

# **EC-444 Parallel & Distributed Computing**

**BS (CE)-2020**

## **Semester Project Report**



### **Submitted By:**

<b>Reg. No.</b>	<b>Name</b>
<b>20-CE-035</b>	<b>Abdullah Javed</b>

**Department of Computer Engineering**

**HITEC University Taxila**

---

**Submitted To:**

**Engr. Fasih Ahmad**

## Contents

Implementation of the 2D-Convolution in CUDA C .....	4
Problem Statement:.....	4
Tools Used:.....	4
Algorithm: .....	5
Parallelization Strategy: .....	7
Version 1 (Serial):.....	7
Version 2 (Open MP):.....	7
Version 3 (CUDA): .....	7
Source Code: .....	8
Version 1 (Serial):.....	8
Version 2 (OpenMP):.....	10
Version 3 (CUDA): .....	12
Detailed Profile & Analysis:.....	14
Version 1: .....	14
Version 2: .....	14
Version 3: .....	15
Comparison Between Version 1 & Version 2:.....	15
Comparison Between Version 2 & Version 3:.....	16
Execution Time & Speed up Results: .....	17
Speedup From Version 1 to Version 2: .....	17
Speedup From Version 2 to Version 3: .....	17
Speedup from Serial to CUDA Version:.....	17
Conclusion:.....	18

## Table of Figures

Figure 1 image matrix x and kernel h.....	5
Figure 2 Inversion of h matrix .....	5
Figure 3 Convolution of pixels Row by Row .....	5
Figure 4 Resultant Matrix.....	6
Figure 5 Version 1 Source Code (Serial) .....	9
Figure 6 Version 2 Source Code (OMP).....	11
Figure 7 Version 3 Source Code (CUDA) .....	13
Figure 8 Serial Execution Time.....	14
Figure 9 Parallel Execution Time Using OMP .....	14
Figure 10 GPU Execution Time.....	15

# Implementation of the 2D-Convolution in CUDA C

## Problem Statement:

Convolution is the most important and fundamental concept in signal processing and analysis. Convolution is the process of adding each element of the image to its local neighbors, weighted by the kernel. Mathematically, it describes a rule of how to combine two functions or pieces of information to form a third function. The feature map (or input data) and the kernel are combined to form a transformed feature map. The convolution algorithm is often interpreted as a filter, where the kernel filters the feature map for certain information.

## Tools Used:

- Linux Ubuntu
- GCC/G++ Compiler
- OpenMP
- CUDA
- "Timer.h" (For Profiling)

## Algorithm:

25	100	75	49	130
50	80	0	70	100
5	10	20	30	0
60	50	12	24	32
37	53	55	21	90
140	17	0	23	222

$x$

1	0	1
0	1	0
0	0	1

$h$

Figure 1 image matrix  $x$  and kernel  $h$

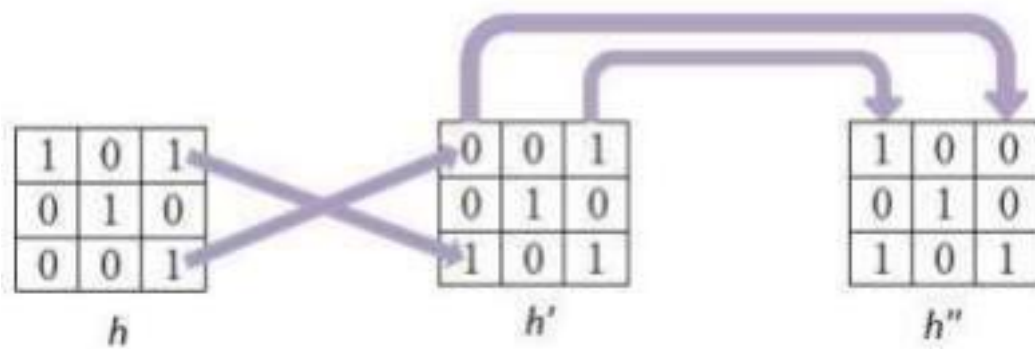


Figure 2 Inversion of  $h$  matrix

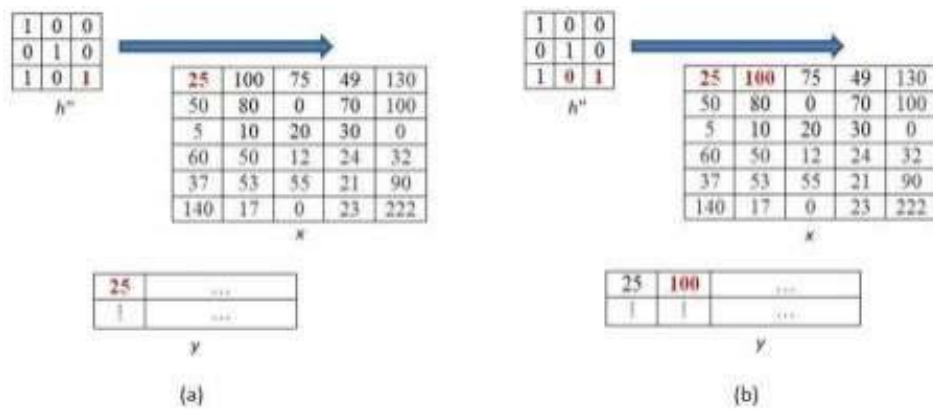
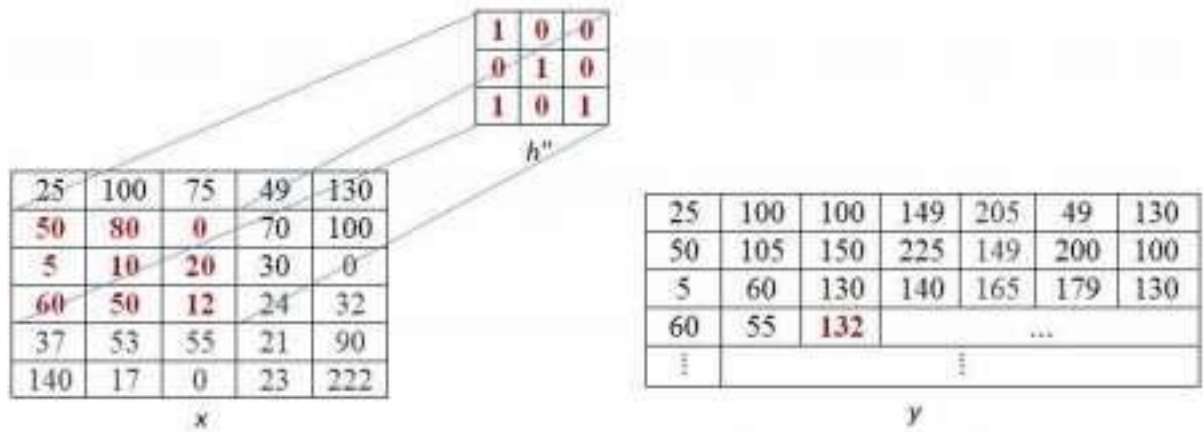
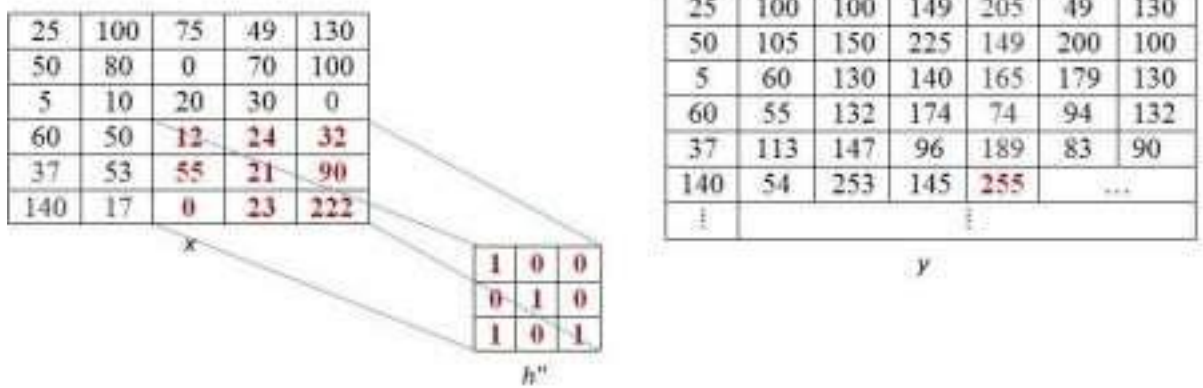


Figure 3 Convolution of pixels Row by Row



$$y(4,3) = 50 \times 1 + 80 \times 0 + 0 \times 0 + 5 \times 0 + 10 \times 1 + 20 \times 0 + 60 \times 1 + 50 \times 0 + 12 \times 1$$

$$= 50 + 0 + 0 + 0 + 10 + 0 + 60 + 0 + 12 = 132$$



$$y(6,5) = 12 \times 1 + 24 \times 0 + 32 \times 0 + 55 \times 0 + 21 \times 1 + 90 \times 0 + 0 \times 1 + 23 \times 0 + 222 \times 1$$

$$= 12 + 0 + 0 + 0 + 21 + 0 + 0 + 0 + 222 = 255$$

25	100	100	149	205	49	130
50	105	150	225	149	200	100
5	60	130	140	165	179	130
60	55	132	174	74	94	132
37	113	147	96	189	83	90
140	54	253	145	255	137	254
0	140	54	53	78	243	90
0	0	140	17	0	23	255

$y$

Figure 4 Resultant Matrix

## Parallelization Strategy:

### Version 1 (Serial):

We need to initialize the **convolution kernel** with the desire or some random values. Also, if we want to print the output matrix, we can do it through a **nested for loop**.

We will include all the necessary **libraries** for this operation. For larger images or more computationally intensive tasks, **parallel execution** using techniques like **multithreading** or **GPU acceleration** may provide significant performance improvements.

### Version 2 (Open MP):

To parallelize the convolution operation, we will **parallelize the loop**. I added the **#pragma omp parallel for** directive just before the **outer loop**. This directive instructs OpenMP to **distribute the iterations of the loop** across multiple threads for parallel execution. We just have to make sure that **OpenMP** is installed and the library is included in the source file.

### Version 3 (CUDA):

In this version, we have a **convolutionKernel** function that performs the 2D convolution operation on the **GPU**. Each thread in the **CUDA thread block** loads a portion of the input image and the convolution kernel into **shared memory**. It then performs the convolution operation and stores the result in the **output image**.

In the **main** function, we allocate memory on the **host** and the **device**, initialize the input image with **random values**, and **copy** it from the **host to device**. We define the **grid and block dimensions** based on the image size and the specified **block size**. Finally, we launch the **convolutionKernel** function on the GPU, copy the output image from the device to the host, and print the first few elements of the output for verification.

We need to make sure that we have an **Nvidia GPU** in the system and install the **CUDA Toolkit** before running this program.

## Source Code:

### Version 1 (Serial):

```
1  #include <iostream>
2  #include <cstdlib>
3  #include <ctime>
4  #include "Timer.h"
5
6  // Constants
7  const int MATRIX_SIZE = 512;
8  const int KERNEL_SIZE = 15;
9
10 // Function to perform convolution
11 void convolution(const unsigned char* inputMatrix, const float* kernel, unsigned char* outputMatrix) {
12     // Loop over the input matrix
13     for (int i = 0; i < MATRIX_SIZE; i++) {
14         for (int j = 0; j < MATRIX_SIZE; j++) {
15             float sum = 0.0;
16
17             // Apply the convolution kernel
18             for (int k = 0; k < KERNEL_SIZE; k++) {
19                 for (int l = 0; l < KERNEL_SIZE; l++) {
20                     int x = i - KERNEL_SIZE / 2 + k;
21                     int y = j - KERNEL_SIZE / 2 + l;
22
23                     // Handle boundary conditions by clamping
24                     if (x < 0) x = 0;
25                     if (x >= MATRIX_SIZE) x = MATRIX_SIZE - 1;
26                     if (y < 0) y = 0;
27                     if (y >= MATRIX_SIZE) y = MATRIX_SIZE - 1;
28
29                     sum += inputMatrix[x * MATRIX_SIZE + y] * kernel[k * KERNEL_SIZE + l];
30                 }
31             }
32
33             // Store the result in the output matrix
34             outputMatrix[i * MATRIX_SIZE + j] = sum;
35         }
36     }
37 }
38
39 int main()
40 {
41     //Declare Timer Variable
42     Timer t1("Serial Execution Time: ");
43
44     // Seed the random number generator
45     std::srand(std::time(nullptr));
46
47     // Create the input grayscale image matrix
48     unsigned char inputMatrix[MATRIX_SIZE * MATRIX_SIZE];
49     for (int i = 0; i < MATRIX_SIZE * MATRIX_SIZE; i++) {
50         inputMatrix[i] = std::rand() % 256; // Random pixel value between 0 and 255
51     }
```



```

51     }
52
53     // Create the convolution kernel
54     float kernel[KERNEL_SIZE * KERNEL_SIZE];
55     // Initialize the kernel with your desired values here
56
57     // Create the output grayscale image matrix
58     unsigned char outputMatrix[MATRIX_SIZE * MATRIX_SIZE];
59
60     //Start Timer
61     t1.Start();
62
63     // Perform convolution
64     convolution(inputMatrix, kernel, outputMatrix);
65
66     //Stop Timer
67     t1.Stop();
68
69     //Print Execution Time
70     t1.Print();
71
72     // Print the output matrix if desired
73     /*
74     for (int i = 0; i < MATRIX_SIZE; i++) {
75         for (int j = 0; j < MATRIX_SIZE; j++) {
76             std::cout << static_cast<int>(outputMatrix[i * MATRIX_SIZE + j]) << " ";
77         }
78         std::cout << std::endl;
79     }
80     */
81
82     return 0;
83 }
84

```

Figure 5 Version 1 Source Code (Serial)

## Version 2 (OpenMP):

```
1  #include <iostream>
2  #include <cstdlib>
3  #include <ctime>
4  #include <omp.h>
5  #include "Timer.h"
6
7  // Constants
8  const int MATRIX_SIZE = 512;
9  const int KERNEL_SIZE = 15;
10
11 // Function to perform convolution
12 void convolution(const unsigned char* inputMatrix, const float* kernel, unsigned char* outputMatrix) {
13     // Loop over the input matrix in parallel
14     #pragma omp parallel for
15     for (int i = 0; i < MATRIX_SIZE; i++) {
16         for (int j = 0; j < MATRIX_SIZE; j++) {
17             float sum = 0.0;
18
19             // Apply the convolution kernel
20             for (int k = 0; k < KERNEL_SIZE; k++) {
21                 for (int l = 0; l < KERNEL_SIZE; l++) {
22                     int x = i - KERNEL_SIZE / 2 + k;
23                     int y = j - KERNEL_SIZE / 2 + l;
24
25                     // Handle boundary conditions by clamping
26                     if (x < 0) x = 0;
27                     if (x >= MATRIX_SIZE) x = MATRIX_SIZE - 1;
28                     if (y < 0) y = 0;
29                     if (y >= MATRIX_SIZE) y = MATRIX_SIZE - 1;
30
31                     sum += inputMatrix[x * MATRIX_SIZE + y] * kernel[k * KERNEL_SIZE + l];
32                 }
33             }
34
35             // Store the result in the output matrix
36             outputMatrix[i * MATRIX_SIZE + j] = sum;
37         }
38     }
39 }
40
41 int main()
42 {
43     //Declare Timer Variable
44     Timer t1("Parallel Execution Time: ");
45
46     // Seed the random number generator
47     std::srand(std::time(nullptr));
48
49     // Create the input grayscale image matrix
50     unsigned char inputMatrix[MATRIX_SIZE * MATRIX_SIZE];
```

```

50     unsigned char inputMatrix[MATRIX_SIZE * MATRIX_SIZE];
51     for (int i = 0; i < MATRIX_SIZE * MATRIX_SIZE; i++) {
52         inputMatrix[i] = std::rand() % 256; // Random pixel value between 0 and 255
53     }
54
55     // Create the convolution kernel
56     float kernel[KERNEL_SIZE * KERNEL_SIZE];
57     // Initialize the kernel with your desired values here
58
59     // Create the output grayscale image matrix
60     unsigned char outputMatrix[MATRIX_SIZE * MATRIX_SIZE];
61
62     //Start Timer
63     t1.Start();
64
65     // Perform convolution
66     convolution(inputMatrix, kernel, outputMatrix);
67
68     //Stop Timer
69     t1.Stop();
70
71     //Print Execution Time
72     t1.Print();
73
74     // Print the output matrix if desired
75     /*
76     for (int i = 0; i < MATRIX_SIZE; i++) {
77         for (int j = 0; j < MATRIX_SIZE; j++) {
78             std::cout << static_cast<int>(outputMatrix[i * MATRIX_SIZE + j]) << " ";
79         }
80         std::cout << std::endl;
81     }
82     */
83
84     return 0;
85 }
86

```

Figure 6 Version 2 Source Code (OMP)

### Version 3 (CUDA):

```
1  #include <stdio.h>
2  #include "Timer.h"
3
4  #define WIDTH 512
5  #define HEIGHT 512
6  #define KERNEL_SIZE 15
7  #define BLOCK_SIZE 16
8
9  __global__ void convolutionKernel(const unsigned char* input, unsigned char* output) {
10     int tx = threadIdx.x;
11     int ty = threadIdx.y;
12     int bx = blockIdx.x;
13     int by = blockIdx.y;
14
15     // Compute the global position of the thread
16     int row = by * blockDim.y + ty;
17     int col = bx * blockDim.x + tx;
18
19     // Shared memory for the convolution kernel
20     __shared__ unsigned char sharedKernel[KERNEL_SIZE][KERNEL_SIZE];
21
22     // Load the convolution kernel into shared memory
23     if (ty < KERNEL_SIZE && tx < KERNEL_SIZE) {
24         sharedKernel[ty][tx] = input[row * WIDTH + col];
25     }
26     __syncthreads();
27
28     // Convolution operation
29     int sum = 0;
30     if (row < HEIGHT && col < WIDTH) {
31         for (int i = 0; i < KERNEL_SIZE; i++) {
32             for (int j = 0; j < KERNEL_SIZE; j++) {
33                 sum += sharedKernel[i][j] * input[(row + i) * WIDTH + (col + j)];
34             }
35         }
36         output[row * WIDTH + col] = sum / (KERNEL_SIZE * KERNEL_SIZE);
37     }
38 }
39
40 int main()
41 {
42     //Declare Timer Variables
43     Timer gputime;
44     initTimer(&gputime, "GPU Execution Time: ");
45
46     unsigned char *h_input, *h_output;
47     unsigned char *d_input, *d_output;
48
49     size_t size = WIDTH * HEIGHT * sizeof(unsigned char);
50
51     // Allocate host memory
```

```

52     h_input = (unsigned char*)malloc(size);
53     h_output = (unsigned char*)malloc(size);
54
55     // Allocate device memory
56     cudaMalloc((void**)&d_input, size);
57     cudaMalloc((void**)&d_output, size);
58
59     // Initialize input data with random values
60     for (int i = 0; i < WIDTH * HEIGHT; i++) {
61         h_input[i] = rand() % 256;
62     }
63
64     // Copy input data from host to device
65     cudaMemcpy(d_input, h_input, size, cudaMemcpyHostToDevice);
66
67     // Define grid and block dimensions
68     dim3 blockDim(BLOCK_SIZE, BLOCK_SIZE);
69     dim3 gridDim((WIDTH + BLOCK_SIZE - 1) / BLOCK_SIZE, (HEIGHT + BLOCK_SIZE - 1) / BLOCK_SIZE);
70
71     //Start Timer
72     startTimer(&gputime);
73
74     // Launch the convolution kernel
75     convolutionKernel<<<gridDim, blockDim>>>(d_input, d_output);
76
77     //Stop Timer
78     stopTimer(&gputime);
79
80     //Print Execution Time
81     printTimer(gputime);
82
83     // Copy output data from device to host
84     cudaMemcpy(h_output, d_output, size, cudaMemcpyDeviceToHost);
85
86     // Print the first few elements of the output for verification
87     //for (int i = 0; i < 10; i++) {
88     //    printf("%u ", h_output[i]);
89     //}
90     printf("\n");
91
92     // Free device memory
93     cudaFree(d_input);
94     cudaFree(d_output);
95
96     // Free host memory
97     free(h_input);
98     free(h_output);
99
100     return 0;
101 }

```

Figure 7 Version 3 Source Code (CUDA)

## Detailed Profile & Analysis:

### Version 1:

We have used the **Timer.h** header file for profiling. We simply initialize a timer variable which will store the execution time. Start the timer before calling the convolution function and end the timer after the execution. Finally, we just print the execution time.

```
abdullah@abdullah-PE70-6QE:~/pdc/project/serial$ g++ code.cpp -o output
abdullah@abdullah-PE70-6QE:~/pdc/project/serial$ ./output
Serial Execution Time: : 304.931 msec
abdullah@abdullah-PE70-6QE:~/pdc/project/serial$
```

*Figure 8 Serial Execution Time*

### Version 2:

We have followed the same procedure for profiling as in version 1. Include **Timer.h** library, start and end the timer before and after the function call for the convolution.

I have tested parallel execution with 2, 4, 6 and 8 threads to see which one gives the best performance.

```
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$ export OMP_NUM_THREADS=2
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$ g++ -o output -fopenmp code.cpp
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$ ./output
Parallel Execution Time: : 156.366 msec
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$ export OMP_NUM_THREADS=4
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$ g++ -o output -fopenmp code.cpp
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$ ./output
Parallel Execution Time: : 82.697 msec
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$ export OMP_NUM_THREADS=6
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$ g++ -o output -fopenmp code.cpp
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$ ./output
Parallel Execution Time: : 98.531 msec
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$ export OMP_NUM_THREADS=8
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$ g++ -o output -fopenmp code.cpp
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$ ./output
Parallel Execution Time: : 78.384 msec
abdullah@abdullah-PE70-6QE:~/pdc/project/omp$
```

*Figure 9 Parallel Execution Time Using OMP*

### Version 3:

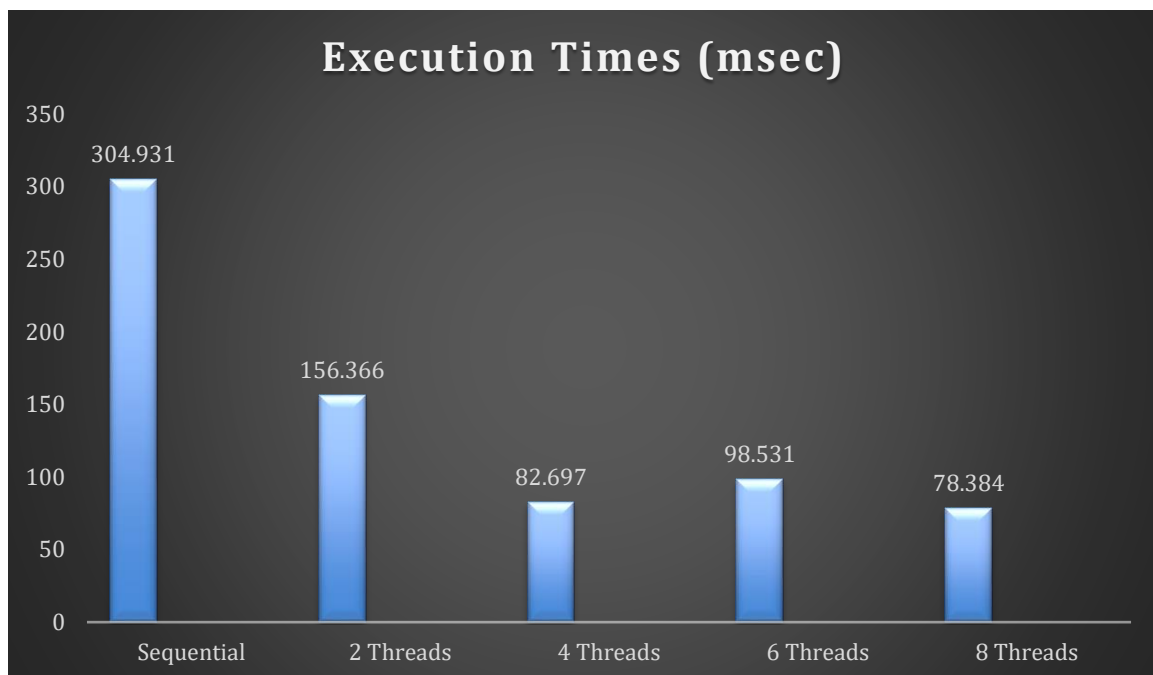
Simple include the **Timer.h** library and start and end the timer before and after the convolutionKernel call.

```
abdullah@abdullah-PE70-6QE:~/pdc/project/cuda$ nvcc -o output code.cu
code.cu(44): warning #2464-D: conversion from a string literal to "char *" is deprecated
code.cu(44): warning #2464-D: conversion from a string literal to "char *" is deprecated
abdullah@abdullah-PE70-6QE:~/pdc/project/cuda$ ./output
GPU Execution Time: = 0.001 msec
abdullah@abdullah-PE70-6QE:~/pdc/project/cuda$
```

Figure 10 GPU Execution Time

### Comparison Between Version 1 & Version 2:

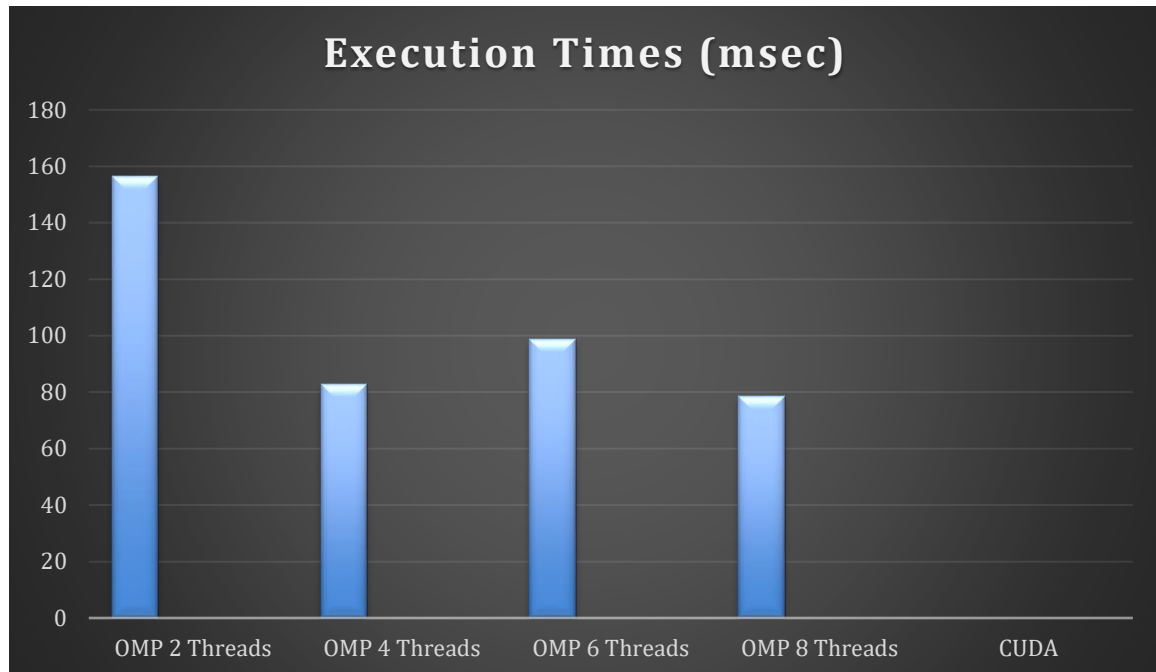
**Version 1** was the **Sequential** version of the convolution program, and **Version 2** was the **Parallel** version of the same program. Both have been executed on **CPU**. The Following figure will give us a comparison between the performances of the two versions:



The best results have been achieved using **8 Threads**

### Comparison Between Version 2 & Version 3:

Version 2 was the parallel version for convolution using OpenMP and Version 3 was the GPU or CUDA version for the same program. The following figure will give us a comparison between the performances of the two versions:

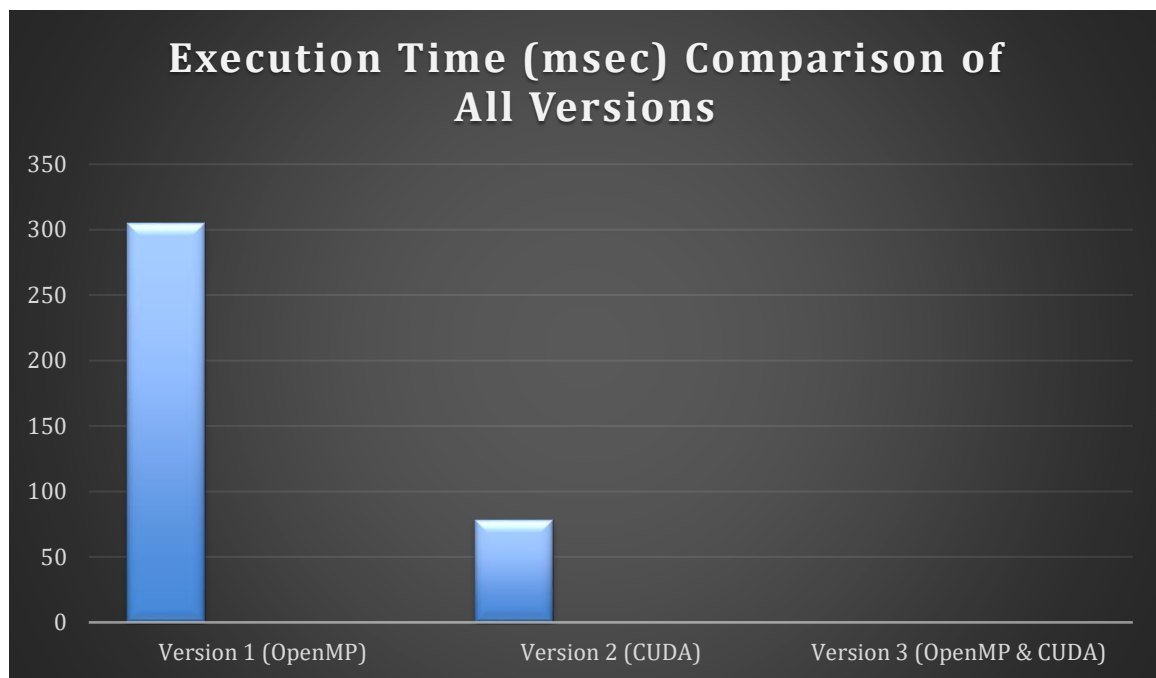


CUDA gives such a huge performance improvement



## Execution Time & Speed up Results:

Versions	Execution Times
Version 1	304.931 msec
Version 2	78.384 msec (best)
Version 3	0.001 msec



### Speedup From Version 1 to Version 2:

$$\text{Speed up} = \text{Execution Time (old)} / \text{Execution Time (new)}$$

$$\text{Speed up} = 304.931 / 78.384$$

$$\text{Speed up} = 3.8902$$

### Speedup From Version 2 to Version 3:

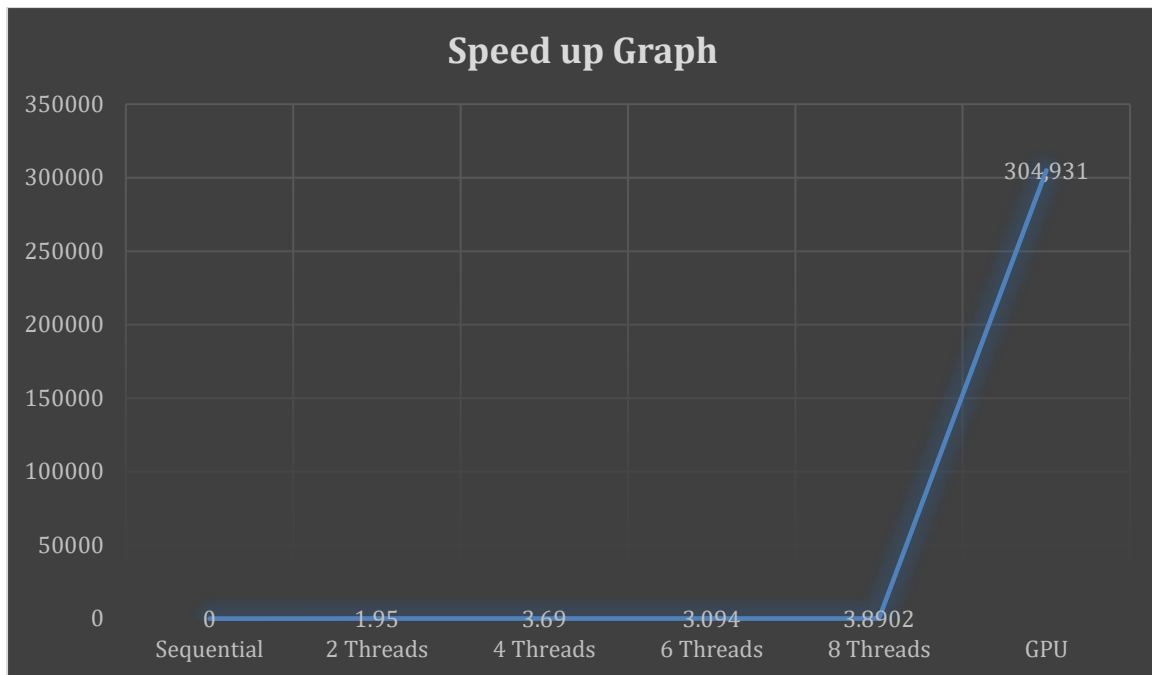
$$\text{Speed up} = 78.384 / 0.001$$

$$\text{Speed up} = 78,384$$

### Speedup from Serial to CUDA Version:

$$\text{Speed up} = 304.931 / 0.001$$

$$\text{Speed up} = 304,931$$



## Conclusion:

We have implemented many possible execution strategies like **sequential, parallel with 2, 4, 6 and 8 Threads** on **CPU** and execution on **GPU**. After testing all of the methods, **repetitive profiling** and testing, calculating **Speedup**, we have come to the conclusion that **GPU** adds exponential increase in a machine's performance. The graphs show that **multithreading on CPU** is good, but it is no way near as fast as **GPU**, because the GPU has a much greater number of **cores & threads** as compared to a **CPU**.