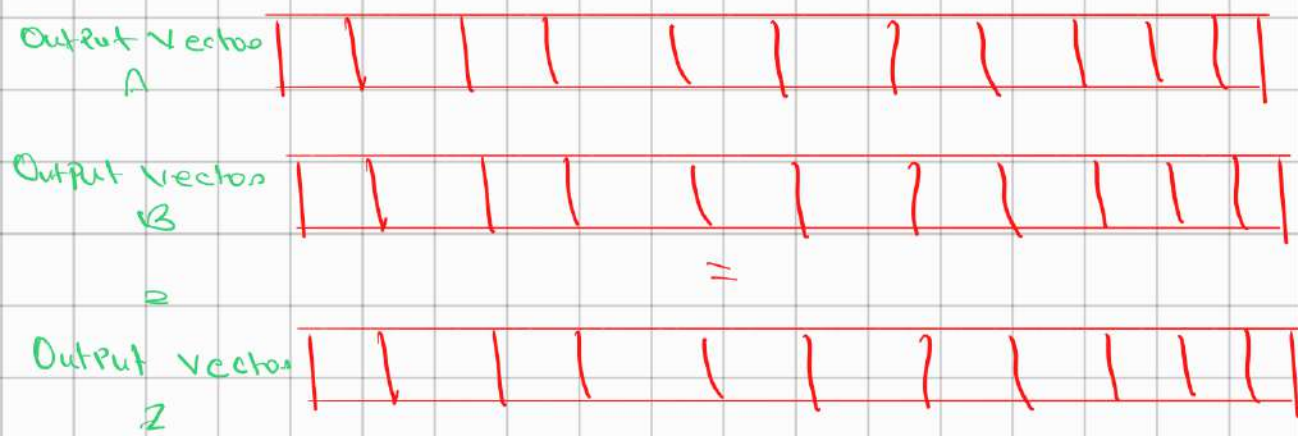# Types of Parallelism:-

## Task Parallelism
→ Different Operations performed on same or different data.

→ usually a modest number of tasks unleash a modest Amount of Parallelism.

## Data Parallelism.
→ Same Operations performed on different data.

→ Potentially Massive amount of data unleash Massive amount of Parallelism.
  ● Most Suitable for GPU.

# Vector Addition :- (Hello World for Cuda programming).

Output Vector A

Output Vector B

=

Output Vector Z

```
int main(int argc, char**argv) {

    cudaDeviceSynchronize();

    // Allocate memory and initialize data
    Timer timer;
    unsigned int N = (argc > 1)?(atoi(argv[1])):(1 << 25);
    float* x = (float*) malloc(N*sizeof(float));
    float* y = (float*) malloc(N*sizeof(float));
    float* z = (float*) malloc(N*sizeof(float));
    for (unsigned int i = 0; i < N; ++i) {
        x[i] = rand();
        y[i] = rand();
    }

    // Vector addition on CPU
    startTime(&timer);
    vecadd_cpu(x, y, z, N);
    stopTime(&timer);
    printElapsedTime(timer, "CPU time", CYAN);

    // Vector addition on GPU
    startTime(&timer);
    vecadd_gpu(x, y, z, N);
    stopTime(&timer);
    printElapsedTime(timer, "GPU time", DGREEN);

    // Free memory
    free(x);
    free(y);
    free(z);

    return 0;
```
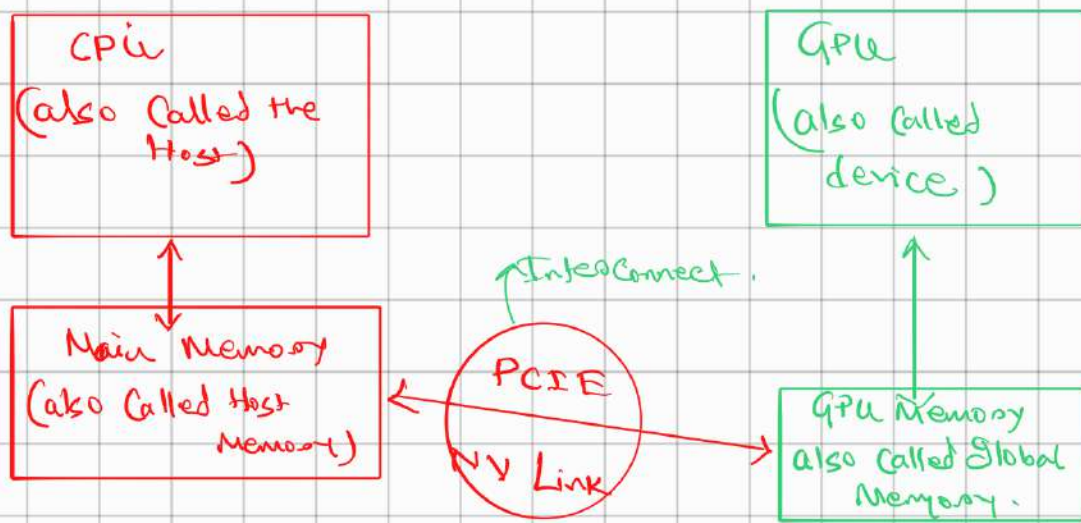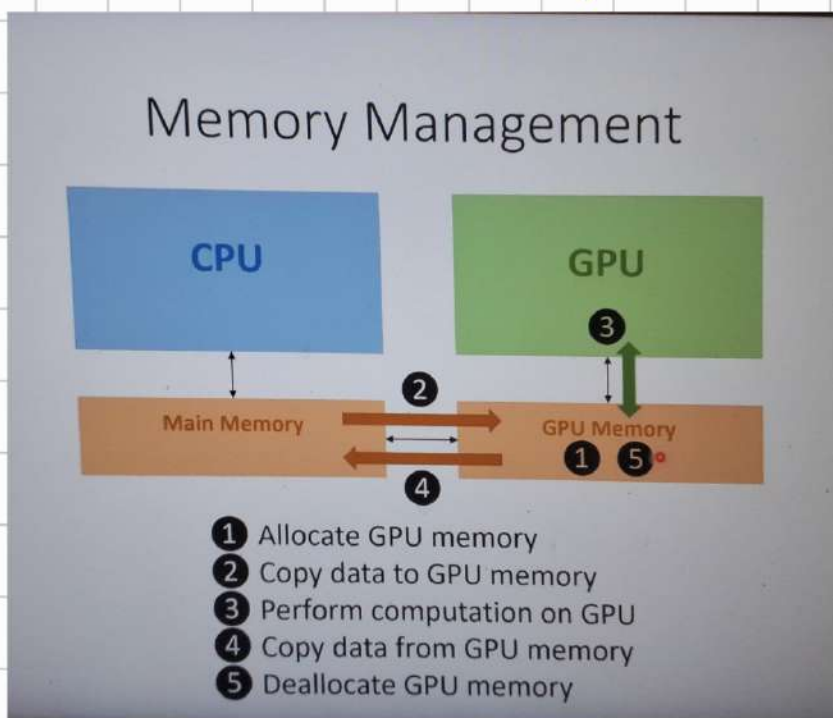
going to implement vec

# System Organization:-



CPU
(also Called the
Host)

GPU
(also Called
device)

↑ Interconnect.

Main Memory
(also Called Host
Memory)

PCIE
NV Link

GPU Memory
also Called Global
Memory.

The Cpu and Gpu have Separate memory and Cannot access each other's Memory.



## Memory Management

CPU          GPU

③

②

Main Memory          GPU Memory

① ⑤

④

① Allocate GPU memory
② Copy data to GPU memory
③ Perform computation on GPU
④ Copy data from GPU memory
⑤ Deallocate GPU memory

## CUDA Memory Management API

• Allocating memory:

    cudaError_t cudaMalloc(void **devPtr, size_t size)
    • devPtr: Pointer to pointer to allocated device memory
    • size: Requested allocation size in bytes

• Deallocating memory:

    cudaError_t cudaFree(void *devPtr)
    • devPtr: Pointer to device memory to free

• Return type: cudaError_t
    • Helps with error checking (discussed later)

# Cuda Memory Management API:-

● Allocating Memory:- Cuda Error_t CudaMalloc (void ** dev Ptr, size_t size)

    • dev Ptr : Pointer to Pointer to Allocate device Memory.
    • Size:- Requested Allocation Size in Bytes.

Void VecAdd _Gpu (float ** X , float **y, float *Z, int *N) {
    // Allocate GPu Memory
    float ** X_d , ** Y_d, *Z_d;
    CudaMalloc ((void**) & x_d, N * sizeof (float));
        ↳ functions Cannot modify Arguments so if I
        Want CudaMalloc to modify X-d. I should
        give a pointer to X_d.
    CudaMalloc ((void**) & Y_d , N * Sizeof (float));
    CudaMalloc ((void**) & Z_d, N * size of (float);

    //Copy to the Gpu.
    // Run the Gpu Code
    // Copy from the GPu.
    // Deallocate GPu Memory. // CudaFree (x_d);

## CUDA Memory Management API

• Copying memory:

```
cudaError_t cudaMemcpy(void *dst, const void *src,
                size_t count, enum cudaMemcpyKind kind)
```

• dst: Destination memory address
• src: Source memory address
• count: Size in bytes to copy
• kind: Type of transfer
  • cudaMemcpyHostToHost
  • cudaMemcpyHostToDevice
  • cudaMemcpyDeviceToHost
  • cudaMemcpyDeviceToDevice

// Copy to the GPU.
cudaMemcpy ( x-d, x , N* sizeof (float), cudaMemCpyHostToDevice).
cudaMemCpy (Y-d, y, N*size of (float), cudaMemCpyHostToDevice).

// Call a GPU kernel function:(launch a grid of threads)

Const Unsigned int numThreadsPerBlock = 512;
Const Unsigned Int numBlocks = (N+512-1)/512 → This is effectively a ceil Operation
vecadd-kernel<<< numBlocks, numThreadsPerBlock >>> (x-d, y-d, z-d, N);
Cuda Device Synchronize ( );
// Copy from GPU to CPU.
cudaMemCpy (Z, z-d, N*size of (float), cudaMemCpyDeviceToHost);

# IF I have N elements in my Vector that means I want to have N threads in total.

// Deallocate GPU Memory
Cuda Free (X-d);
cuda Free (Y-d);
Cuda Free ( Z-d);

## Code Example

```c
void vecadd(float* x, float* y, float* z, int N) {

    // Allocate GPU memory
    float *x_d, *y_d, *z_d;
    cudaMalloc((void**) &x_d, N*sizeof(float));
    cudaMalloc((void**) &y_d, N*sizeof(float));
    cudaMalloc((void**) &z_d, N*sizeof(float));

    // Copy data to GPU memory
    cudaMemcpy(x_d, x, N*sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(y_d, y, N*sizeof(float), cudaMemcpyHostToDevice);

    // Perform computation on GPU
    ...

    // Copy data from GPU memory
    cudaMemcpy(z, z_d, N*sizeof(float), cudaMemcpyDeviceToHost);

    // Deallocate GPU memory
    cudaFree(x_d);
    cudaFree(y_d);
    cudaFree(z_d);

}
```
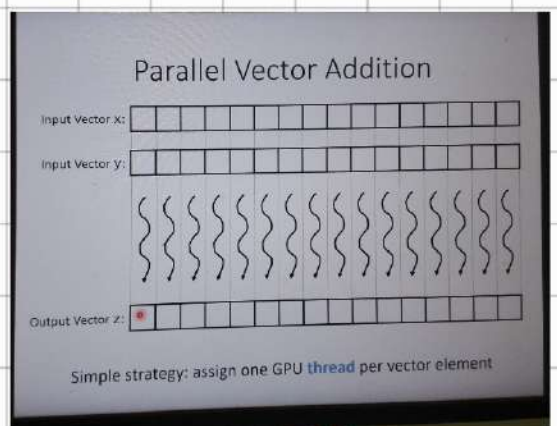
### Parallel Vector Addition

Input Vector X:
Input Vector Y:

Output Vector Z:

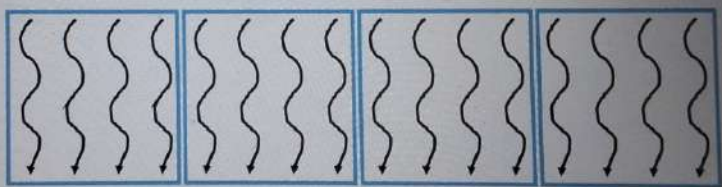Simple strategy: assign one GPU thread per vector element

# How threads on the GPU are Organized?
→ An array of GPU threads is called the grid.
From the CPU we create whole bunch of threads on the gpu. to Perform Some Operation these threads are Organised on the grid.

Blocks

Threads in a grid are grouped into thread **blocks**
(significance: threads in the same block can collaborate in ways
that threads in different blocks cannot – discussed later)

I would like to launch a grid of threads and to launch a grid of threads. I actually needs to specify How many blocks I need in the grid. and How many threads I want in each of my grid.

# LAUNCHING A GRID

* Threads in the same grid execute the same function known as a kernel. and the way to launch a grid is by calling this special function a kernel and telling this function what is the grid size in other Words the Number of the blocks size. Number of threads in each block.



Launching a Grid

• Threads in the same grid execute the same function known as a **kernel**

• A grid can be launched by calling a kernel and configuring it with appropriate grid and block sizes
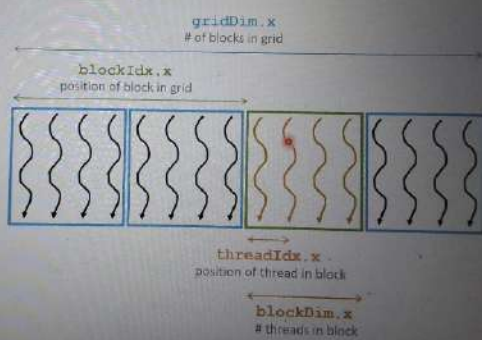
• Example:

```
const unsigned int numThreadsPerBlock = 512;
const unsigned int numBlocks = N/numThreadsPerBlock;
vecadd_kernel <<< numBlocks, numThreadsPerBlock >>> (x_d, y_d, z_d, N);
```



Implementing a Kernel

• A kernel is similar to a C/C++ function

• It is preceded by the keyword __global__ to indicate that it is a GPU kernel

• It uses special keywords to distinguish different threads from each other



Thread Index

gridDim.x
# of blocks in grid

blockIdx.x
position of block in grid

threadIdx.x
position of thread in block

blockDim.x
# threads in block

```
__global__ void vecadd_kernel (float * x, float * y, float * z, int N)
{
    unsigned int i = blockDim.x * blockIdx.x + threadIdx.x;

    if (i < N){

        z[i] = x[i] + y[i];

    }
}
```
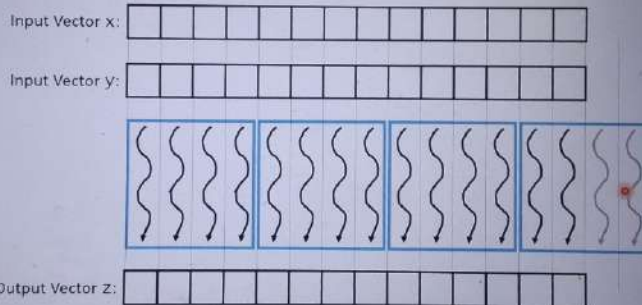→ Single Program Multiple Data.

## Boundary Conditions

- Recall launch configurations:
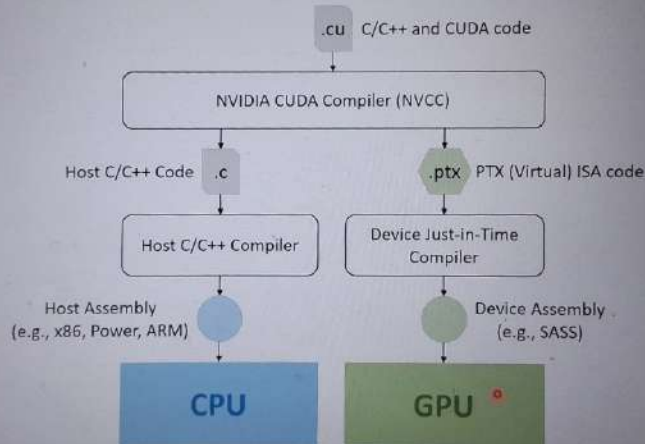
```
const unsigned int numThreadsPerBlock = 512;
const unsigned int numBlocks = N/numThreadsPerBlock;
vecadd_kernel <<< numBlocks, numThreadsPerBlock >>> (x_d, y_d, z_d, N);
```

- If N is not a multiple of numThreadsPerBlock, fewer threads will be launched than desired

- Solution: use the ceiling to launch extra threads then omit the threads after the boundary

## Execution with Boundary Checks

Input Vector x:

Input Vector y:

Output Vector z:

## Compilation

```
.cu  C/C++ and CUDA code
```

NVIDIA CUDA Compiler (NVCC)

```
Host C/C++ Code  .c          .ptx  PTX (Virtual) ISA code
```

Host C/C++ Compiler        Device Just-in-Time Compiler

Host Assembly            Device Assembly
(e.g., x86, Power, ARM)      (e.g., SASS)

**CPU**            **GPU**

## Function Declarations

- Keywords that differentiate where a function is intended to run

| Keyword | Callable From | Executed On |
|---|---|---|
| __host__ (default) | Host | Host |
| __global__ | Host (or Device) | Device |
| __device__ | Device | Device |

## Function Declarations

- Keywords that differentiate where a function is intended to run

| Keyword | Callable From | Executed On |
|---|---|---|
| __host__ (default) | Host | Host |
| __global__ | Host (or Device) | Device |
| __device__ | Device | Device |

- Why need __host__ if it is the default?

```
#include "timer.h"

__host__ __device__ float f(float a, float b) {       → We can use this function is both
    return a + b;                                          CPU and GPU.
}

void vecadd_cpu(float* x, float* y, float* z, int N) {
    for(unsigned int i = 0; i < N; ++i) {
        z[i] = f(x[i], y[i]);
    }
}

__global__ void vecadd_kernel(float* x, float* y, float* z, int N) {
    unsigned int i = blockDim.x*blockIdx.x + threadIdx.x;
    if(i < N) {
        z[i] = f(x[i], y[i]);
    }
}
```

## Function Declarations

- The keyword __host__ is useful when needing to mark a function as executable on both the host and the device

```
__host__ __device__ float f(float a, float b) {
    return a + b;
}

void vecadd(float* x, float* y, float* z, int N) {
    for(unsigned int i = 0; i < N; ++i) {
        z[i] = f(x[i], y[i]);
    }
}

__global__ void vecadd_kernel(float* x, float* y, float* z, int N) {
    int i = blockDim.x*blockIdx.x + threadIdx.x;
    if (i < N) {
        z[i] = f(x[i], y[i]);
    }
}
```

## Asynchronous Kernel Calls

- By default, kernel calls are asynchronous
  - Useful for overlapping GPU computations with CPU computations

- Use the following API function to synchronize
  ```
  cudaError_t cudaDeviceSynchronize()
  ```
  - Blocks until the device has completed all preceding requested tasks

## Error Checking

- All CUDA API calls return an error code `cudaError_t` that can be used to check if any errors occurred
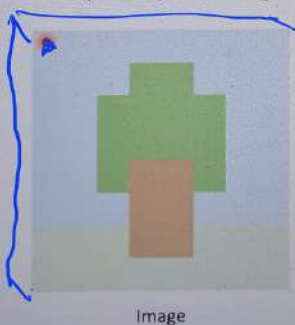  - Example:
    ```
    cudaError_t err = ...;
    if(err != cudaSuccess) {
        ... // Error handling code
    }
    ```

- For kernel calls, one can check the error returned by `cudaDeviceSynchronize()` or call the following API function:
  ```
  cudaError_t cudaGetLastError()
  ```

# LECTURE 3r Multi Dimensional GRIDS and Data.

## Multidimensional Grids

- CUDA supports multidimensional grids (up to 3D)
  - Simplifies processing multidimensional data

→ Block.

Image

Block and Thread Assignment

```
Timer timer;
// Allocate GPU memory
startTime(&timer);
unsigned char *red_d, *green_d, *blue_d, *gray_d;
cudaMalloc((void**) &red_d, width*height*sizeof(unsigned char));
cudaMalloc((void**) &green_d, width*height*sizeof(unsigned char));
cudaMalloc((void**) &blue_d, width*height*sizeof(unsigned char));
cudaMalloc((void**) &gray_d, width*height*sizeof(unsigned char));
cudaDeviceSynchronize();
stopTime(&timer);
printElapsedTime(timer, "Allocation time");

// Copy data to GPU
startTime(&timer);
cudaMemcpy(red_d, red, width*height*sizeof(unsigned char), cudaMemcpyHostToDevice);
cudaMemcpy(green_d, green, width*height*sizeof(unsigned char), cudaMemcpyHostToDevice);
cudaMemcpy(blue_d, blue, width*height*sizeof(unsigned char), cudaMemcpyHostToDevice);
cudaDeviceSynchronize();
stopTime(&timer);
printElapsedTime(timer, "Copy to GPU time");

// Call kernel
startTime(&timer);

cudaDeviceSynchronize();
stopTime(&timer);
printElapsedTime(timer, "Kernel time", GREEN);

// Copy data from GPU
startTime(&timer);
cudaMemcpy(gray, gray_d, width*height*sizeof(unsigned char), cudaMemcpyDeviceToHost);
cudaDeviceSynchronize();
stopTime(&timer);
printElapsedTime(timer, "Copy from GPU time");
```
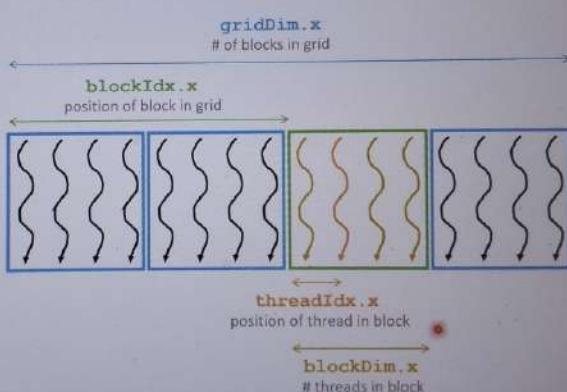
## Configuring Multidimensional Grids

- Use built-in `dim3` type
  - The rest is the same

```
dim3 numThreadsPerBlock(32, 32);
dim3 numBlocks((width + numThreadsPerBlock.x - 1)/numThreadsPerBlock.x,
               (height + numThreadsPerBlock.y - 1)/numThreadsPerBlock.y);
rgb2gray_kernel <<< numBlocks, numThreadsPerBlock >>>
                (red_d, green_d, blue_d, gray_d, width, height);
```

## Previously: One Dimensional Indexing

gridDim.x
# of blocks in grid

blockIdx.x
position of block in grid

threadIdx.x
position of thread in block

blockDim.x
# threads in block

## Multidimensional Indexing

- Built-in dimension and index variables each have thee components x, y, and z

blockIdx.x threadIdx.x

blockIdx.y

gridDim.y

threadIdx.y

blockDim.y

gridDim.x   blockDim.x

blockDim.x = 4    blockDim.y = 4

blockIdx.x * blockDim.x + threadIdx.x

## Layout of Multidimensional Data

- Convention is C is to store data in **row major order**
  - Elements in the same row are contiguous in memory

Logical view of data

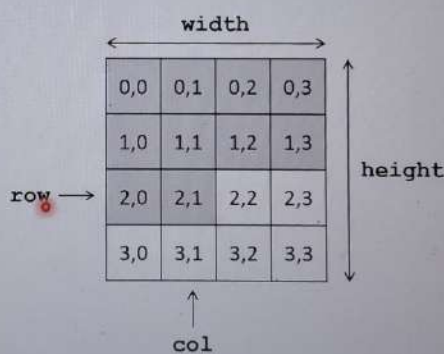| 0,0 | 0,1 | 0,2 | 0,3 |
|-----|-----|-----|-----|
| 1,0 | 1,1 | 1,2 | 1,3 |
| 2,0 | 2,1 | 2,2 | 2,3 |
| 3,0 | 3,1 | 3,2 | 3,3 |

Actual layout in memory

| 0,0 | 0,1 | 0,2 | 0,3 | 1,0 | 1,1 | 1,2 | 1,3 | 2,0 | 2,1 | 2,2 | 2,3 | 3,0 | 3,1 | 3,2 | 3,3 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|

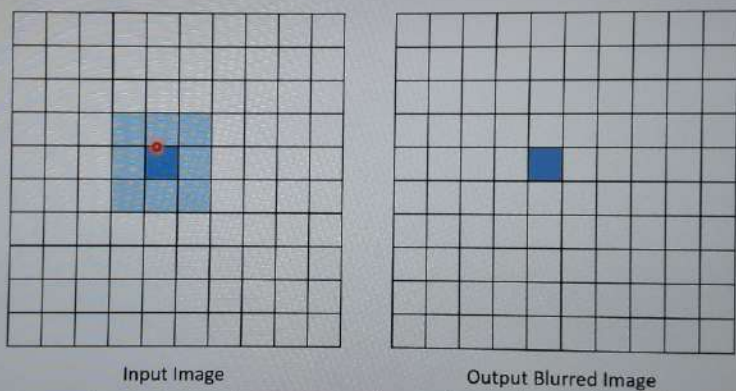So Now the question is How I will Convert My Row and Column and Convert Such that it Converts into a Row Major Array,

## Accessing Multidimensional Data

width

| 0,0 | 0,1 | 0,2 | 0,3 |
|-----|-----|-----|-----|
| 1,0 | 1,1 | 1,2 | 1,3 |
| 2,0 | 2,1 | 2,2 | 2,3 |
| 3,0 | 3,1 | 3,2 | 3,3 |

height

row →

col

Row * Width + Column.
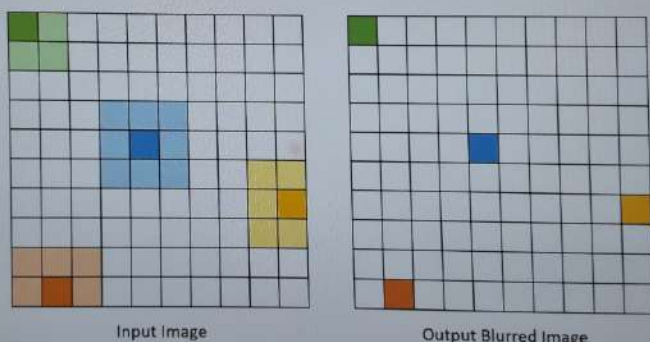
Example: Blur

## Example: Blur

Input Image

Output Blurred Image

☆ We can think of thread as being responsible for an Output Pixel. So what we Can do is We can have a thread for every Output Pixel. and each thread is responsible for finding the Corresponding input Pixel and looping over the Surrounding Pixels. and Computing the Average.

there May be More Great Approach for solving this Problem

## Boundary Conditions

Input Image

Output Blurred Image

## Boundary Conditions

```
__global__ void blur_kernel(unsigned char* image, unsigned char* blurred,
                            unsigned int width, unsigned int height) {

    int outRow = blockIdx.y*blockDim.y + threadIdx.y;
    int outCol = blockIdx.x*blockDim.x + threadIdx.x;

    if (outRow < height && outCol < width) {

        unsigned int average = 0;
        for(int inRow = outRow - BLUR_SIZE; inRow < outRow + BLUR_SIZE + 1; ++inRow) {
            for(int inCol = outCol - BLUR_SIZE; inCol < outCol + BLUR_SIZE + 1; ++inCol) {
                if(inRow >= 0 && inRow < height && inCol >= 0 && inCol < width) {
                    average += image[inRow*width + inCol];
                }                                         boundary check
            }
        }
        blurred[outRow*width + outCol] =
                (unsigned char)(average/((2*BLUR_SIZE + 1)*(2*BLUR_SIZE + 1)));

    }
}
```

**Rule of thumb:** every memory access must have a corresponding
guard the compares its indexes to the array dimensions

One of the key thing Rule of thumb is for boundary Condition is every time You have a Memory Access. You should have a Corresponding guard that checks if the index of the axis is in bounds.

# Example's MATRIX MATRIX Multiplication



Example: Matrix-Matrix Multiplication

$$C = A \times B$$



Example: Matrix-Matrix Multiplication

$$C = A \times B$$

Row →

Parallelization approach: assign one thread to each element in the output matrix (C)

```c
#include "common.h"
#include "timer.h"

__global__ void mm_kernel(float* A, float* B, float* C, unsigned int N) {

    unsigned int row = blockIdx.y*blockDim.y + threadIdx.y;
    unsigned int col = blockIdx.x*blockDim.x + threadIdx.x;

    float sum = 0.0f;
    for(unsigned int i = 0; i < N; ++i) {
        sum += A[row*N + i]*B[i*N + col];
    }
    C[row*N + col] = sum;

}

void mm_gpu(float* A, float* B, float* C, unsigned int N) {

    Timer timer;

    // Allocate GPU memory
    startTime(&timer);
    float *A_d, *B_d, *C_d;
    cudaMalloc((void**) &A_d, N*N*sizeof(float));
    cudaMalloc((void**) &B_d, N*N*sizeof(float));
    cudaMalloc((void**) &C_d, N*N*sizeof(float));
    cudaDeviceSynchronize();
    stopTime(&timer);
    printElapsedTime(timer, "Allocation time");

    // Copy data to GPU
    startTime(&timer);
    cudaMemcpy(A_d, A, N*N*sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(B_d, B, N*N*sizeof(float), cudaMemcpyHostToDevice);
    cudaDeviceSynchronize();
    stopTime(&timer);
"kernel.cu" 66L, 1875C
```

# Different Matrix Dimensions

$$C = A \times B$$



Assignments.