

Name: Muhammad Allah Rakha

Roll No: P19-0006 BCS-6A (PDC)

# Assignment: NVIDIA CUDA on Google Colab

GitHub Link Code: <https://github.com/aaaastark/NVIDIA-CUDA-Google-Colab.git>

## CUDA

### Kernel, Grid, Block, Threads and Dimensional of a Block/Grid

CUDA C extends C by allowing the programmer to define C functions, called kernels, that, when called, are executed N times in parallel by N different CUDA threads, as opposed to only once like regular C functions.

A kernel is defined using the `__global__` declaration specifier and the number of CUDA threads that execute that kernel for a given kernel call is specified using a new `<<<...>>>` execution configuration syntax. Each thread that executes the kernel is given a unique thread ID that is accessible within the kernel through the built-in `threadIdx` variable.

**Syntax:** `Kernel_Name<<< GridSize, BlockSize, SMEMSize, Stream >>> (arg,...);`

**SMEMsize:** is the size of Shared Memory at Runtime.

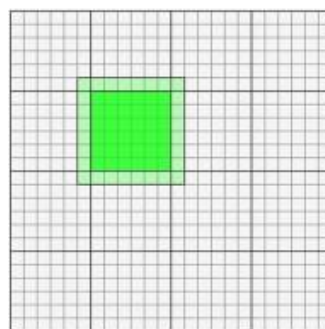
**Stream:** is a stream on which kernel will execute.

### Sample Example:

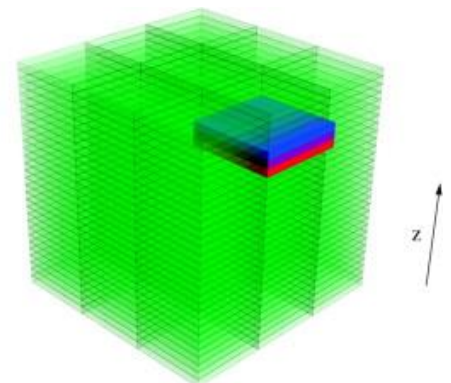
Here, each of the N threads that execute `VecAdd()` performs one pair-wise addition.

```
// Kernel definition
__global__ void VecAdd(float* A, float* B, float* C)
{
    int i = threadIdx.x; C[i] = A[i] + B[i];
}
int main()
{
    ...
    // Kernel invocation with N threads
    VecAdd<<<1, N>>>(A, B, C);
    ...
}
```

For convenience, `threadIdx` is a 3-component vector, so that threads can be identified using a one-dimensional, two-dimensional, or three-dimensional thread index, forming a one-dimensional, two-dimensional, or three-dimensional thread block. This provides a natural way to invoke computation across the elements in a domain such as a vector, matrix, or volume.

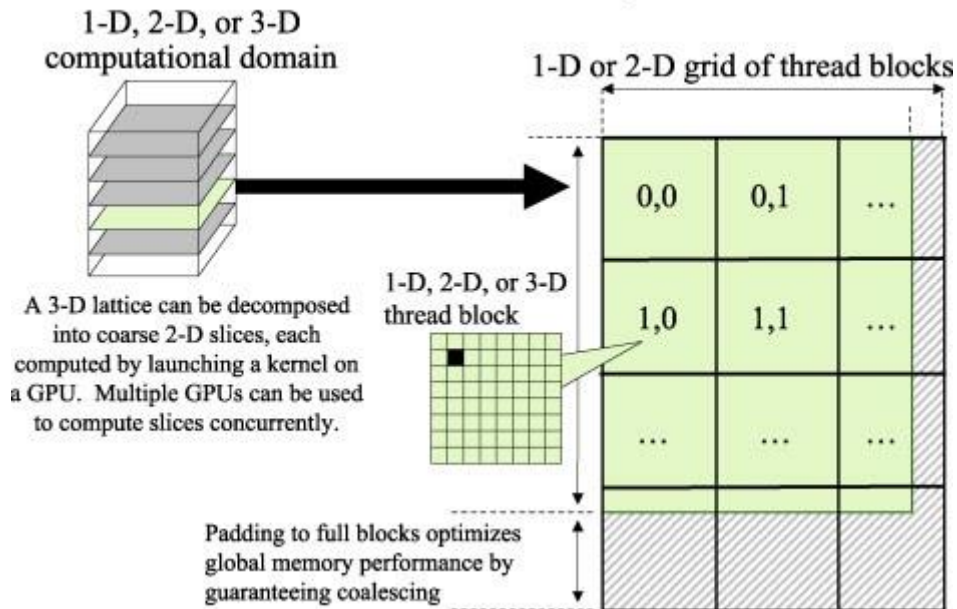


(a)



(b)

## CUDA Parallel Decomposition



The index of a thread and its thread ID relate to each other in a straightforward way: For a one-dimensional block, they are the same; for a two-dimensional block of size  $(D_x, D_y)$ , the thread ID of a thread of index  $(x, y)$  is  $(x + y D_x)$ ; for a three-dimensional block of size  $(D_x, D_y, D_z)$ , the thread ID of a thread of index  $(x, y, z)$  is  $(x + y D_x + z D_x D_y)$ .

As an example, the following code adds two matrices A and B of size  $N \times N$  and stores the result into matrix.

**C Language:**

```
// Kernel definition
__global__ void MatAdd(float A[N][N], float B[N][N], float C[N][N])
{
    int i = threadIdx.x;
    int j = threadIdx.y;
    C[i][j] = A[i][j] + B[i][j];
}
int main()
{
    ... // Kernel invocation with one block of N * N * 1 threads
    int numBlocks = 1;
    dim3 threadsPerBlock(N, N);
    MatAdd<<<numBlocks, threadsPerBlock>>>(A, B, C);
    ...
}
```

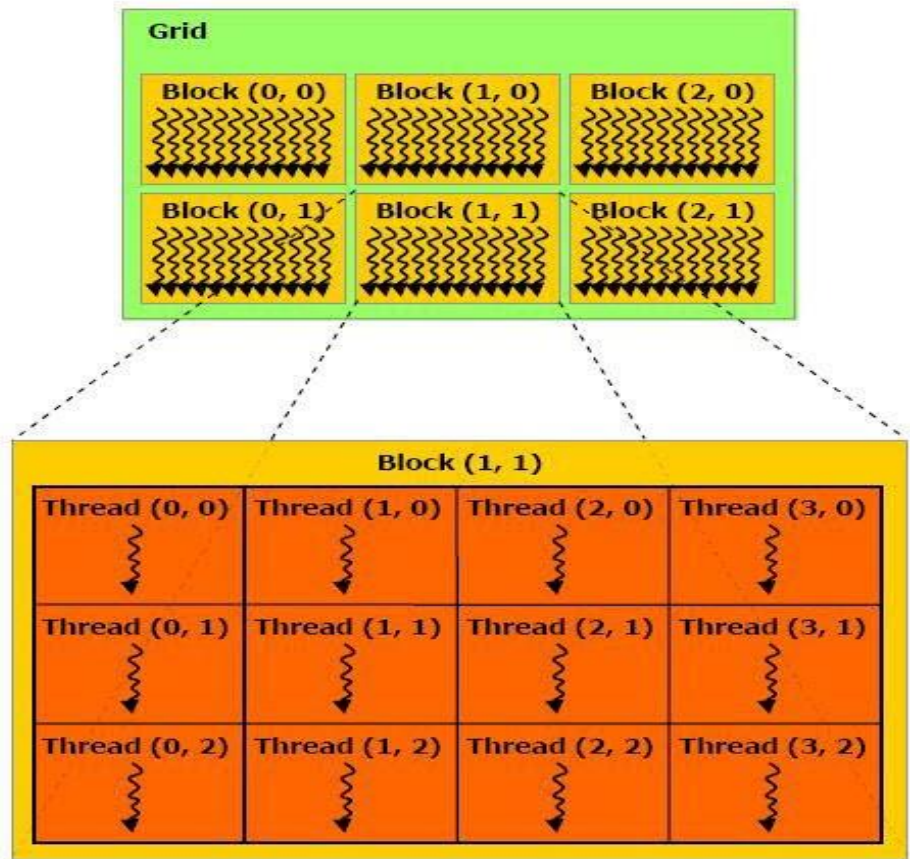
Threads of a block are expected to reside on the same processor core and must share the limited memory resources of that core. On current GPUs, a thread block may contain up to 1024 threads.

However, a kernel can be executed by multiple equally-shaped thread blocks, so that the total number of threads is equal to the number of threads per block times the number of blocks.

Blocks are organized into a one-dimensional, two-dimensional, or three-dimensional grid of thread blocks as illustrated by **Figures** the number of thread blocks in a grid is usually dictated by the size of the data being processed or the number of processors in the system, which it can greatly exceed.

The number of threads per block and the number of blocks per grid specified in the <<<...>>> syntax can be of type int or dim3. Two-dimensional blocks or grids can be specified as in the example above.

Each block within the grid can be identified by a one-dimensional, two-dimensional, or three-dimensional index accessible within the kernel through the built-in blockIdx variable. The dimension of the thread block is accessible within the kernel through the built-in blockDim variable.



Extending the previous MatAdd() example to handle multiple blocks, the code becomes as follows.

A thread block size of 16x16 (256 threads), although arbitrary in this case, is a common choice. The grid is created with enough blocks to have one thread per matrix element as before. For simplicity, this example assumes that the number of threads per grid in each dimension is evenly divisible by the number of threads per block in that dimension, although that need not be the case.

// Kernel definition

```

__global__ void MatAdd(float A[N][N], float B[N][N], float C[N][N])
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;

    int j = blockIdx.y * blockDim.y + threadIdx.y;

    if (i < N && j < N)
        C[i][j] = A[i][j] + B[i][j];
}

int main()
{
    ...
    // Kernel invocation
    dim3 threadsPerBlock(16, 16);

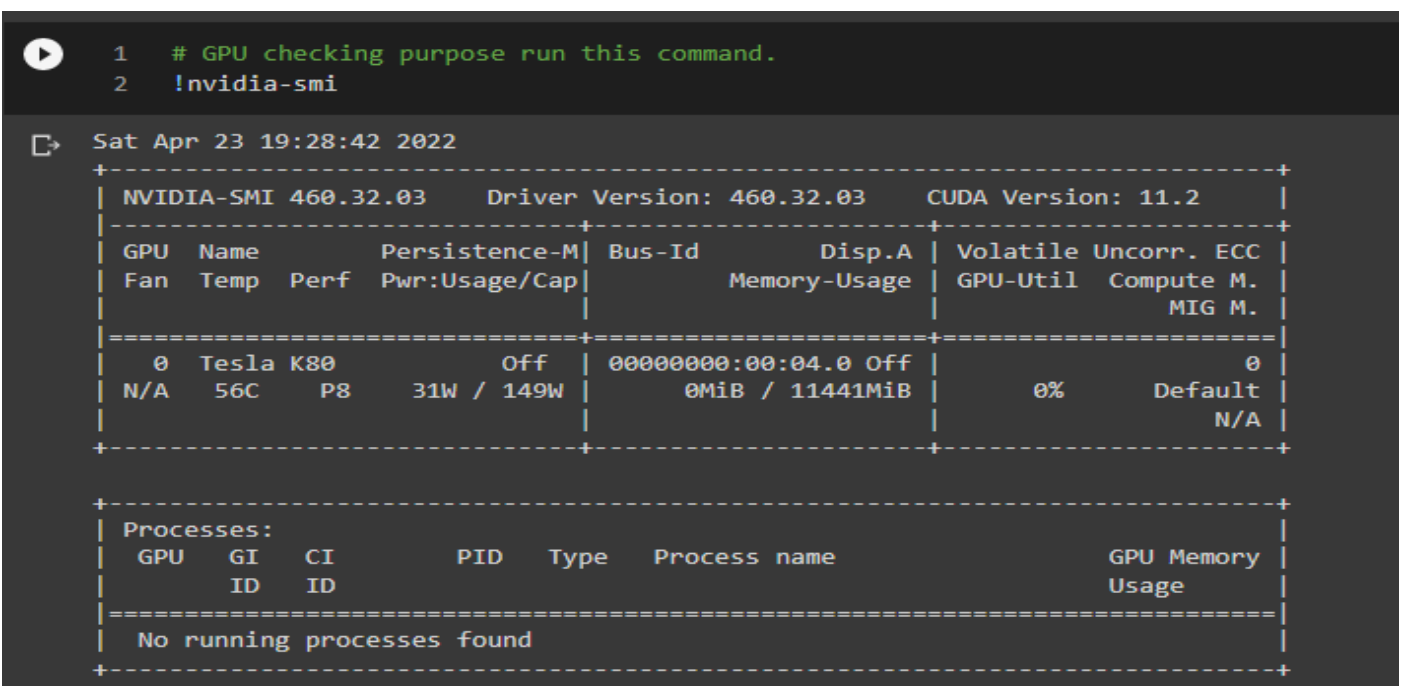
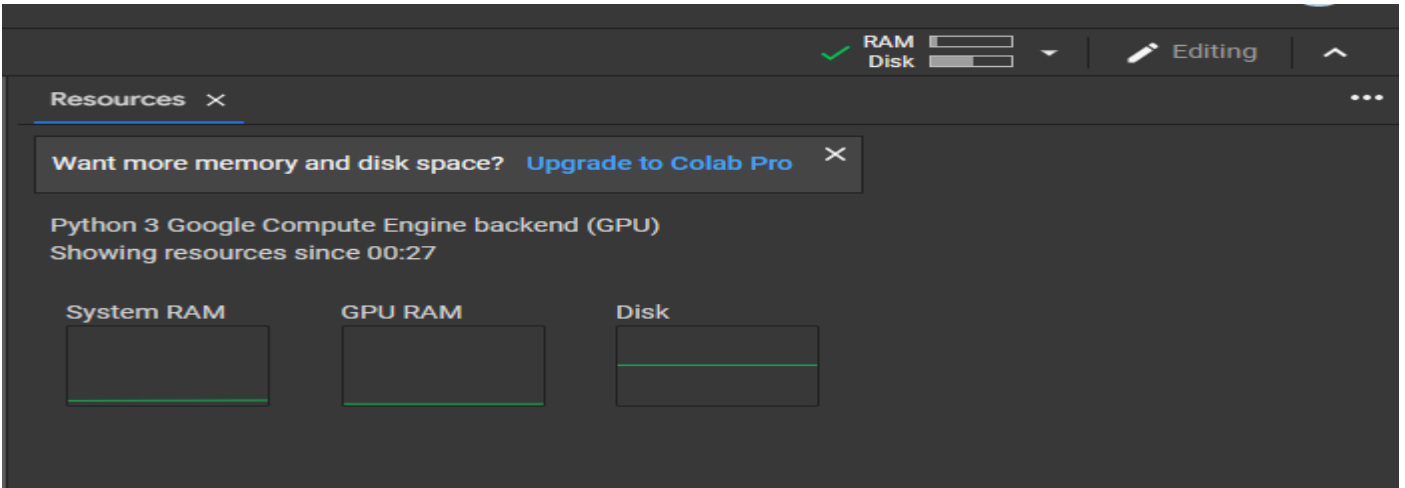
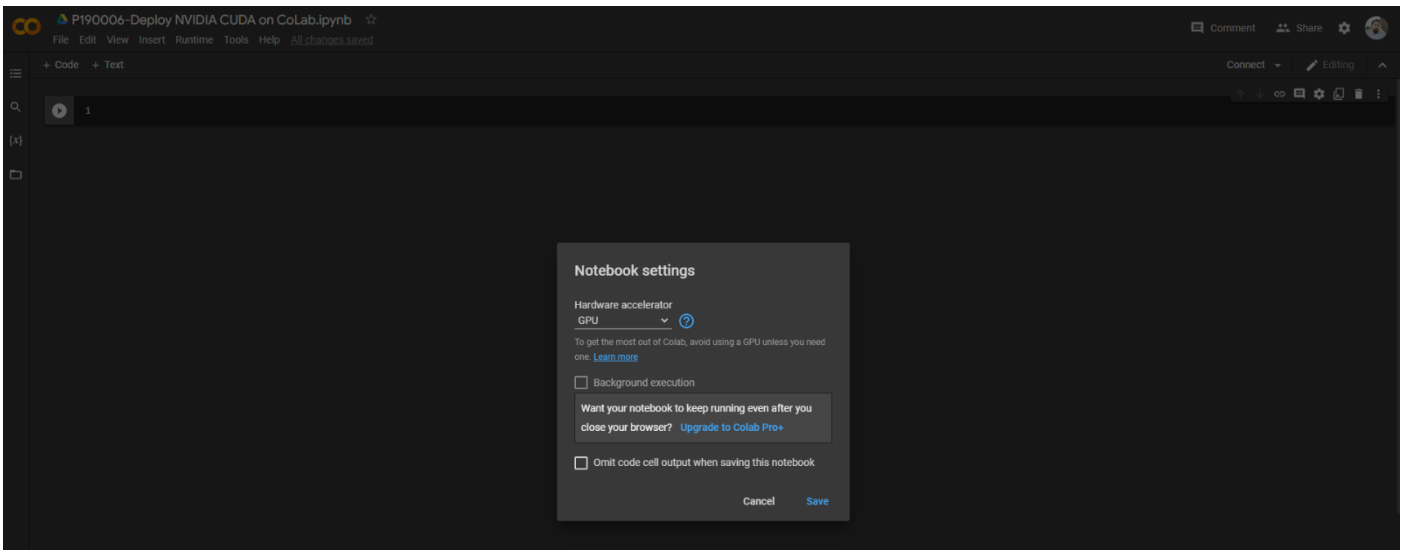
    dim3 numBlocks(N / threadsPerBlock.x, N / threadsPerBlock.y);

    MatAdd<<<numBlocks, threadsPerBlock>>>(A, B, C);
    ...
}

```

# CUDA

## Deploy the NVIDIA CUDA on Google Colab (Using GPU)



```
[2] 1 # Refresh the Cloud Instance of CUDA on Server
2
3 !apt-get --purge remove cuda nvidia* libnvidia-*
4 !dpkg -l | grep cuda- | awk '{print $2}' | xargs -n1 dpkg --purge
5 !apt-get remove cuda-*
6 !apt autoremove
7 !apt-get update
```

```
Reading package lists... Done
Building dependency tree
Reading state information... Done
Note, selecting 'nvidia-kernel-common-418-server' for glob 'nvidia*'
Note, selecting 'nvidia-325-updates' for glob 'nvidia*'
Note, selecting 'nvidia-346-updates' for glob 'nvidia*'
Note, selecting 'nvidia-driver-binary' for glob 'nvidia*'
```

```
Get:19 http://ppa.launchpad.net/c2d4u.team/c2d4u4.0+/ubuntu bionic/main Sources [1,950 kB]
Get:20 http://security.ubuntu.com/ubuntu bionic-security/restricted amd64 Packages [907 kB]
Get:21 http://security.ubuntu.com/ubuntu bionic-security/universe amd64 Packages [1,496 kB]
Get:22 http://security.ubuntu.com/ubuntu bionic-security/main amd64 Packages [2,728 kB]
Get:23 http://ppa.launchpad.net/c2d4u.team/c2d4u4.0+/ubuntu bionic/main amd64 Packages [999 kB]
Fetched 14.7 MB in 4s (3,899 kB/s)
Reading package lists... Done
```

```
[5] 1 # Install CUDA Version 9
2
3 !wget https://developer.nvidia.com/compute/cuda/9.2/Prod/local\_installers/cuda-repo-ubuntu1604-9-2-local\_9.2.88-1\_amd64 -O cuda-repo-ubuntu1604-9-2-local_9.2.88-1_amd64.deb
4 !dpkg -i cuda-repo-ubuntu1604-9-2-local_9.2.88-1_amd64.deb
5 !apt-key add /var/cuda-repo-9-2-local/7fa2af80.pub
6 !apt-get update
7 !apt-get install cuda-9.2
```

depmod...

```
DKMS: install completed.
Setting up cuda-cuobjdump-9-2 (9.2.88-1) ...
Setting up x11-utils (7.7+3build1) ...
Setting up cuda-cusparse-9-2 (9.2.88-1) ...
Setting up libxfont2:amd64 (1:2.0.3-1) ...
Setting up cuda-nvgraph-9-2 (9.2.88-1) ...
Setting up xfonts-utils (1:7.7+6) ...
Setting up cuda-gpu-library-advisor-9-2 (9.2.88-1) ...
Setting up fakeroot (1.22-2ubuntu1) ...
update-alternatives: using /usr/bin/fakeroot-sysv to provide /usr/bin/fakeroot (fakeroot) in auto mode
Setting up cuda-gdb-9-2 (9.2.88-1) ...
```

```
Processing triggers for dbus (1.12.2-1ubuntu1.2) ...
Processing triggers for hicolor-icon-theme (0.17-2) ...
Processing triggers for fontconfig (2.12.6-0ubuntu2) ...
Processing triggers for mime-support (3.60ubuntu1) ...
Processing triggers for libc-bin (2.27-3ubuntu1.3) ...
/sbin/ldconfig.real: /usr/local/lib/python3.7/dist-packages/ideep4py/lib/libmkldnn.so.0 is not a symbolic link
```

```
1 # Check the version of CUDA.
2
3 !nvcc --version
```

```
nvcc: NVIDIA (R) Cuda compiler driver
Copyright (c) 2005-2018 NVIDIA Corporation
Built on Wed Apr 11 23:16:29 CDT 2018
Cuda compilation tools, release 9.2, V9.2.88
```

```
[8] 1 # Execute the given command to Install a Small Extension to Run NVCC from Notebook cells.  
2  
3 !pip install git+https://github.com/andreinechaev/nvcc4jupyter.git
```

```
Collecting git+https://github.com/andreinechaev/nvcc4jupyter.git  
Cloning https://github.com/andreinechaev/nvcc4jupyter.git to /tmp/pip-req-build-7ujcw79x  
Running command git clone -q https://github.com/andreinechaev/nvcc4jupyter.git /tmp/pip-req-build-7ujcw79x  
Building wheels for collected packages: NVCCPlugin  
Building wheel for NVCCPlugin (setup.py) ... done  
Created wheel for NVCCPlugin: filename=NVCCPlugin-0.0.2-py3-none-any.whl size=4306 sha256=2ef1f55e5495dafa31bd8b629125fb80c1d766ddc7917a4143c64f92c9ec863  
Stored in directory: /tmp/pip-ephem-wheel-cache-s6wvfovl/wheels/ca/33/8d/3c86eb85e97d2b6169d95c6e8f2c297fdec60db6e84cb56f5e  
Successfully built NVCCPlugin  
Installing collected packages: NVCCPlugin  
Successfully installed NVCCPlugin-0.0.2
```

```
[10] 1 # Load the Extension using this code.  
2  
3 %load_ext nvcc_plugin
```

```
created output directory at /content/src  
Out bin /content/result.out
```

```
1 %%cu  
2 #include <stdio.h>  
3 #include <stdlib.h>  
4 global void add(int *a, int *b, int *c) {  
5     *c = *a + *b;  
6 }  
7 int main() {  
8     int a, b, c;  
9     // host copies of variables a, b & c  
10    int *d_a, *d_b, *d_c;  
11    // device copies of variables a, b & c  
12    int size = sizeof(int);  
13    // Allocate space for device copies of a, b, c  
14    cudaMalloc((void **)&d_a, size);  
15    cudaMalloc((void **)&d_b, size);  
16    cudaMalloc((void **)&d_c, size);  
17    // Setup input values  
18    c = 0;  
19    a = 3;  
20    b = 5;  
21    // Copy inputs to device  
22    cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);  
23    cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);  
24    // Launch add() kernel on GPU  
25    add<<<1,1>>>(d_a, d_b, d_c);  
26    // Copy result back to host  
27    cudaError err = cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);  
28    if(err!=cudaSuccess) {  
29        printf("CUDA error copying to Host: %s\n", cudaGetErrorString(err));  
30    }  
31    printf("result is %d\n",c);  
32    // Cleanup  
33    cudaFree(d_a);  
34    cudaFree(d_b);  
35    cudaFree(d_c);  
36    return 0;  
37 }
```

```
result is 8
```

# CUDA

## Matrix Multiplicaton

### Matrix Multiplication

```
1  %%cu
2
3  #include <stdio.h>
4  #include <math.h>
5  #define TILE_WIDTH 2
6
7  /*matrix multiplication kernels*/
8  //non shared
9
10 global void MatrixMul( float *Md , float *Nd , float *Pd , const int WIDTH )
11 {
12     // calculate thread id
13
14     unsigned int col = TILE_WIDTH*blockIdx.x + threadIdx.x ;
15     unsigned int row = TILE_WIDTH*blockIdx.y + threadIdx.y ;
16     for (int k = 0 ; k<WIDTH ; k++ )
17     {
18         Pd[row*WIDTH + col]+= Md[row * WIDTH + k ] * Nd[ k * WIDTH + col] ;
19     }
20 }
21
22 // shared
23 global void MatrixMulSh( float *Md , float *Nd , float *Pd , const int WIDTH )
24 {
25     //Taking shared array to break the Matrix in Tile width and fetch them in that array per ele
26
27     __shared__ float Mds [TILE_WIDTH][TILE_WIDTH] ;
28     __shared__ float Nds [TILE_WIDTH][TILE_WIDTH] ;
29
30     // calculate thread id
31     unsigned int col = TILE_WIDTH*blockIdx.x + threadIdx.x ;
32     unsigned int row = TILE_WIDTH*blockIdx.y + threadIdx.y ;
33
```

```

34 for (int m = 0 ; m<WIDTH/TILE_WIDTH ; m++ ) // m indicate number of phase
35 {
36     Mds[threadIdx.y][threadIdx.x] = Md[row*WIDTH + (m*TILE_WIDTH + threadIdx.x)] ;
37     Nds[threadIdx.y][threadIdx.x] = Nd[ ( m*TILE_WIDTH + threadIdx.y ) * WIDTH + col] ;
38     __syncthreads() ; // for synchronizeing the threads
39
40     // Do for tile
41     for ( int k = 0 ; k<TILE_WIDTH ; k++ )
42         Pd[row*WIDTH + col]+= Mds[threadIdx.x][k] * Nds[k][threadIdx.y] ;
43     __syncthreads() ; // for synchronizeing the threads
44
45 }
46 }
47
48 // main routine
49 int main ()
50 {
51     const int WIDTH = 6 ;
52     float array1_h[WIDTH][WIDTH] ,array2_h[WIDTH][WIDTH],
53         result_array_h[WIDTH][WIDTH] ,M result_array_h[WIDTH][WIDTH] ;
54     float *array1_d , *array2_d ,*result_array_d ,*M_result_array_d ; // device array
55     int i , j ;
56     //input in host array
57     for ( i = 0 ; i<WIDTH ; i++ )
58     {
59         for ( j = 0 ; j<WIDTH ; j++ )
60         {
61             array1_h[i][j] = 1 ;
62             array2_h[i][j] = 2 ;
63         }
64     }
65
66     //create device array cudaMalloc ( (void **)&array_name, sizeofmatrixinbytes) ;
67     cudaMalloc((void **) &array1_d , WIDTH*WIDTH*sizeof (int) ) ;
68     cudaMalloc((void **) &array2_d , WIDTH*WIDTH*sizeof (int) ) ;
69
70     //copy host array to device array; cudaMemcpy ( dest , source , WIDTH , direction )
71     cudaMemcpy ( array1_d , array1_h , WIDTH*WIDTH*sizeof (int) , cudaMemcpyHostToDevice ) ;
72     cudaMemcpy ( array2_d , array2_h , WIDTH*WIDTH*sizeof (int) , cudaMemcpyHostToDevice ) ;
73
74     //allocating memory for resultent device array
75     cudaMalloc((void **) &result_array_d , WIDTH*WIDTH*sizeof (int) ) ;
76     cudaMalloc((void **) &M_result_array_d , WIDTH*WIDTH*sizeof (int) ) ;
77

```





# CUDA

## Vector Addition

### Vector Addition CUDA

```
1  %%CU
2  #include <stdio.h>
3
4  #define HANDLE_ERROR( err ) ( HandleError( err, __FILE__, __LINE__ ) )
5
6  static void HandleError( cudaError_t err, const char *file, int line )
7  {
8      if (err != cudaSuccess)
9      {
10         printf( "%s in %s at line %d\n", cudaGetErrorString( err ),
11             file, line );
12         exit( EXIT_FAILURE );
13     }
14 }
15
16 const short N = 10 ;
17
18 // CUDA Kernel for Vector Addition
19 __global__ void Vector Addition ( const int *dev_a , const int *dev_b , int *dev_c )
20 {
21     //Get the id of thread within a block
22     unsigned short tid = blockIdx.x ;
23
24     if ( tid < N ) // check the boundry condition for the threads
25         dev_c [tid] = dev_a[tid] + dev_b[tid] ;
26
27 }
28
```

```

29 int main (void)
30 {
31     //Host array
32     int Host_a[N], Host_b[N], Host_c[N];
33
34     //Device array
35     int *dev_a , *dev_b, *dev_c ;
36
37     //Allocate the memory on the GPU
38     HANDLE_ERROR ( cudaMalloc((void **)&dev_a , N*sizeof(int) ) );
39     HANDLE_ERROR ( cudaMalloc((void **)&dev_b , N*sizeof(int) ) );
40     HANDLE_ERROR ( cudaMalloc((void **)&dev_c , N*sizeof(int) ) );
41
42     //fill the Host array with random elements on the CPU
43     for ( int i = 0; i <N ; i++ )
44     {
45         Host_a[i] = -i ;
46         Host_b[i] = i*i ;
47     }
48
49     //Copy Host array to Device array
50     HANDLE_ERROR ( cudaMemcpy (dev_a , Host_a , N*sizeof(int) , cudaMemcpyHostToDevice));
51     HANDLE_ERROR ( cudaMemcpy (dev_b , Host_b , N*sizeof(int) , cudaMemcpyHostToDevice));
52
53     //Make a call to GPU kernel
54     Vector_Addition <<< N, 1 >>> (dev_a , dev_b , dev_c ) ;
55
56     //Copy back to Host array from Device array
57     HANDLE_ERROR ( cudaMemcpy(Host_c , dev_c , N*sizeof(int) , cudaMemcpyDeviceToHost));
58
59     //Display the result
60     for ( int i = 0; i<N; i++ )
61         printf ("%d + %d = %d\n", Host_a[i] , Host_b[i] , Host_c[i] ) ;
62     //Free the Device array memory
63     cudaFree (dev_a) ;
64     cudaFree (dev_b) ;
65     cudaFree (dev_c) ;
66
67     system("pause");
68     return 0 ;
69 }

```

```

sh: 1: pause: not found
0 + 0 = 0
-1 + 1 = 0
-2 + 4 = 2
-3 + 9 = 6
-4 + 16 = 12
-5 + 25 = 20
-6 + 36 = 30
-7 + 49 = 42
-8 + 64 = 56
-9 + 81 = 72

```