# GPU accelerated computing for Finite Element Method
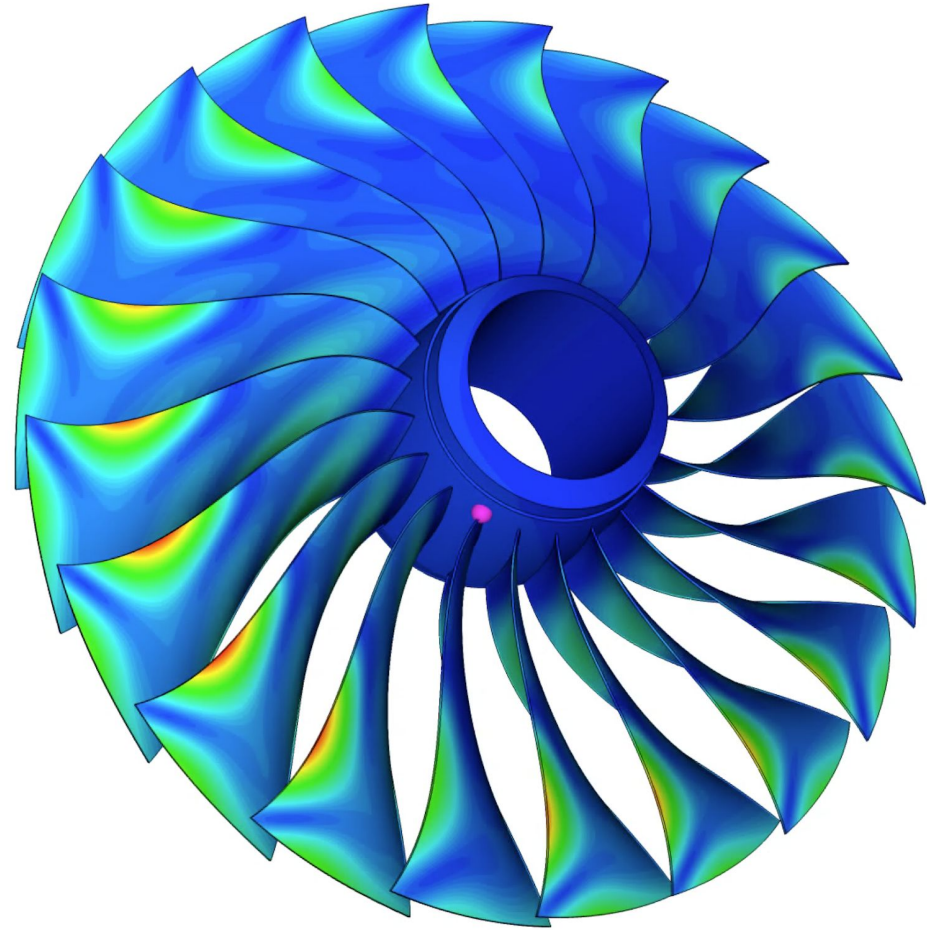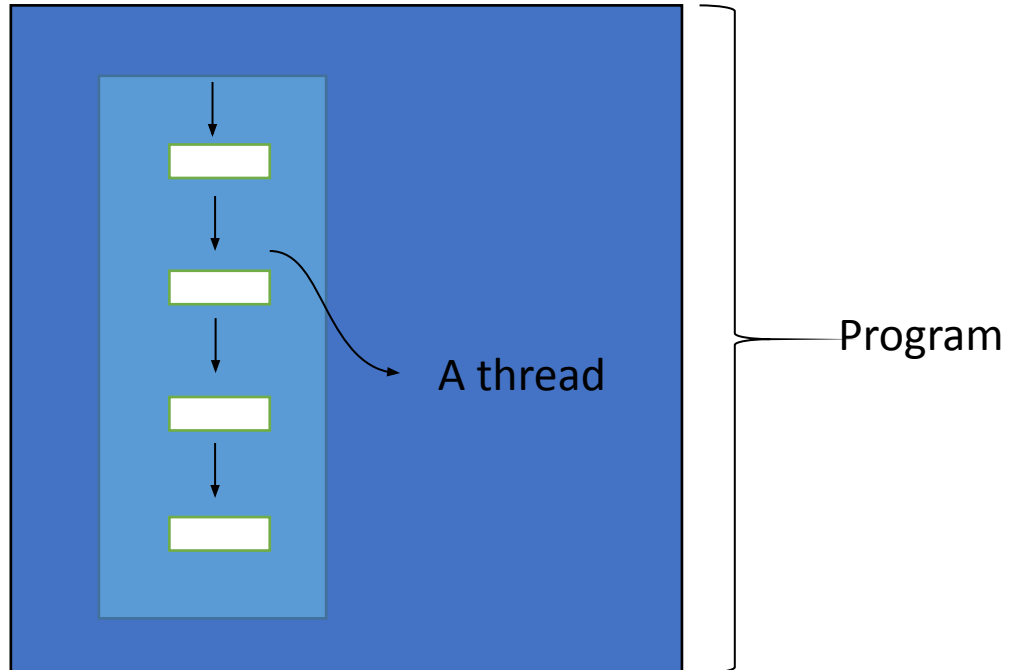
## GPU programming - Basics

Thread

- A thread is a single sequential flow of control within a program.

- A sequential code in one processor has one thread.

**GPU programming -** C library functions

Pointer

- A variable that points to the storage/memory address of another variable.
  - A variable of type certain type will store a value

$$\text{int } v = 0;$$

  - This variable has its address (where it is located the memory). This address can be obtained by using '**&**'

$$\&v$$

  - A pointer stores the address of the variable

$$\text{int } *y = \&v;$$

  - The value of the variable can be accessed using the variable or the pointer
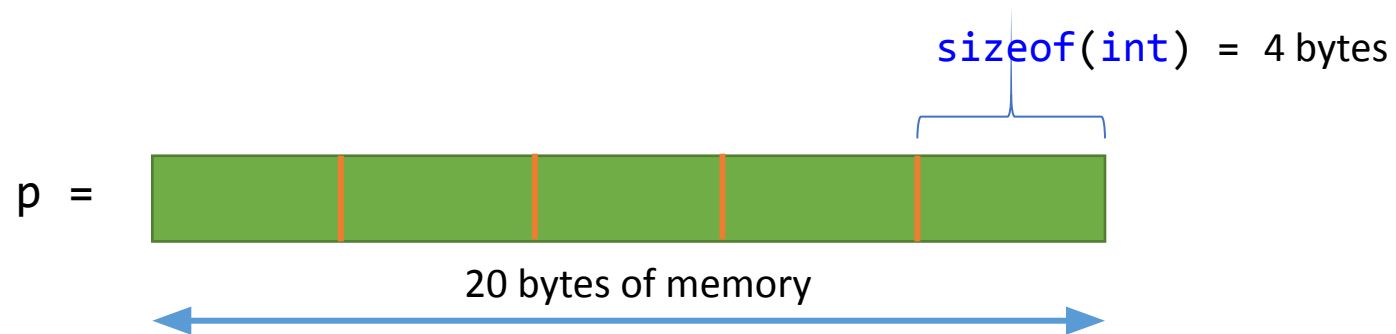
$$v$$
$$*y$$

malloc()

- Dynamically allocates a single large block of memory
    - Syntax

```
pointer = (type*) malloc(byte size)
```

    - Example

```
n = 5;
int *p;
p = (int*)malloc(n * sizeof(int));
```

sizeof(int) = 4 bytes

p =



20 bytes of memory

# Introduction to GPU

EXERCISE 0 : Vector Addition using C program

Source code
https://github.com/sivasanarul/FEMwithGPU/tree/master/EX1_vector_addition

vector_add.c

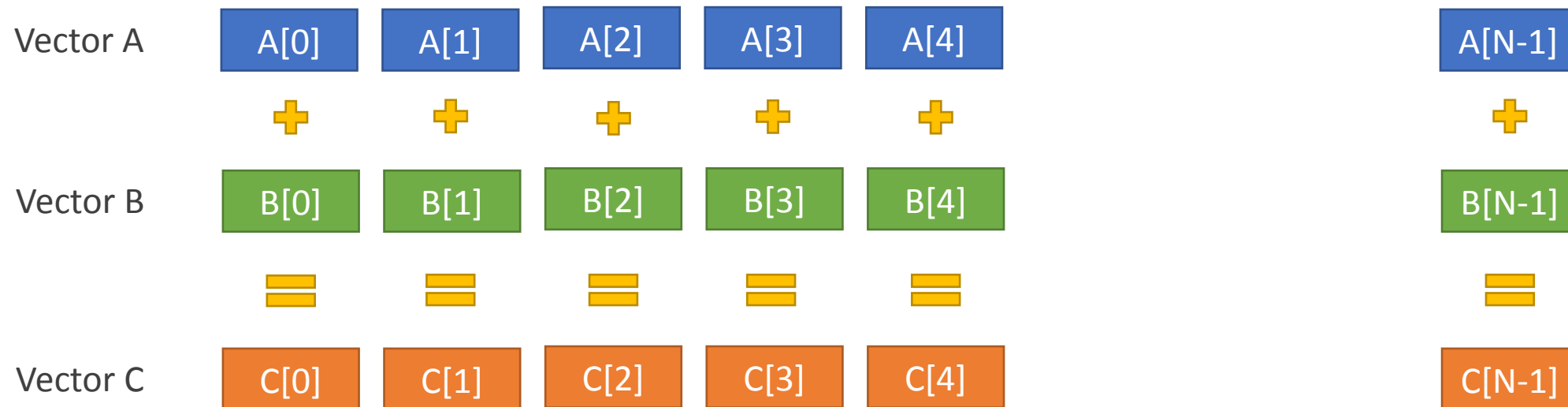# Introduction to GPU

Vector Addition : vector_addition()

$$A + B = C$$

**GPU programming -** C program for vector addition

C programming – Vector addition

```
6    float *a, *b, *c;
7
8    a    = (float*)malloc(sizeof(float) * N);
9    b    = (float*)malloc(sizeof(float) * N);
10   c    = (float*)malloc(sizeof(float) * N);
```

```
12   // Initialize array
13   for(int i = 0; i < N; i++){
14       a[i] = 1.0f; b[i] = 2.0f;
15   }
```

```
1    #include <stdio.h>
2    #include <stdlib.h>
3    #define N 100000000
4
5    int main(){
6        float *a, *b, *c;
7
8        a    = (float*)malloc(sizeof(float) * N);
9        b    = (float*)malloc(sizeof(float) * N);
10       c    = (float*)malloc(sizeof(float) * N);
11
12       // Initialize array
13       for(int i = 0; i < N; i++){
14           a[i] = 1.0f; b[i] = 2.0f;
15       }
16
17
18
19       clock_t t;
20       t = clock();
21       // Main function
22       for(int i=0;i<N;i++){
23           c[i] = a[i] + b[i];}
24       t = clock() - t;
25       double time_taken = ((double)t)/CLOCKS_PER_SEC; // in seconds
26       printf("fun() took %f seconds to execute \n", time_taken);
27
28       free(a);  free(b); free(c);
29
30   }
```

Memory Allocation for the variables

➡️

Initializing the variables

⬇️

```
21   // Main function
22   for(int i=0;i<N;i++){
23       c[i] = a[i] + b[i];}
```

```
28   free(a);  free(b); free(c);
```
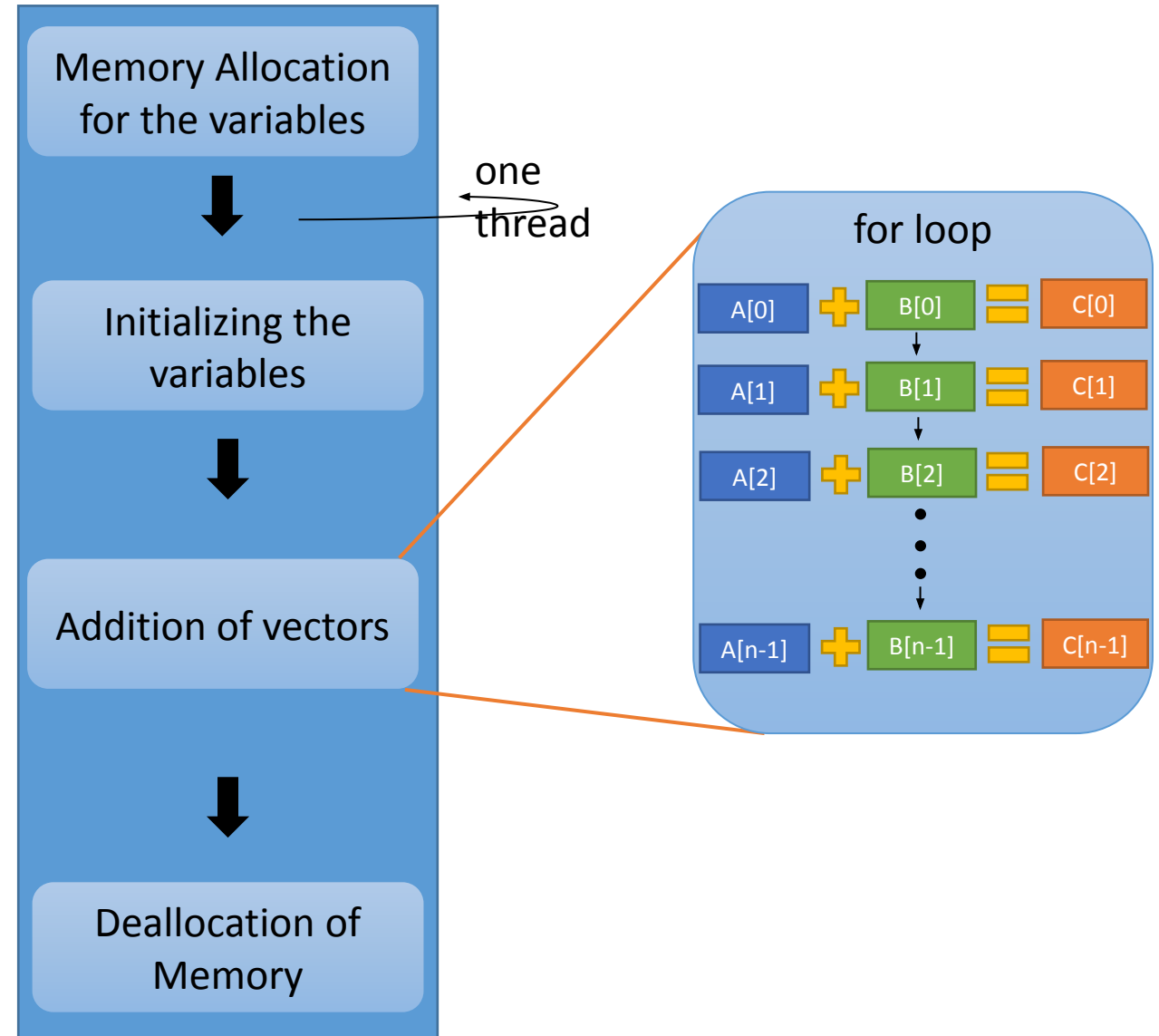
⬅️

Deallocation of Memory

Addition of vectors

# Introduction to GPU

C programming – Vector addition

```c
1   #include <stdio.h>
2   #include <stdlib.h>
3   #define N 100000000
4
5   int main(){
6       float *a, *b, *c;
7
8       a   = (float*)malloc(sizeof(float) * N);
9       b   = (float*)malloc(sizeof(float) * N);
10      c   = (float*)malloc(sizeof(float) * N);
11
12      // Initialize array
13      for(int i = 0; i < N; i++){
14          a[i] = 1.0f; b[i] = 2.0f;
15      }
16
17
18
19      clock_t t;
20      t = clock();
21      // Main function
22      for(int i=0;i<N;i++){
23          c[i] = a[i] + b[i];}
24      t = clock() - t;
25      double time_taken = ((double)t)/CLOCKS_PER_SEC; // in seconds
26      printf("fun() took %f seconds to execute \n", time_taken);
27
28      free(a);  free(b); free(c);
29
30  }
```
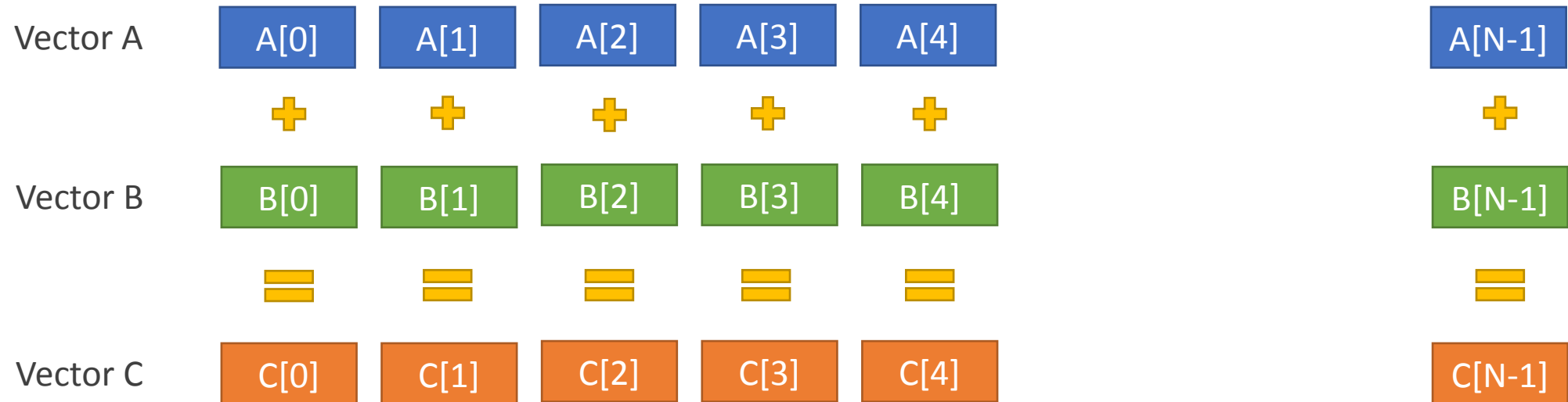
Serial Code



Memory Allocation for the variables

one thread

Initializing the variables

Addition of vectors

Deallocation of Memory

for loop

A[0] + B[0] = C[0]
A[1] + B[1] = C[1]
A[2] + B[2] = C[2]
⋮
A[n-1] + B[n-1] = C[n-1]

# Introduction to GPU

Vector Addition : vector_addition()

$$A + B = C$$



- Divide array into blocks  cuda
- threads into multiple blocks
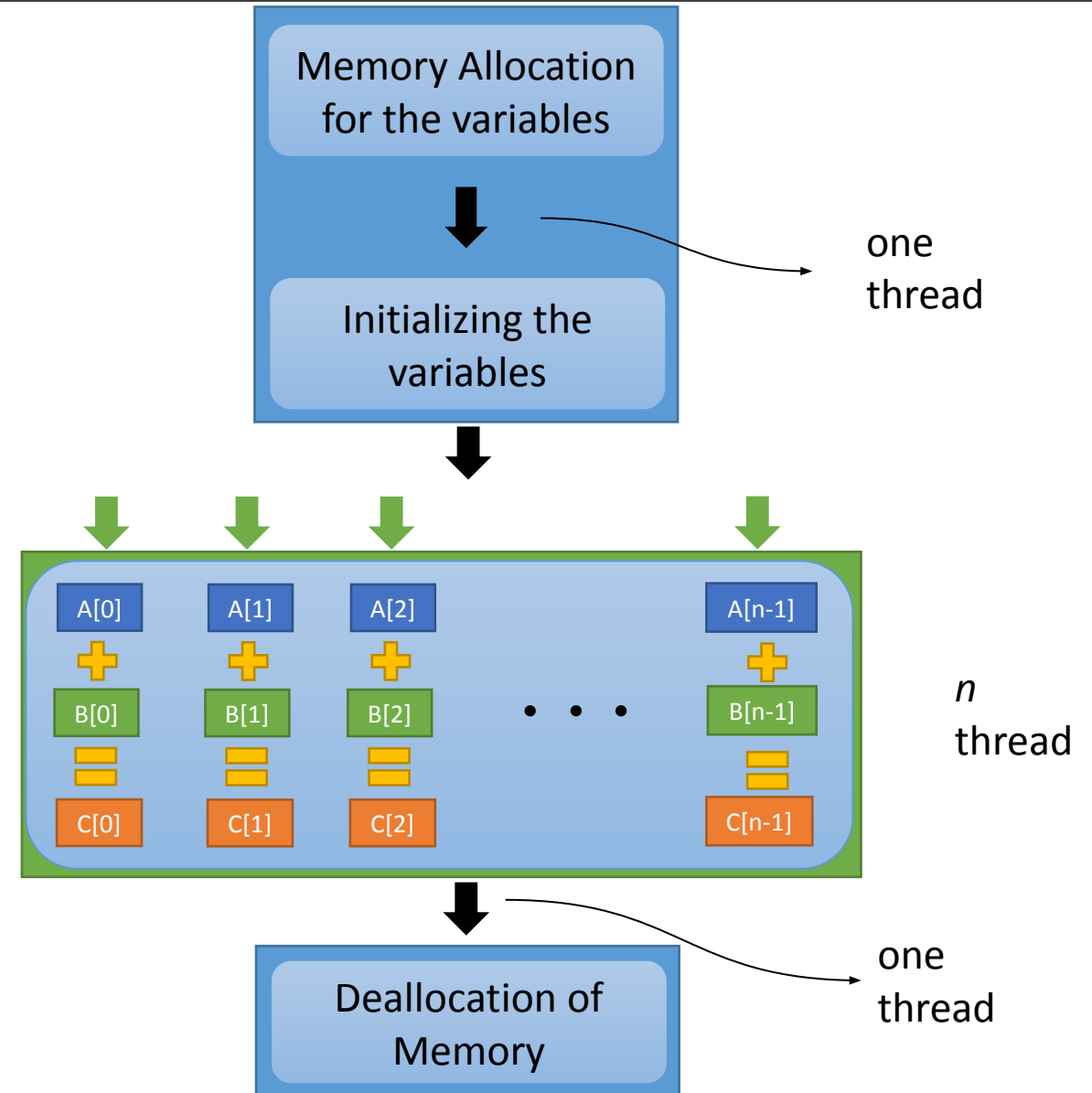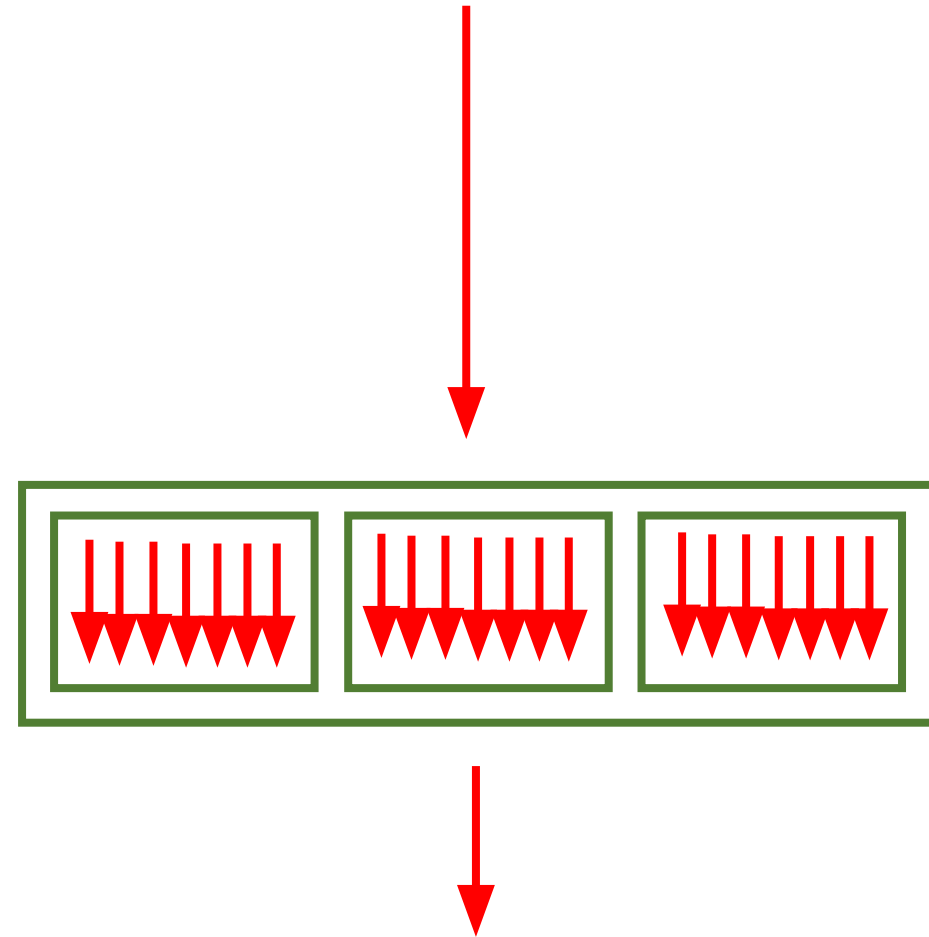
# Introduction to GPU

CUDA programming Model – Heterogenous Computing

```c
16  int main(){
17      float *a, *b, *out;
18      float *d_a, *d_b, *d_out;
19
20      a   = (float*)malloc(sizeof(float) * N);
21      b   = (float*)malloc(sizeof(float) * N);
22      out = (float*)malloc(sizeof(float) * N);
23
24      // Initialize array
25      for(int i = 0; i < N; i++){
26          a[i] = 1.0f;
27          b[i] = 2.0f;
28      }
29
30      // Allocate device memore for a
31      cudaMalloc((void**)&d_a,sizeof(float)*N);
32      cudaMalloc((void**)&d_b,sizeof(float)*N);
33      cudaMalloc((void**)&d_out,sizeof(float)*N);
34
35      // Transfer data from host to device memory
36      cudaMemcpy(d_a,a, sizeof(float)*N, cudaMemcpyHostToDevice);
37      cudaMemcpy(d_b,b, sizeof(float)*N, cudaMemcpyHostToDevice);
38
39      // Main function
40      int block_size = 256;
41      int grid_size  = (N+block_size)/block_size;
42      vector_add<<<grid_size,block_size>>>(d_out, d_a, d_b, N);
43
44      cudaMemcpy(out, d_out, sizeof(float)*N, cudaMemcpyDeviceToHost);
45
46      // Deallocate device memory
47      cudaFree(d_a);
48      cudaFree(d_b);
49      cudaFree(d_out);
50
51      // Deallocate host memory
52      free(a);
53      free(b);
54      free(out);
55
56  }
```

Serial Code

Parallel Code

Serial Code
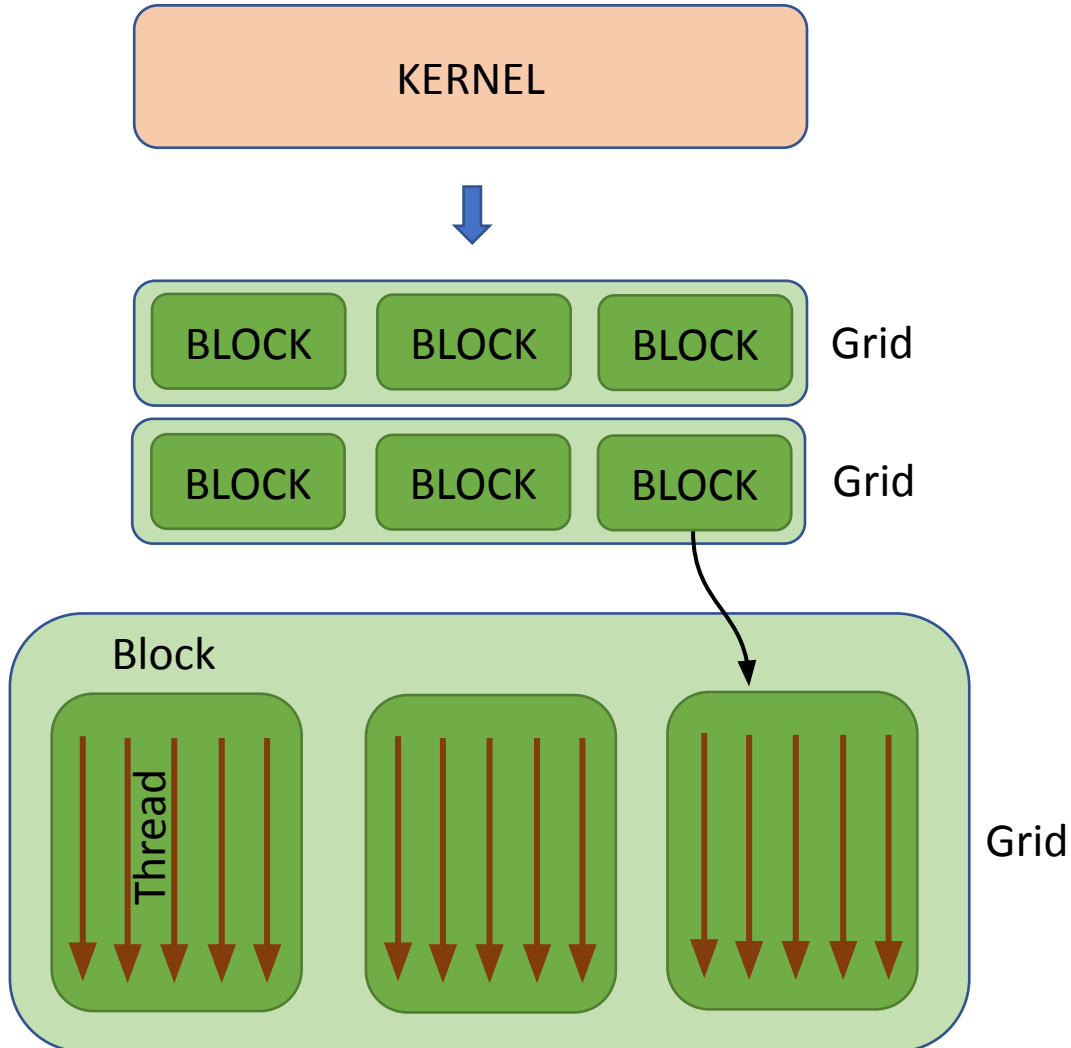


Memory Allocation for the variables

one thread

Initializing the variables

| A[0] | A[1] | A[2] | A[n-1] |
|------|------|------|--------|
| B[0] | B[1] | B[2] | B[n-1] |
| C[0] | C[1] | C[2] | C[n-1] |

$n$ thread

Deallocation of Memory

one thread

# Introduction to GPU

CUDA programming Model – Heterogenous Computing

```
16  int main(){
17      float *a, *b, *out;
18      float *d_a, *d_b, *d_out;
19
20      a   = (float*)malloc(sizeof(float) * N);
21      b   = (float*)malloc(sizeof(float) * N);
22      out = (float*)malloc(sizeof(float) * N);
23
24      // Initialize array
25      for(int i = 0; i < N; i++){
26          a[i] = 1.0f;
27          b[i] = 2.0f;
28      }
29
30      // Allocate device memore for a
31      cudaMalloc((void**)&d_a,sizeof(float)*N);
32      cudaMalloc((void**)&d_b,sizeof(float)*N);
33      cudaMalloc((void**)&d_out,sizeof(float)*N);
34
35      // Transfer data from host to device memory
36      cudaMemcpy(d_a,a, sizeof(float)*N, cudaMemcpyHostToDevice);
37      cudaMemcpy(d_b,b, sizeof(float)*N, cudaMemcpyHostToDevice);
38
39      // Main function
40      int block_size = 256;
41      int grid_size  = (N+block_size)/block_size;
42      vector_add<<<grid_size,block_size>>>(d_out, d_a, d_b, N);
43
44      cudaMemcpy(out, d_out, sizeof(float)*N, cudaMemcpyDeviceToHost);
45
46      // Deallocate device memory
47      cudaFree(d_a);
48      cudaFree(d_b);
49      cudaFree(d_out);
50
51      // Deallocate host memory
52      free(a);
53      free(b);
54      free(out);
55
56  }
```

Serial Code

Parallel Code

Serial Code

Serial Code in host

Parallel Code in device

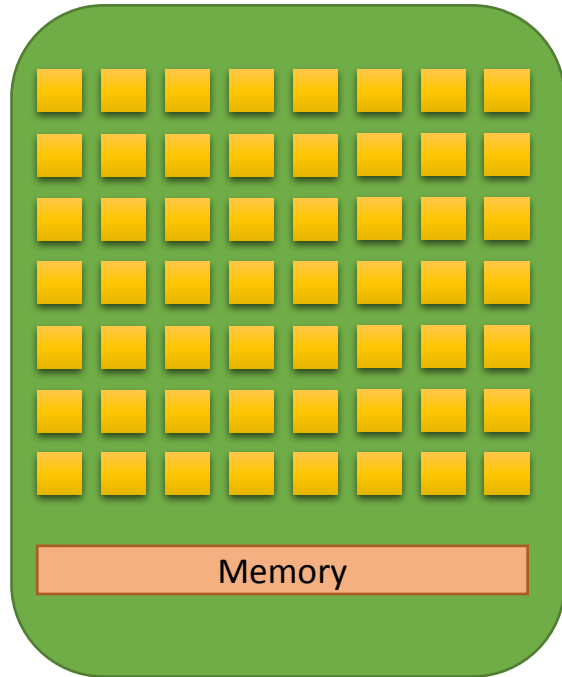Serial Code in host

# Introduction to GPU

A kernel is a function that is run in a GPU



- A kernel is split into grids
- A grid is split into blocks
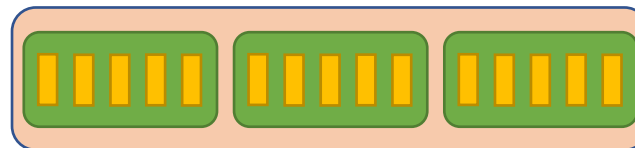- Each block is collection of threads

# Introduction to GPU
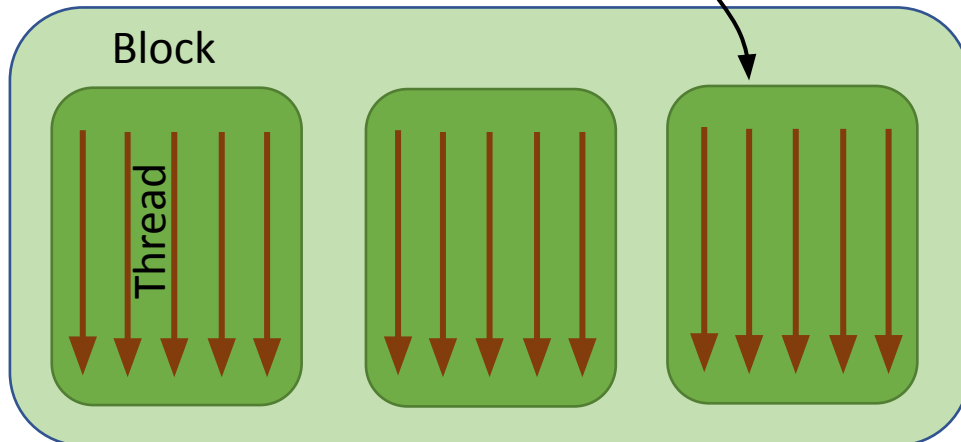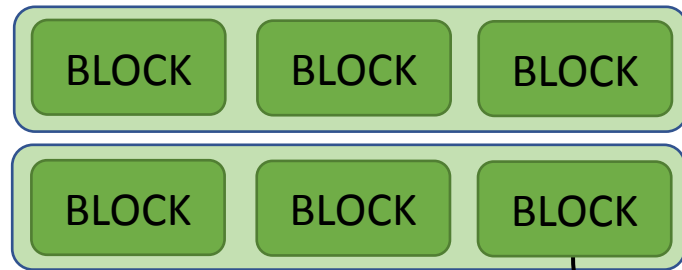
cores

Streaming Multiprocessor – collection of cores

CUDA-enabled GPU – Collection of Streaming Multiprocessor

Memory

# Introduction to GPU

A kernel is a function that is run in a GPU

KERNEL

BLOCK   BLOCK   BLOCK

BLOCK   BLOCK   BLOCK

Block

Thread

A kernel is executed in a CUDA-enabled GPU

KERNEL ➡

Multiple kernels can be executed in one CUDA-enabled GPU

A block is executed in a Streaming Multiprocessor

BLOCK ➡

Multiple blocks can be executed in one Streaming multiprocessor

Thread

A thread is executed in a core

➡

**Introduction of GPU programming**

Vector Addition : vector_addition()

- In one core as one thread
- In one streaming multiprocessor as one block
- In the entire GPU device as multiple blocks

# Introduction to GPU

Vector Addition : vector_addition()

CUDA-enabled GPU

In one core as one thread

for loop

A[0] + B[0] = C[0]

A[1] + B[1] = C[1]

A[2] + B[2] = C[2]

A[n-1] + B[n-1] = C[n-1]

```
10    __global__ void vector_add(float *out, float *a, float *b, int n){
11        for(int i=0;i<n;i++){
12            out[i] = a[i] + b[i];}
13    }
```

# Introduction to GPU

Vector Addition : vector_addition()

CUDA-enabled GPU

In one streaming multiprocessor (SM) as one block

- When the kernel is called, the number of blocks and the number of threads in each block is specified

$$\texttt{vector\_add <<<1,256>>> (d\_out, d\_a, d\_b, N)}$$

- The block runs in a stream multiprocessor (SM)
- The SM will have 256 thread.
- Each thread runs in a core.
- Each thread will have an unique id: `threadIdx.x`

`threadIdx.x`     0   1   2       254  255

**Introduction of GPU programming -** vector addition - one streaming multiprocessor

Vector Addition : vector_addition()

threadIdx.x    0     1     2               255

blockDim.x      256

All the threads in the cores run the same function

```
10    __global__ void vector_add(float *out, float *a, float *b, int n){
11         int index = threadIdx.x;
12         int stride = blockDim.x;
13
14         for(int i=index;i<n;i+=stride){
15             out[i] = a[i] + b[i];}
16    }
```
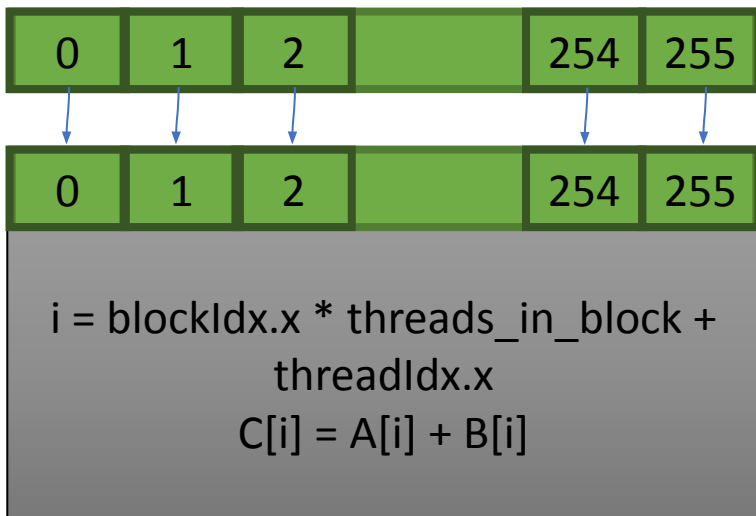
# Introduction to GPU

Vector Addition : vector_addition()



```
10    __global__ void vector_add(float *out, float *a, float *b, int n){
11        int index = threadIdx.x;
12        int stride = blockDim.x;
13
14        for(int i=index;i<n;i+=stride){
15            out[i] = a[i] + b[i];}
16    }
```

# Introduction to GPU

Vector Addition Example : vector_addition()

| Number of blocks : N | blockIdx.x   : 0, 1, 2, 3 … N − 1 |
| --- | --- |
| Number of threads in each block : 256 | threadIdx.x : 0, 1, 2, 3 … 255 |

### Block 0

| 0 | 1 | 2 | | 254 | 255 |
| --- | --- | --- | --- | --- | --- |

| 0 | 1 | 2 | | 254 | 255 |
| --- | --- | --- | --- | --- | --- |

i = blockIdx.x * threads_in_block + threadIdx.x
C[i] = A[i] + B[i]

### Block 1

| 0 | 1 | 2 | | 254 | 255 |
| --- | --- | --- | --- | --- | --- |

| 256 | 257 | 258 | | 510 | 511 |
| --- | --- | --- | --- | --- | --- |

i = blockIdx.x * threads_in_block + threadIdx.x
C[i] = A[i] + B[i]

### Block N-1

| 0 | 1 | 2 | | 254 | 255 |
| --- | --- | --- | --- | --- | --- |

| … | … | … | | … | … |
| --- | --- | --- | --- | --- | --- |

i = blockIdx.x * threads_in_block + threadIdx.x
C[i] = A[i] + B[i]

CUDA programming Model – Heterogenous Computing

Host – CPU, Device - GPU

# Introduction to GPU

- Functions that run on GPU are usually enclosed in "<<<   >>>".

- It has extension '.cu'.

- Complied using nvcc compiler driver.

# Introduction to GPU

The host device CPU manages memory

- It allocates the memory in the GPU
    - cudaMalloc (void **pointer, size_t nbytes)
    - cudaMemset (void *point, int value, size_t count)

- Data is transferred from host memory to device memory
    - cudaMemcpy(device_variable, host_variable, size_of_variable, CudaMemcpyHosttoDevice)

- After the kernel execution and data is transferred from device to host memory
    - cudaFree (void *pointer)

- Finally, the memory is deallocated
    - cudaMemcpy(host_variable, device_variable, size_of_variable, CudaMemcpyDevicetoHost)

# Introduction to GPU

EXERCISE 1 : Vector Addition using GPU program

Source code:
https://github.com/sivasanarul/FEMwithGPU/tree/master/EX1_vector_addition

**Introduction of GPU**

Code profiling – nvprof ./####

1 thread

```
==28571== Profiling application: ./vector_add_onethread
==28571== Profiling result:
            Type  Time(%)      Time     Calls       Avg       Min       Max  Name
 GPU activities:   94.28%   14.4019s         1   14.4019s  14.4019s  14.4019s  vector_add(float*, float*, float*, int)
                    3.14%   479.91ms         2   239.95ms  239.17ms  240.74ms  [CUDA memcpy HtoD]
                    2.58%   394.51ms         1   394.51ms  394.51ms  394.51ms  [CUDA memcpy DtoH]
```

1 block with 256 threads

```
==28894== Profiling application: ./vector_add_oneblock
==28894== Profiling result:
            Type  Time(%)      Time     Calls       Avg       Min       Max  Name
 GPU activities:   38.34%   478.81ms         2   239.40ms  239.19ms  239.62ms  [CUDA memcpy HtoD]
                   30.86%   385.37ms         1   385.37ms  385.37ms  385.37ms  [CUDA memcpy DtoH]
                   30.81%   384.72ms         1   384.72ms  384.72ms  384.72ms  vector_add(float*, float*, float*, int)
```
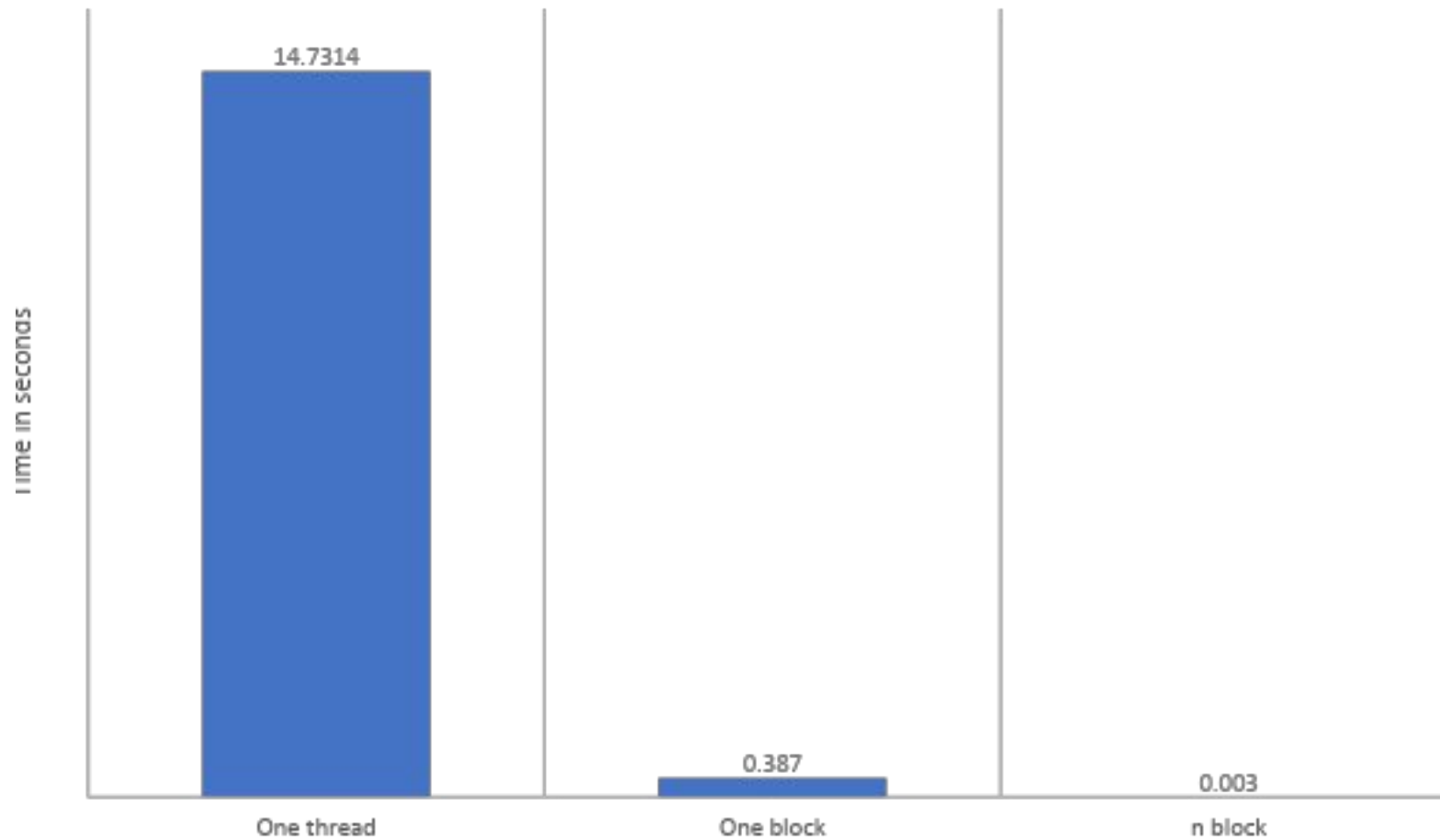
N blocks with all threads

```
==29270== Profiling application: ./vector_add_nblock
==29270== Profiling result:
            Type  Time(%)      Time     Calls       Avg       Min       Max  Name
 GPU activities:   55.40%   474.87ms         2   237.43ms  237.09ms  237.77ms  [CUDA memcpy HtoD]
                   44.14%   378.35ms         1   378.35ms  378.35ms  378.35ms  [CUDA memcpy DtoH]
                    0.45%   3.8855ms         1   3.8855ms  3.8855ms  3.8855ms  vector_add(float*, float*, float*, int)
```

# Introduction to GPU

Vector Addition Example : vector_addition()

# Introduction to GPU

**Introduction of GPU**

Code profiling

```
==29270== Profiling application: ./vector_add_nblock
==29270== Profiling result:
           Type  Time(%)      Time     Calls       Avg       Min       Max  Name
 GPU activities:   55.40%  474.87ms         2  237.43ms  237.09ms  237.77ms  [CUDA memcpy HtoD]
                   44.14%  378.35ms         1  378.35ms  378.35ms  378.35ms  [CUDA memcpy DtoH]
                    0.45%  3.8855ms         1  3.8855ms  3.8855ms  3.8855ms  vector_add(float*, float*, float*, int)
```

Expensive step is the memory transfer

To circumvent:

 Perform  simultaneous data transfer and computation

 Some streaming multiprocessors perform computation, some perform data transfer