

# Implementation of Algorithms using CUDA, and their Performance Comparison

<i>Anuja B. S.</i>	<i>1NH10CS019</i>
<i>Chaitra J. R.</i>	<i>1NH10CS022</i>
<i>Kusuma Chandrika K.</i>	<i>1NH10CS047</i>
<i>Nitin Kaveriappa U. M.</i>	<i>1NH10CS061</i>



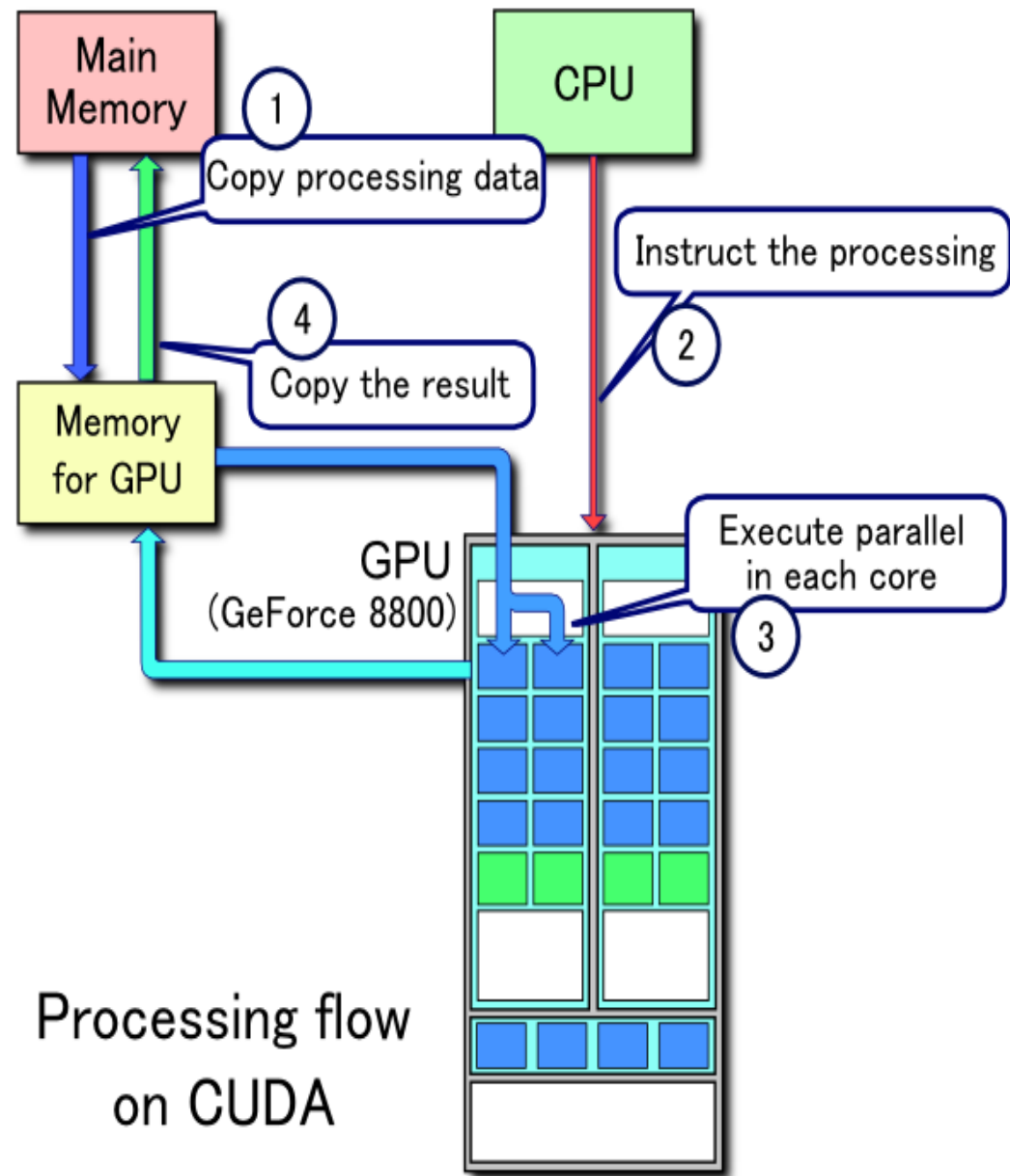
# Objective

- ▶ This project is focused on exploiting GPU (Graphics Processing Unit) for general purpose computing based on CUDA (Compute Unified Device Architecture).
- ▶ It aims at utilization of GPU cores through direct access of GPU computing using CUDA that enables straightforward implementation of parallel algorithms.
- ▶ Diamond Square, Quickhull and Smith Waterman algorithms were implemented in CUDA C and C language and the algorithms were compared on performance and efficiency.

# Introduction

- ▶ **Compute Unified Device Architecture (CUDA)** is a parallel computing platform
- ▶ Programming model created by NVIDIA
- ▶ Implemented by the graphics processing units (GPUs) that they produce.

# Processing Flow on CUDA



# General Keywords and Functions

- ▶ Header file: The header file used for CUDA is `#include<cuda.h>`
- ▶ `__global__`:  
The `__global__` qualifier declares a function as being a kernel.  
Executed on the device, callable from the host,  
`__global__` functions must have void return type.  
Ex: `__global__ void subhtKernel( );`
- ▶ `cudaMalloc()`:  
A single address space is used for the host and all the devices.  
Syntax: `cudaMalloc(void ** pointer, size_t nbytes)`  
Ex: `cudaMalloc((void**) &d_j, sizeof(int));`

## ► cudaMemcpy():

Memory is typically allocated using **cudaMalloc()** and freed using **cudaFree()**.

Data transfer between host memory and device memory are done through this function.

Syntax: `cudaMemcpy(void *dst, void *src, size_t nbytes, enum cudaMemcpyKind direction);`

Ex: `cudaMemcpy(d_j, &j, sizeof(int), cudaMemcpyHostToDevice);`

`cudaMemcpy(d_j, &j, sizeof(int), cudaMemcpyDeviceToHost);`

Syntax: `cudaFree(void* pointer)`

Ex: `cudaFree(d_j);`

► Function call(kernel):

CUDA threads that execute that kernel for a given kernel call is specified using a new

`<<<...>>>` *execution configuration*.

Syntax: `Function_name<<< no.of blocks, no.of threads/blocks>>>`

Ex: `subhtKernel<<<1, j>>>( );`

Each thread that executes the kernel is given a unique *thread ID* that is accessible within the kernel through the built-in **threadIdx** variable.

► cudaDeviceSynchronize():

Blocks until the device has completed all preceding requested tasks.

CudaDeviceSynchronize() returns an error if one of the preceding tasks has failed.

## ► `__shared__`:

The `__shared__` qualifier, declares a variable that:

- Resides in the shared memory space of a thread block,
- Has the lifetime of the block,
- Is only accessible from all the threads within the block.

When declaring a variable in shared memory as an external array such as

Ex: `extern __shared__ float shared[];`



▶ threadIdx:

threadIdx is a 3-component vector

threads can be identified using a one-dimensional, two-dimensional, or three-dimensional thread index.

Ex: `int tid = threadIdx.x;(1-d)`

▶ \_\_syncthreads():

If the threads belong to the same block, in which case they should use \_\_syncthreads() and share data through shared memory.

# CUDA Function

```
__global__ void subhtKernel(int *limit)
{
    extern __shared__ kernel_shared_type sdata[];
    unsigned int tid = threadIdx.x;
    if (tid < *limit)
    {
        __syncthreads();
    }
}

points subht(int j)
{
    int *d_j;
    cudaMalloc((void**) &d_j, sizeof(int));
    cudaMemcpy(d_j, &j, sizeof(int), cudaMemcpyHostToDevice);
    subhtKernel<<<1, j>>>(d_j,);
    cudaDeviceSynchronize();
    cudaMemcpy(&j, d_j, sizeof(int), cudaMemcpyDeviceToHost);
    cudaFree(d_j);
}
```

# Non - CUDA Function

```
points subht(int j)
{
    double h[200];
    for(i=0;i<j;i++)
    {
        //find the point with max distance from the baseline
    }
}
```

# CUDA AND IT'S DOMAINS

- ▶ *Graphical Rendering*

Volume rendering using 2-3 swap image compositing

GPU accelerated volume rendering on an interactive light field display

Diamond Square Algorithm

- ▶ *Geometric Computation*

Geometric problems on 3D meshes

Parallel Convex Hull Algorithm

Quickhull Algorithm

- ▶ *Genomics*

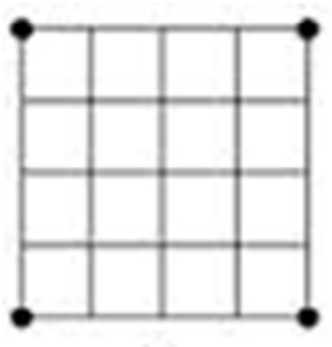
Smith Waterman Algorithm

# Introduction to Diamond Square Algorithm

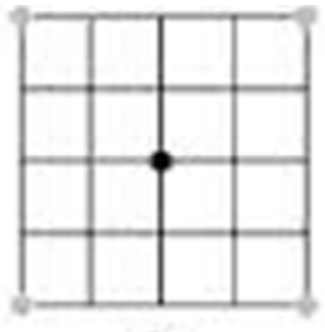
- ▶ The **diamond-square algorithm** is a method for generating height-maps for computer graphics.

# Diamond Square Algorithm

- ▶ Step 1: Create a 2D array. The size of the array is determined by the amount of iterations (i.e. 1 iteration-3x3 array, 2 iterations-5x5 array etc).
- ▶ Step 2: Assign a start value to the four corners of the map.



- ▶ Step 3: Call the Diamond() method.

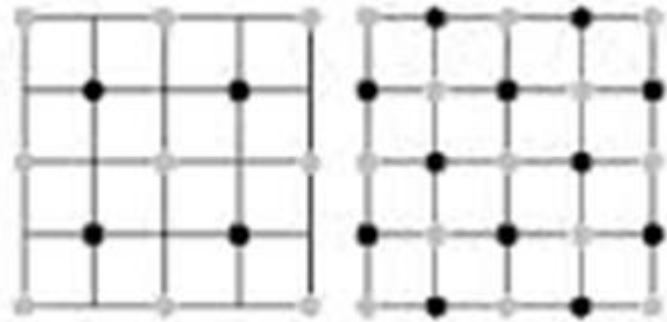


# Diamond Square Algorithm

- ▶ Step 4: Call the SquareEven() method.
- ▶ Step 5: Call the SquareOdd() method.

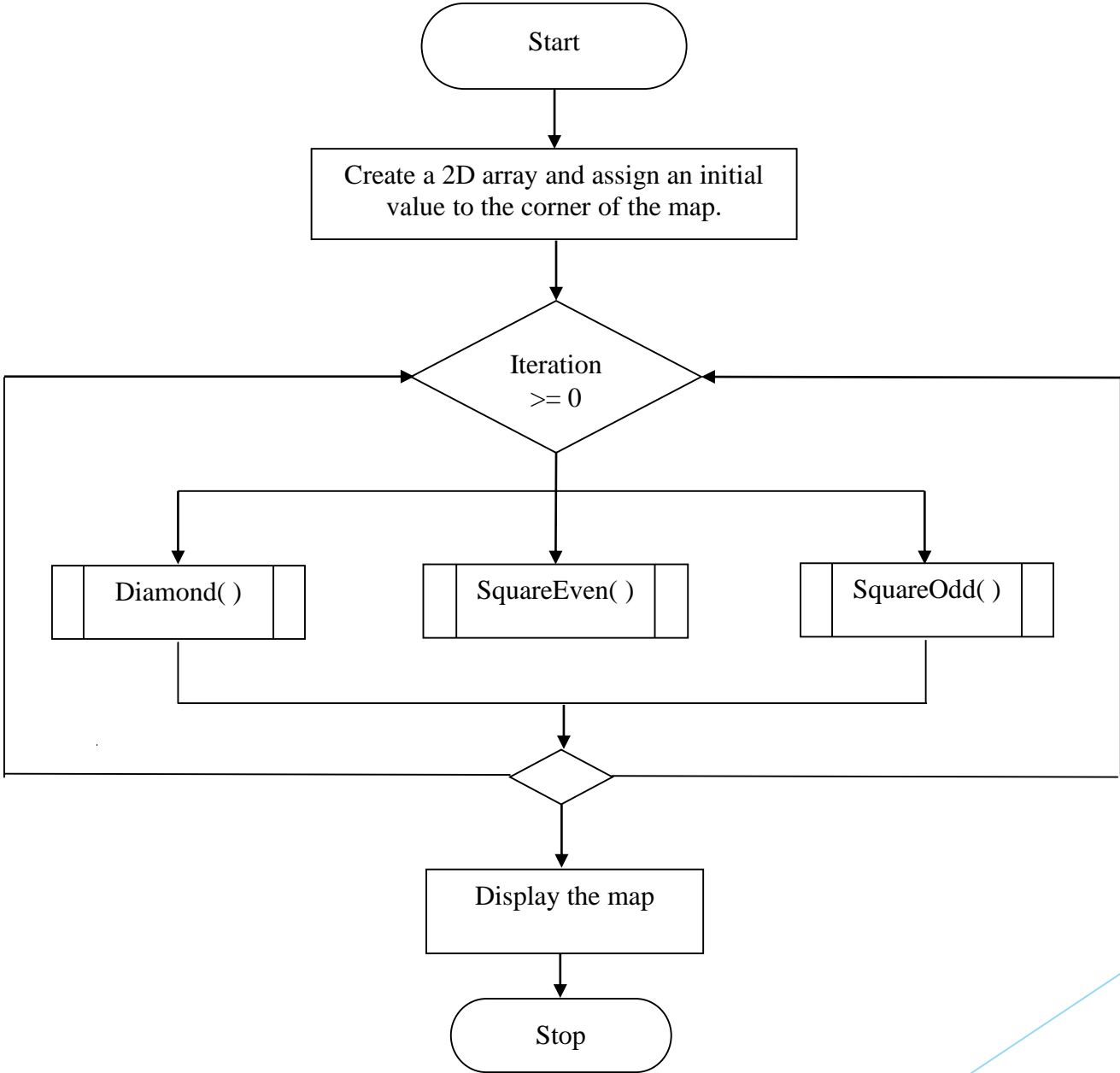


- ▶ Step 6: Go back to step 3 and repeat the steps until iterations  $\geq 0$ .

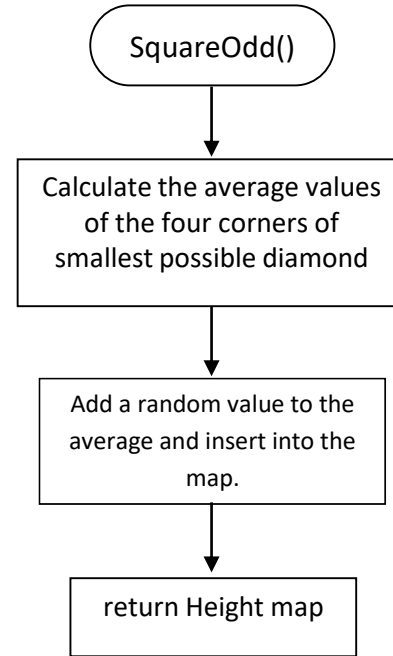
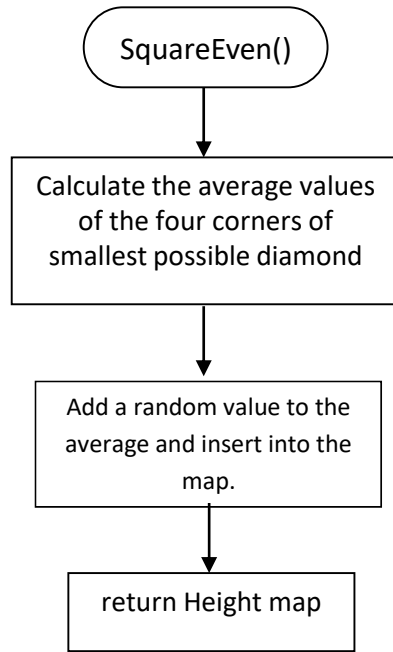
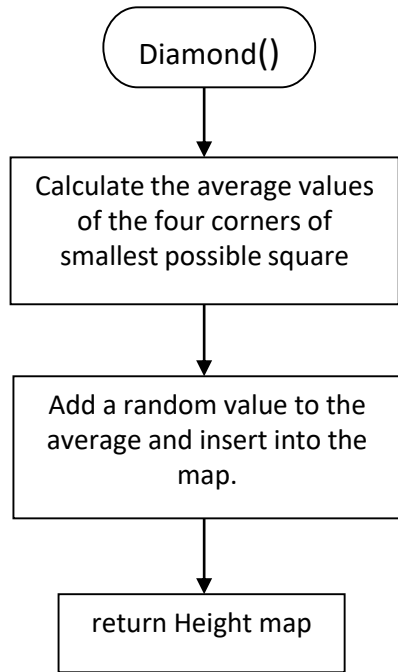


- ▶ Step 7: Display the map.

# Flowchart for Diamond Square Algorithm

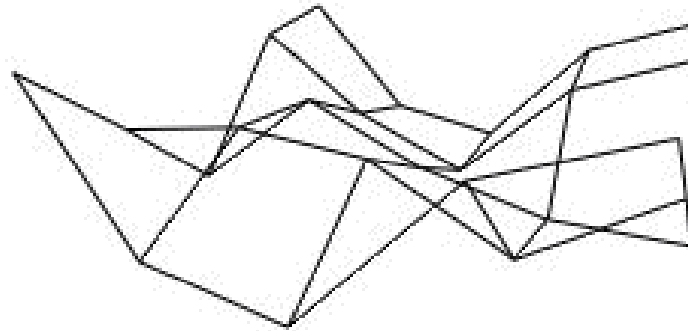
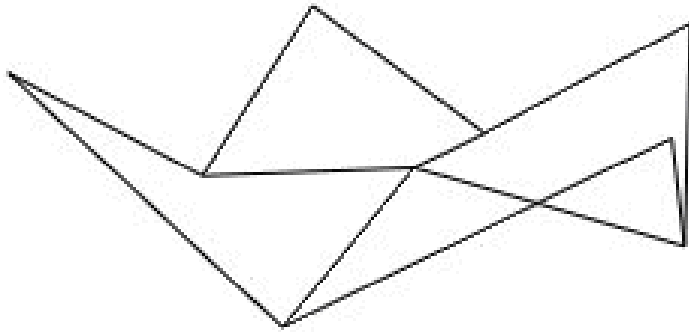






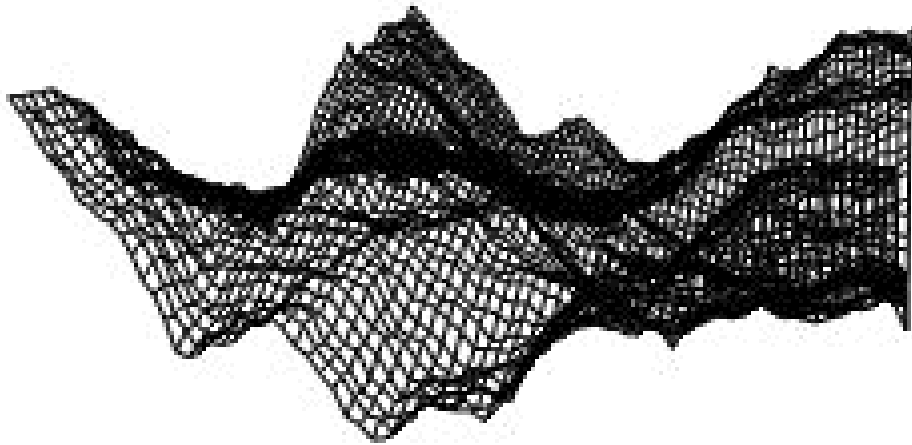
# Implementation

- ▶ After the first pass if we were to connect the nine points , we might get a wireframe surface which looks like this:



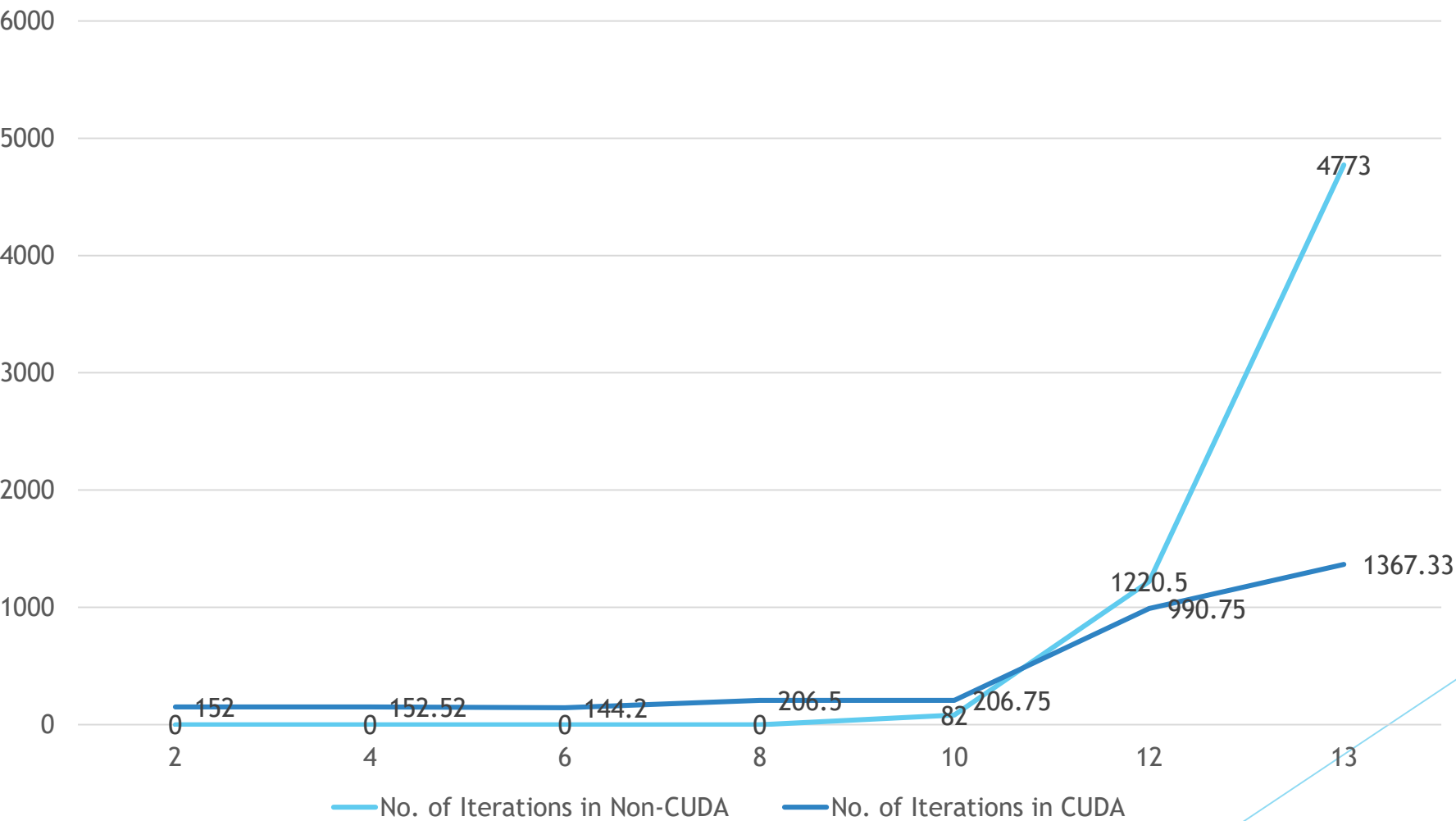
- ▶ When we perform the second pass with 12 diamond centres ,we need to calculate 12 new values . 25 elements of the array have been arranged and the wireframe would look like above:

- ▶ Had a larger array been allocated we could add more detail in each pass and the surface would like this



# Non-CUDA Performance Graph

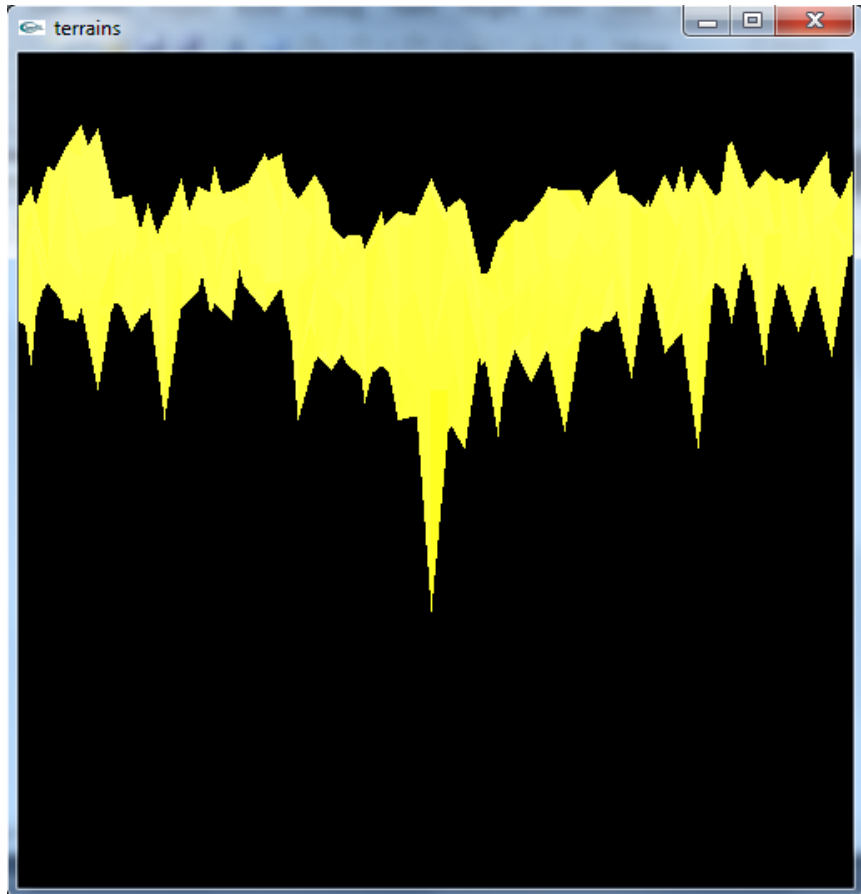
GRAPH FOR DIAMOND SQUARE NON-CUDA PROGRAM



# Performance Comparison

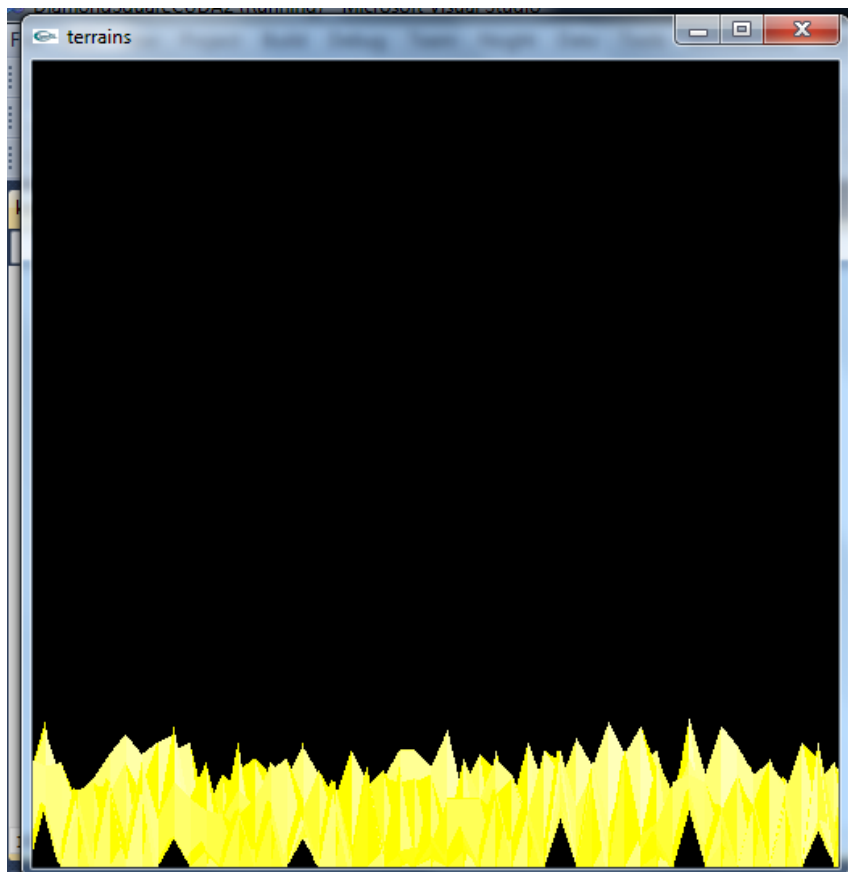
No. of iterations	Non CUDA execution time(ms)	CUDA execution time(ms)
2	0.00	152.00
4	0.00	152.25
6	0.00	144.25
8	0.00	156.00
10	82.00	206.75
12	1220.50	990.75
13	4779.00	1367.33

# Non-CUDA Output



```
C:\Users\Nitin Kaveriappa\Documents\Visual Studio 2010\Projects\DiamondSquareC\Debug\Diam...  
size is 8193  
maxIndex:8192size=8192  
size=4096  
size=2048  
size=1024  
size=512  
size=256  
size=128  
size=64  
size=32  
size=16  
size=8  
size=4  
size=2  
end  
Time elapsed in ms: 4820.00
```

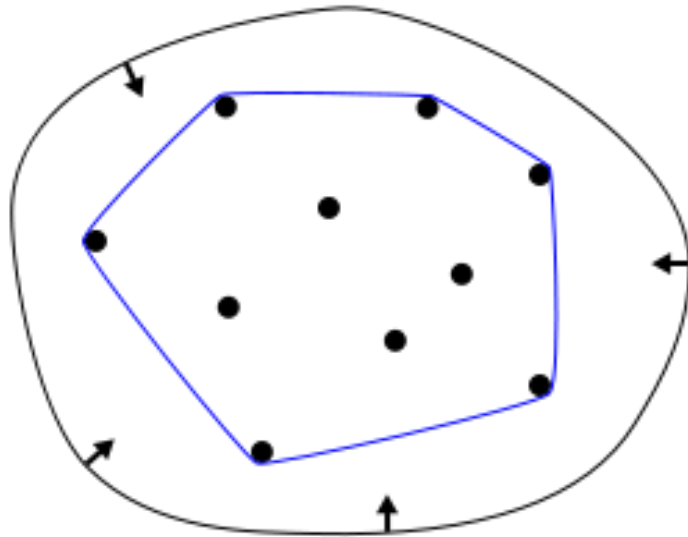
# CUDA Output



```
C:\Users\Nitin Kaveriappa\Documents\Visual Studio 2010\Projects\DiamondSquareCUDA2\Debug...  
size is 8193  
size=8192  
size=4096  
size=2048  
size=1024  
size=512  
size=256  
size=128  
size=64  
size=32  
size=16  
size=8  
size=4  
size=2  
Time elapsed in ms: 3323.00
```

# Introduction to QuickHull Algorithm

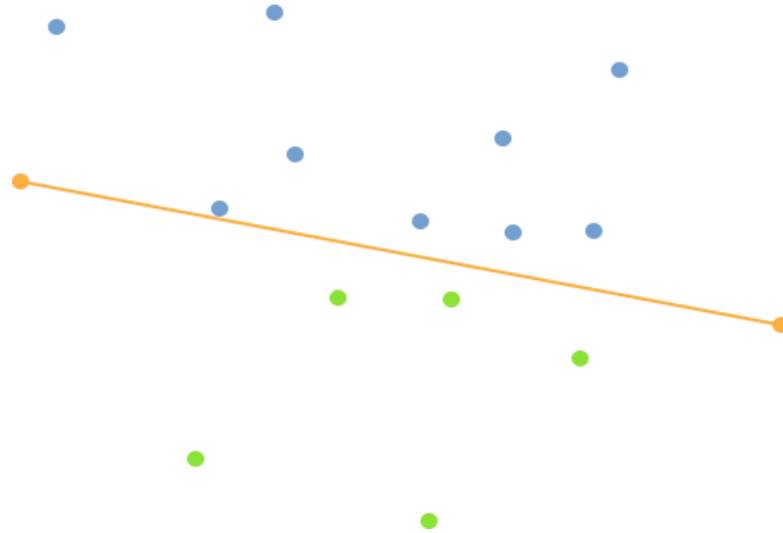
- ▶ **QuickHull** is a method of computing the convex hull of a finite set of points in the plane.
- ▶ It uses a divide and conquer approach.



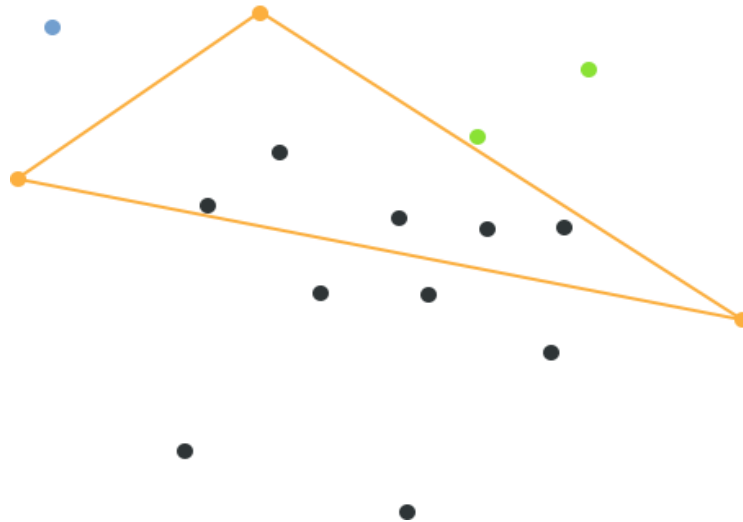


# QuickHull Algorithm

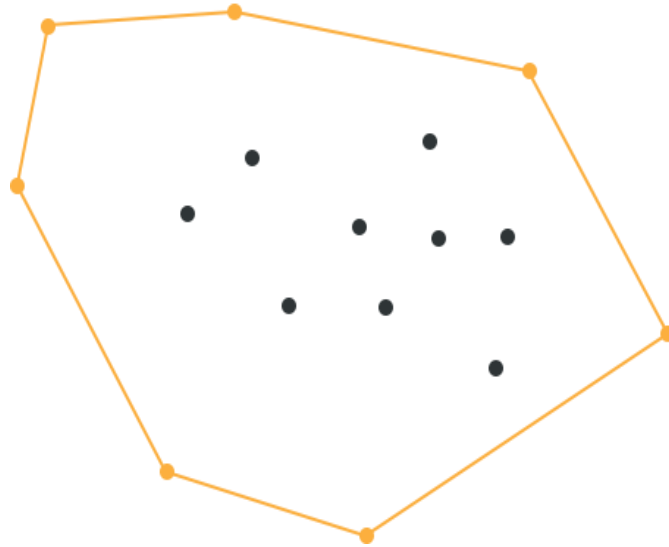
- ▶ Step 1: Points consists of x and y coordinates.
- ▶ Step 2: Find two points furthest away from each other and this is baseline PQ.
- ▶ Step 3: This PQ line divides points into 2 subsets  $S_1$  and  $S_2$ .



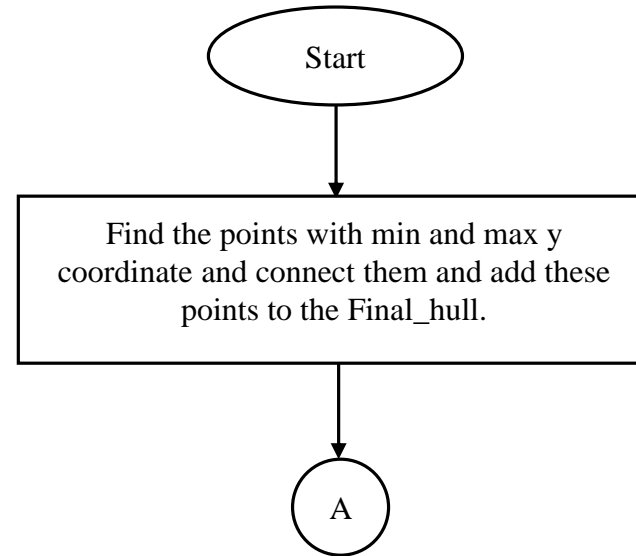
- ▶ Step 4: From line segment PQ find the point with the maximum distance from the line.  
Let it be 'A'.
- ▶ A triangle is formed and the points inside the triangle is ignored.

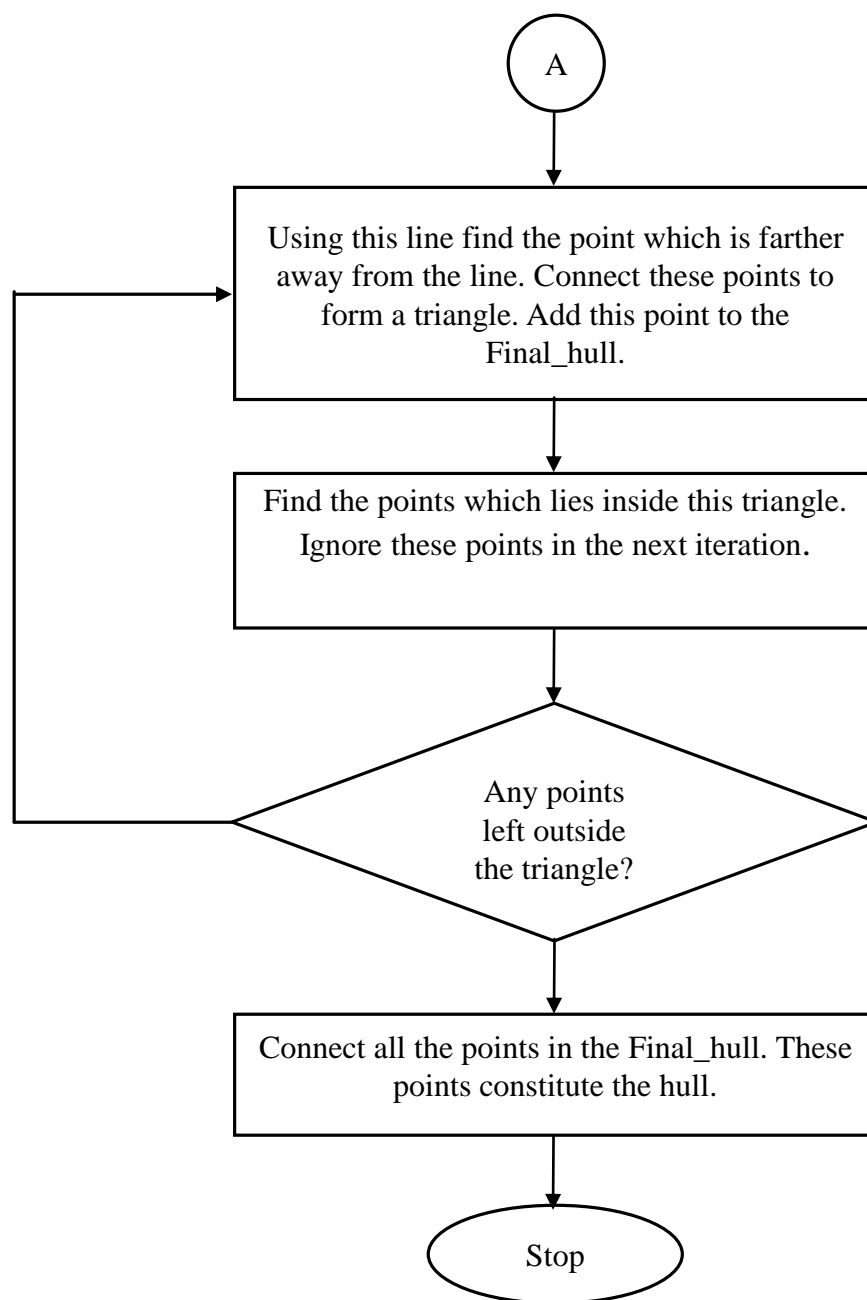


- ▶ Step 5: Repeat these steps until no points are left and connect these points.

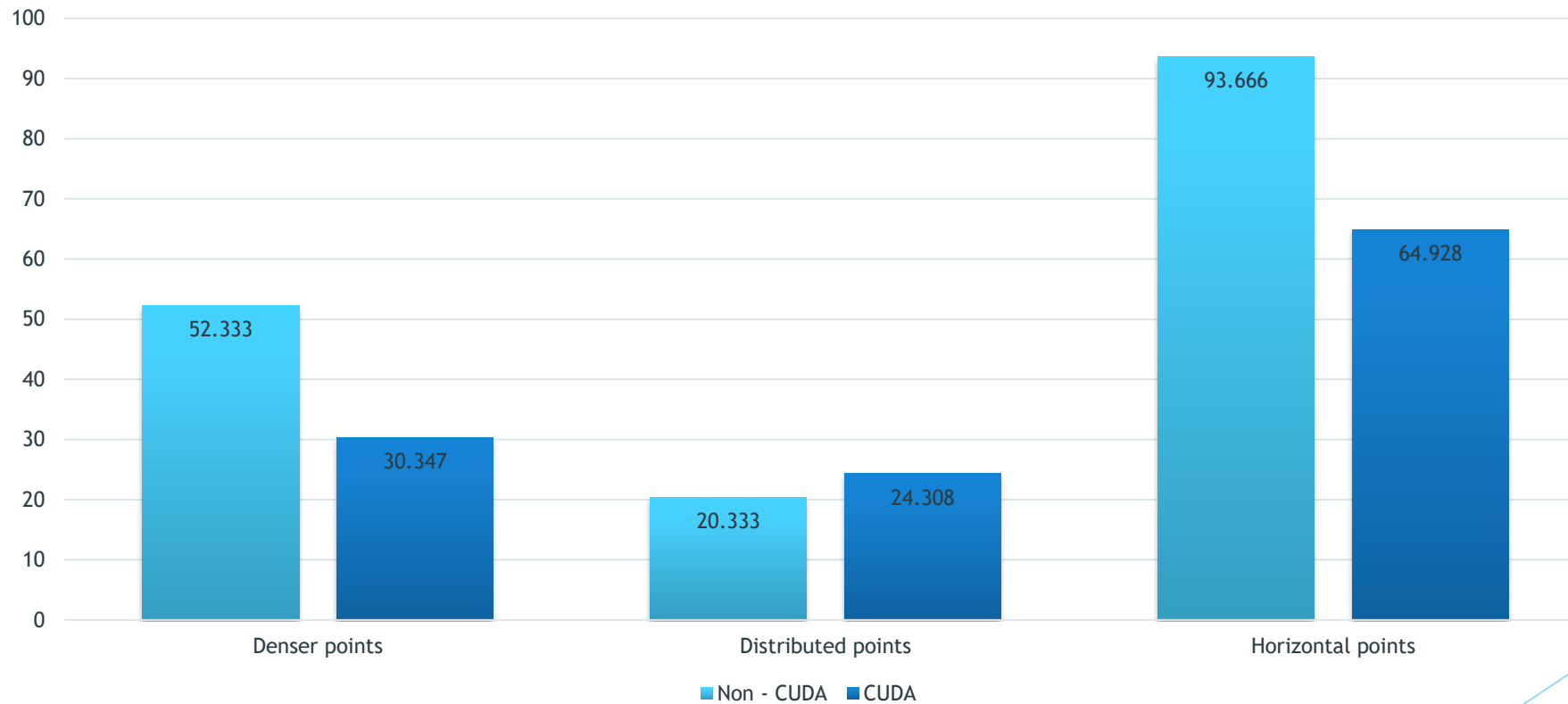


# Flowchart for QuickHull Algorithm

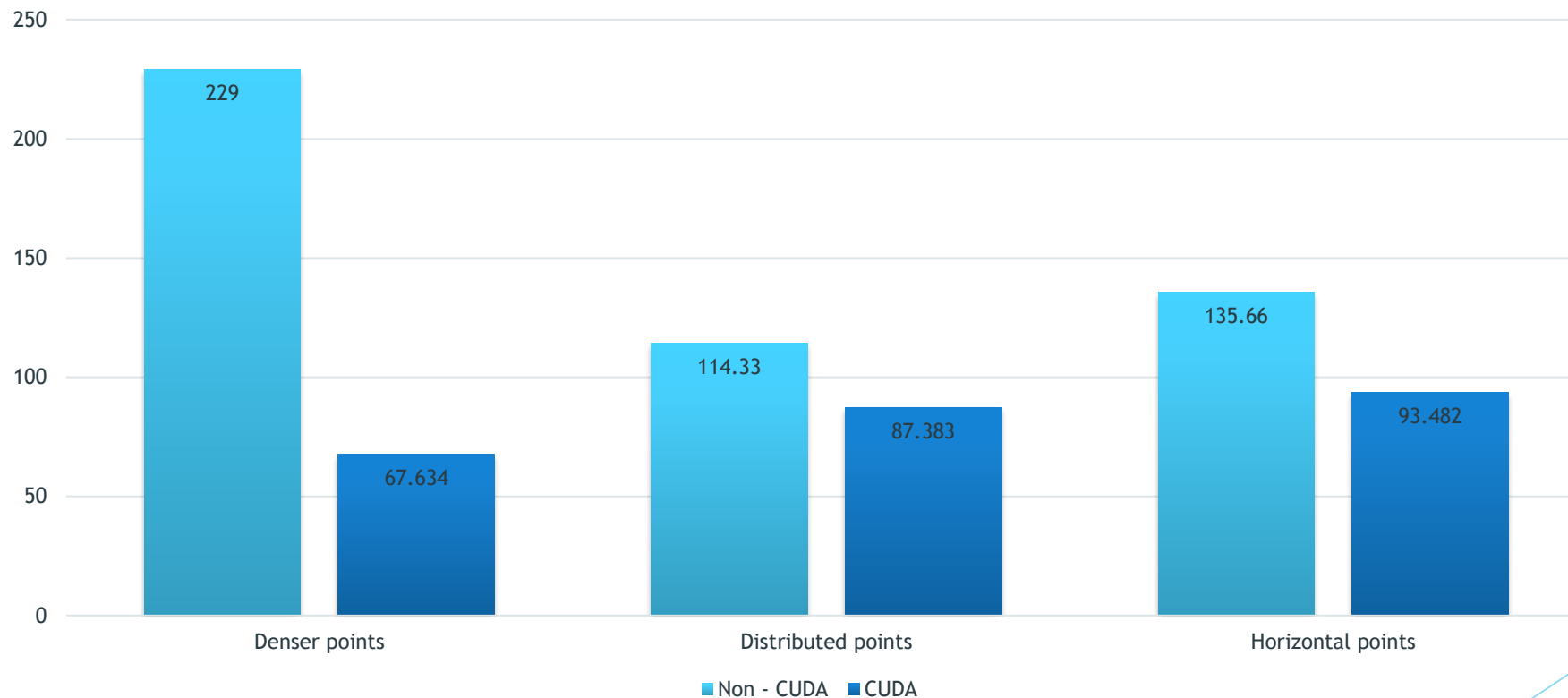




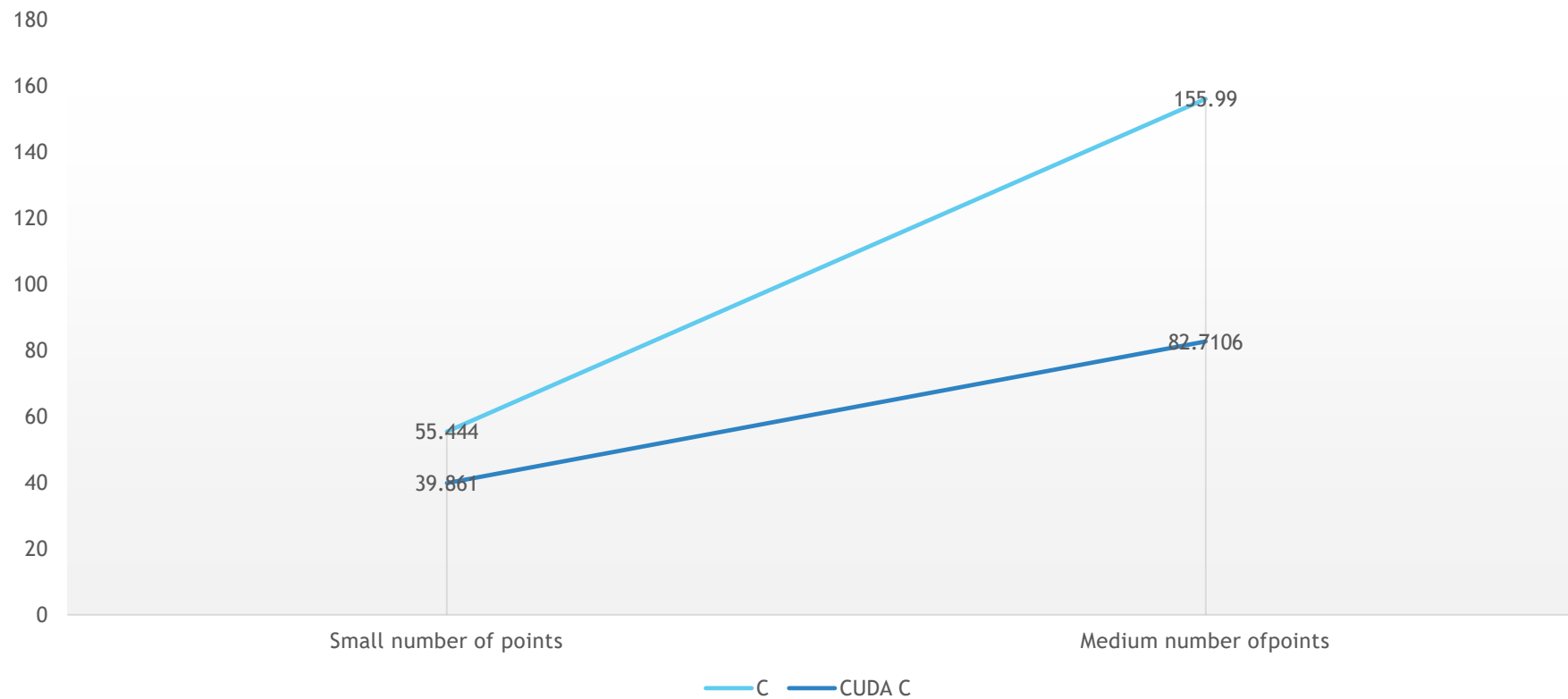
# Graph of execution time for Quickhull C and CUDA C program for smaller number of points



# Graph of execution time for Quickhull C and CUDA C program for medium number of points

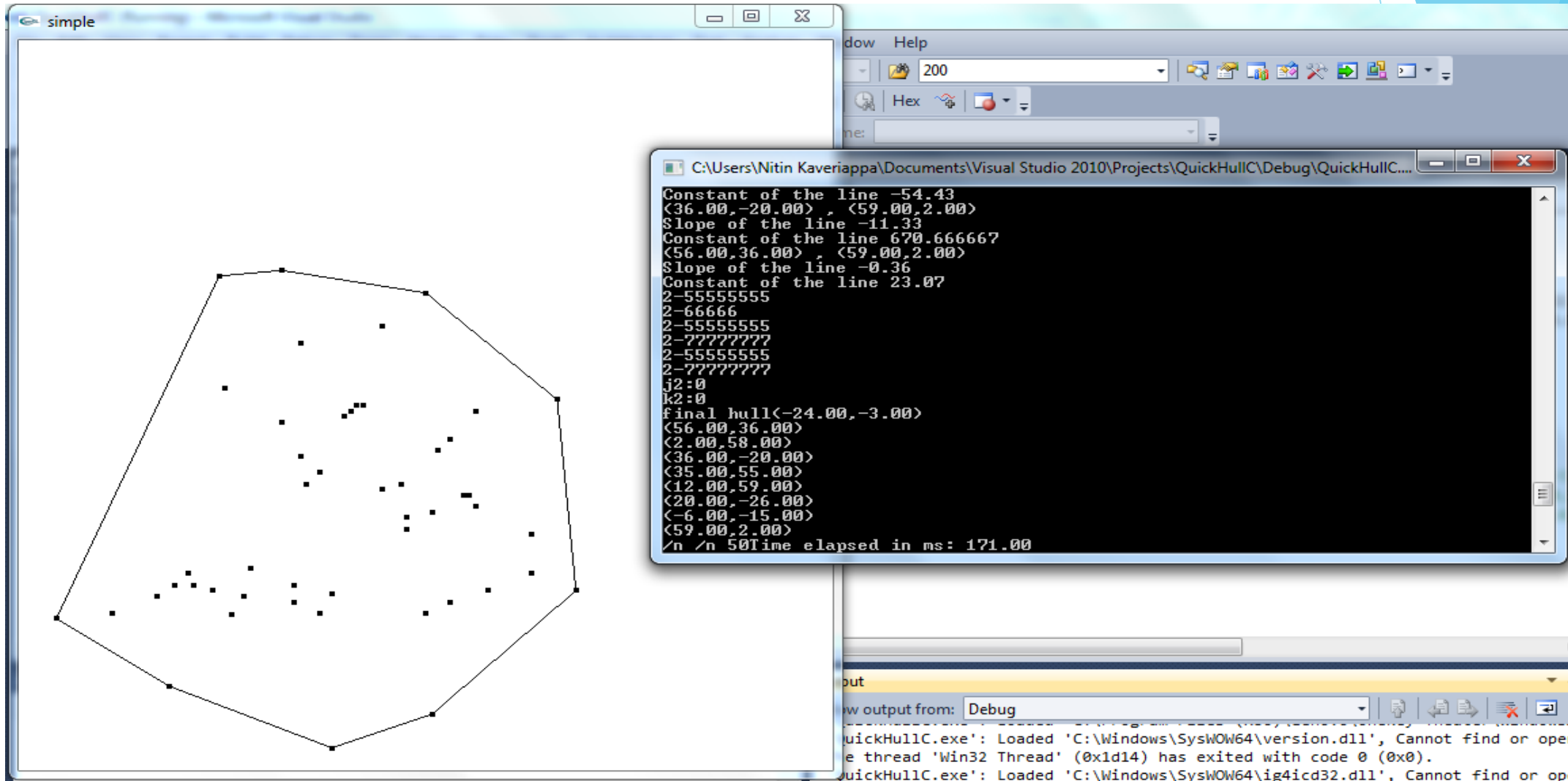


# Graph of execution time for Quickhull C and Quickhull CUDA C program

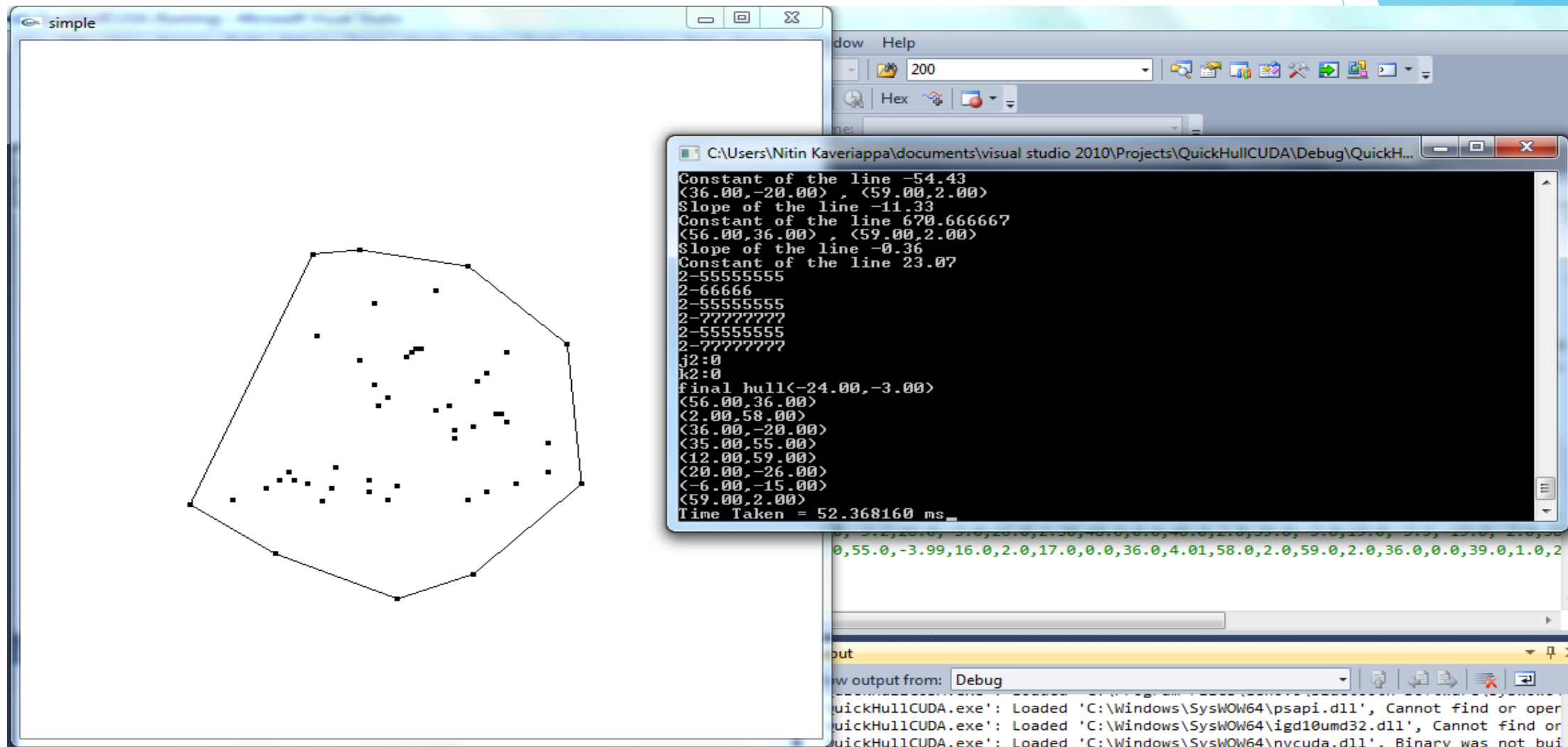




# Non-CUDA Output



# CUDA Output



# Introduction to Smith-Waterman Algorithm

- Optimal local alignment of two sequences
- Performs an exhaustive search for the optimal local Alignment
- Based on the 'dynamic programming' algorithm
  - ✓ Fill the matrix
  - ✓ Find the maximal value (score) in the matrix
  - ✓ Trace back from the score until a 0 value is reached

# Smithwaterman Algorithm

- ▶ Initialisation:  $a, b \rightarrow$  Strings over the alphabet  $\Sigma$   
 $m \rightarrow$  length of  $(a)$   
 $n \rightarrow$  length of  $(b)$   
 $H(i, j) \rightarrow$  is the max similarity score between a suffix of  $a[1 \dots i]$  and a suffix of  $b[1 \dots j]$   
 $W(c, d)$ ,  $c, d \in \Sigma \cup \{-\}$ ,  $'-\'$  is the gap scoring scheme
- ▶ Steps: Matrix  $H$  is built as follows:  
 $H(i, 0) = 0$ ,  $0 \leq i \leq m$   
 $H(0, j) = 0$ ,  $0 \leq j \leq n$   
If  $a_i = b_j$ , then  $W(a_i, b_j) = W(\text{match})$  or  
If  $a_i \neq b_j$ , then  $W(a_i, b_j) = W(\text{mismatch})$

$H(i,j) = \max \{ H(i-1,j-1) + w(a_i,b_j) \text{ Match/mismatch}$

$H(i-1,j) + W(a_i,-) \text{ Deletion}$

$H(i,j-1) + W(-,b_j) \text{ Insertion} \}$  for

$1 \leq i \leq m$  and  $1 \leq j \leq n$

Where  $W(\text{match}) = +2$

$W(a,-) = W(-,b) = W(\text{mismatch}) = -1$

► Output: Scoring matrix H

	-	A	C	A	C	C	A	C	A
-	0	0	0	0	0	0	0	0	0
A	0	2	1	2					
G	0	1	1	1					
C	0	0	3	2					
A	0								
C	0								
A	0								
C	0								
A	0								

- ▶ Input: Scoring matrix  $H$
- ▶ Output: Alignment matrix  $T$

$$H = \begin{pmatrix} - & A & C & A & C & A & C & T & A \\ - & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ A & 0 & 2 & 1 & 2 & 1 & 2 & 1 & 0 & 2 \\ G & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 1 \\ C & 0 & 0 & 3 & 2 & 3 & 2 & 3 & 2 & 1 \\ A & 0 & 2 & 2 & 5 & 4 & 5 & 4 & 3 & 4 \\ C & 0 & 1 & 4 & 4 & 7 & 6 & 7 & 6 & 5 \\ A & 0 & 2 & 3 & 6 & 6 & 9 & 8 & 7 & 8 \\ C & 0 & 1 & 4 & 5 & 8 & 8 & 11 & 10 & 9 \\ A & 0 & 2 & 3 & 6 & 7 & 10 & 10 & 10 & 12 \end{pmatrix}$$

- ▶ Input: Alignment matrix T.
- ▶ Output: Optimal local alignment of sequences a and b i.e. for determining similar regions between two strings.

$$H = \begin{pmatrix} - & - & A & C & A & C & A & C & T & A \\ - & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ A & 0 & 2 & 1 & 2 & 1 & 2 & 1 & 0 & 2 \\ G & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 1 \\ C & 0 & 0 & 3 & 2 & 3 & 2 & 3 & 2 & 1 \\ A & 0 & 2 & 2 & 5 & 4 & 5 & 4 & 3 & 4 \\ C & 0 & 1 & 4 & 4 & 7 & 6 & 7 & 6 & 5 \\ A & 0 & 2 & 3 & 6 & 6 & 9 & 8 & 7 & 8 \\ C & 0 & 1 & 4 & 5 & 8 & 8 & 11 & 10 & 9 \\ A & 0 & 2 & 3 & 6 & 7 & 10 & 10 & 10 & 12 \end{pmatrix}$$

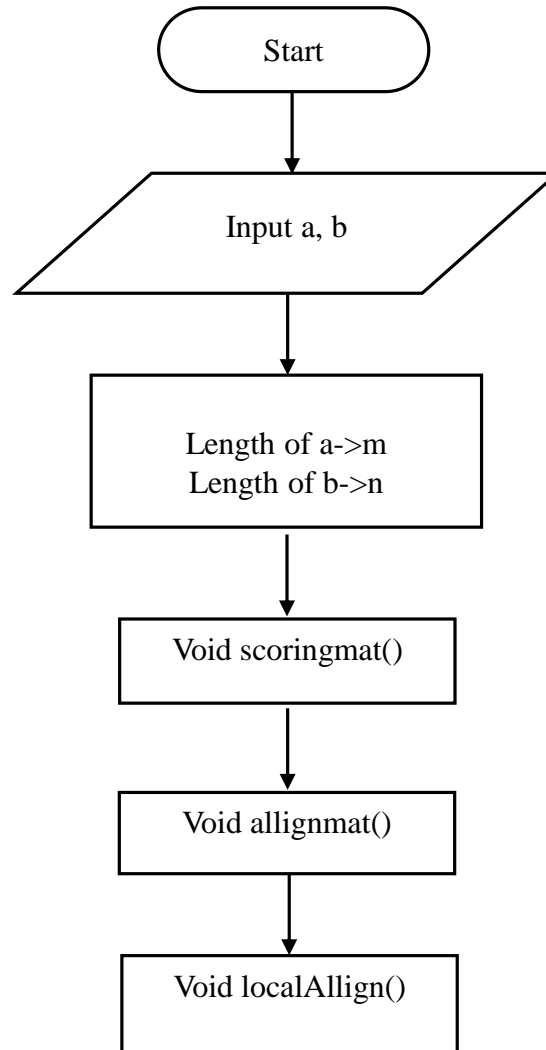
$$T = \begin{pmatrix} - & - & A & C & A & C & A & C & T & A \\ - & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ A & 0 & \swarrow & \leftarrow & \swarrow & \leftarrow & \swarrow & \leftarrow & \leftarrow & \swarrow \\ G & 0 & \uparrow & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \swarrow & \uparrow \\ C & 0 & \uparrow & \swarrow & \leftarrow & \swarrow & \leftarrow & \swarrow & \leftarrow & \leftarrow \\ A & 0 & \swarrow & \uparrow & \swarrow & \leftarrow & \swarrow & \leftarrow & \leftarrow & \swarrow \\ C & 0 & \uparrow & \swarrow & \uparrow & \swarrow & \leftarrow & \swarrow & \leftarrow & \leftarrow \\ A & 0 & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \leftarrow & \leftarrow & \swarrow \\ C & 0 & \uparrow & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \leftarrow & \leftarrow \\ A & 0 & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \uparrow & \swarrow & \swarrow \end{pmatrix}$$

For the example, the results are:

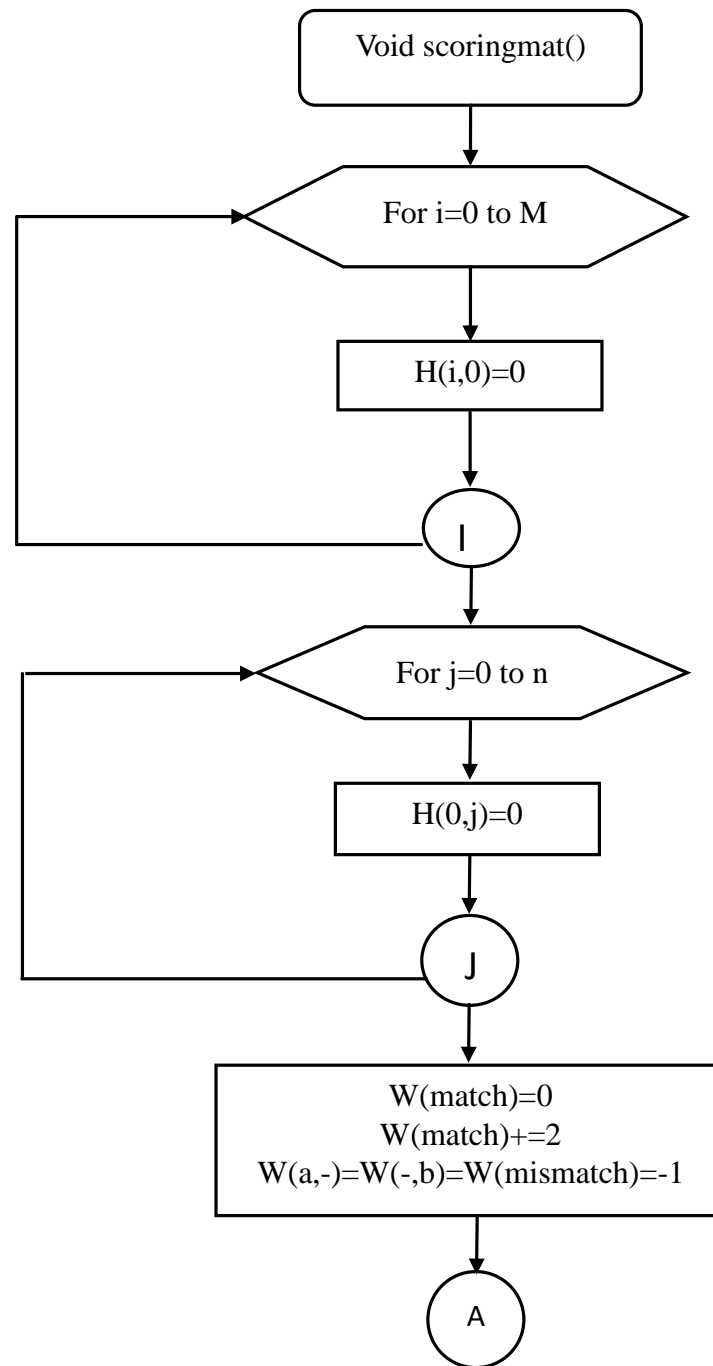
Sequence 1 = A-CACACTA

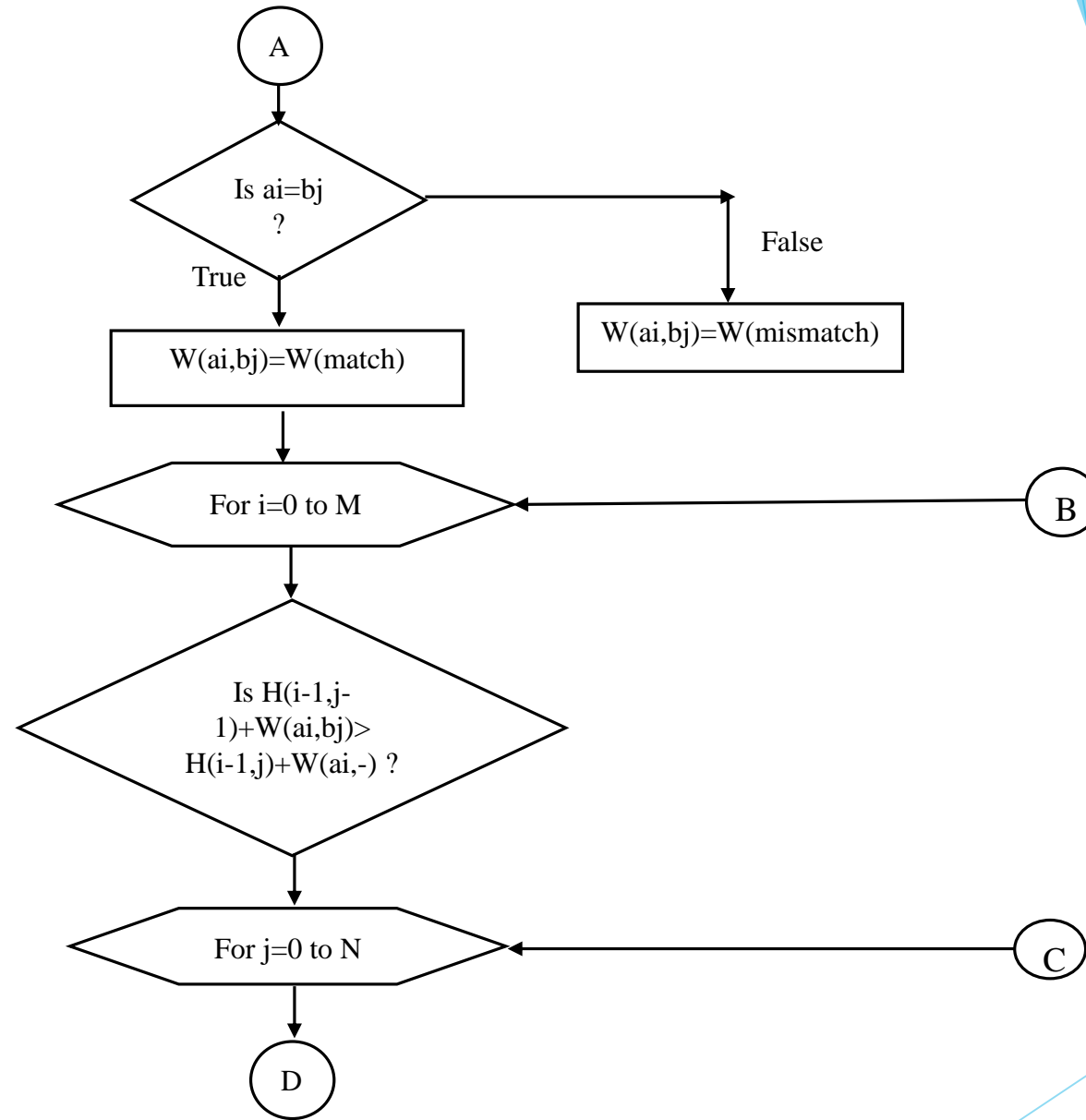
Sequence 2 = AGCACAC-A

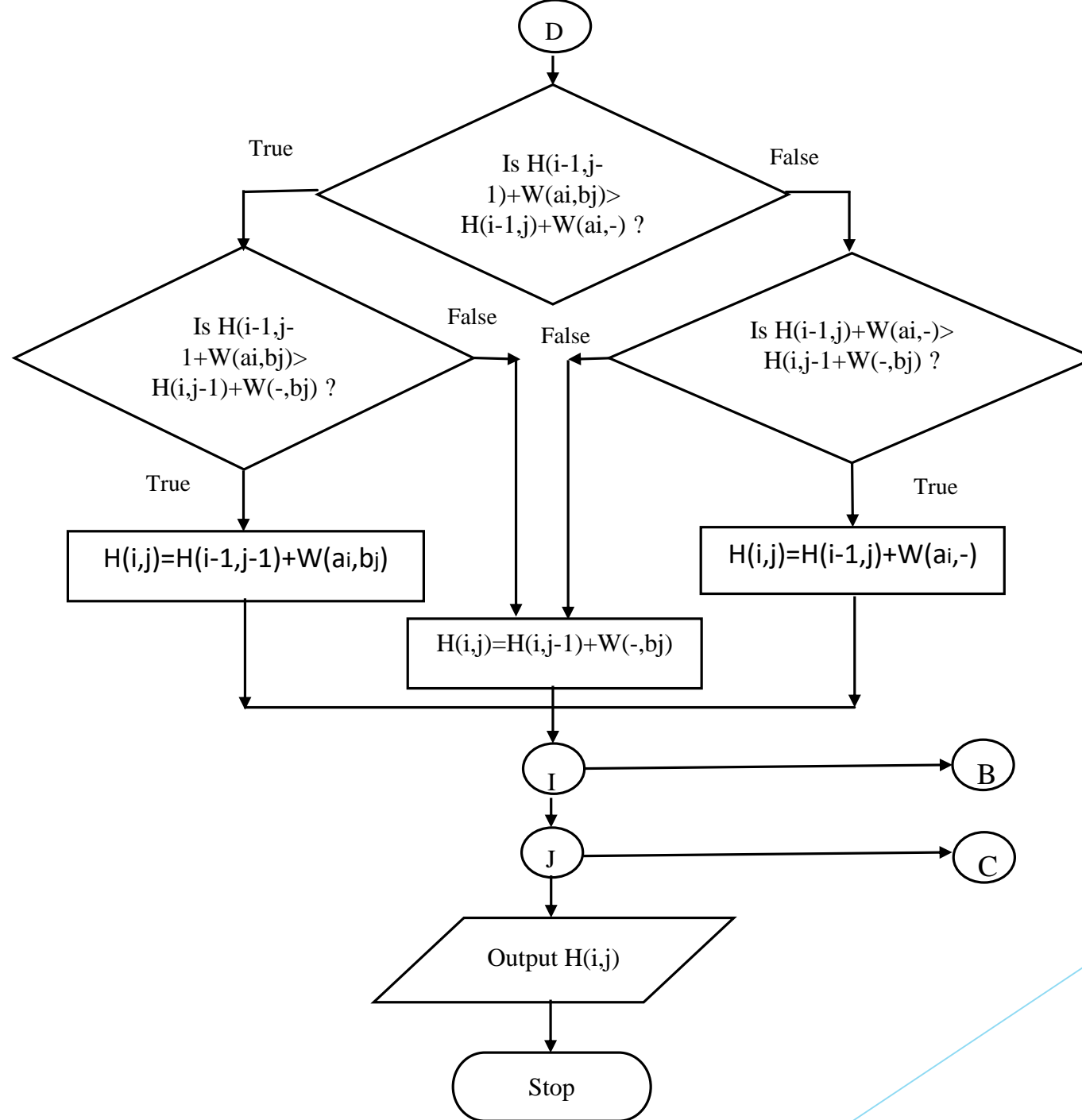
## Flowchart for Smith Waterman Algorithm

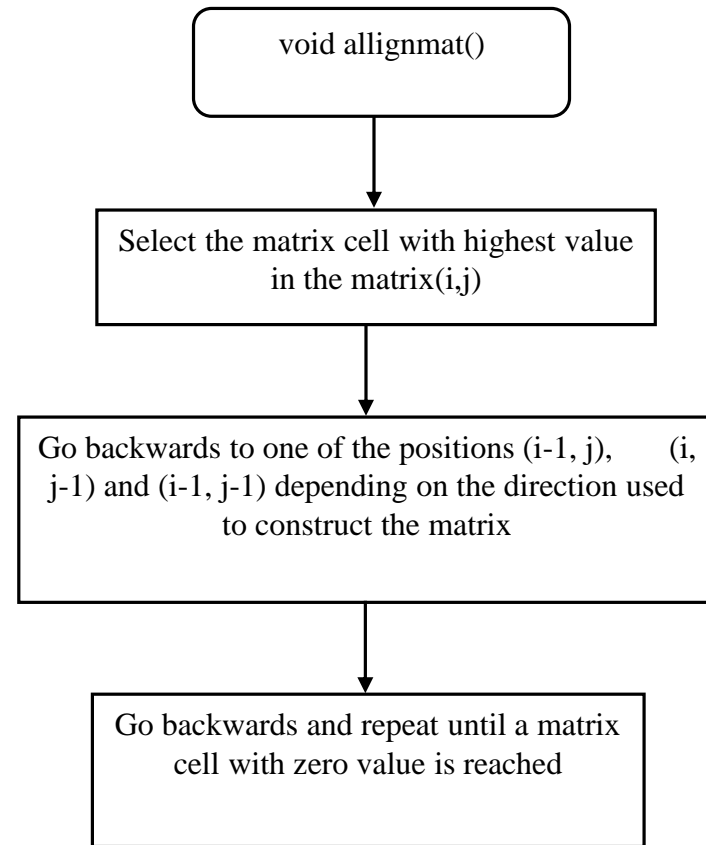


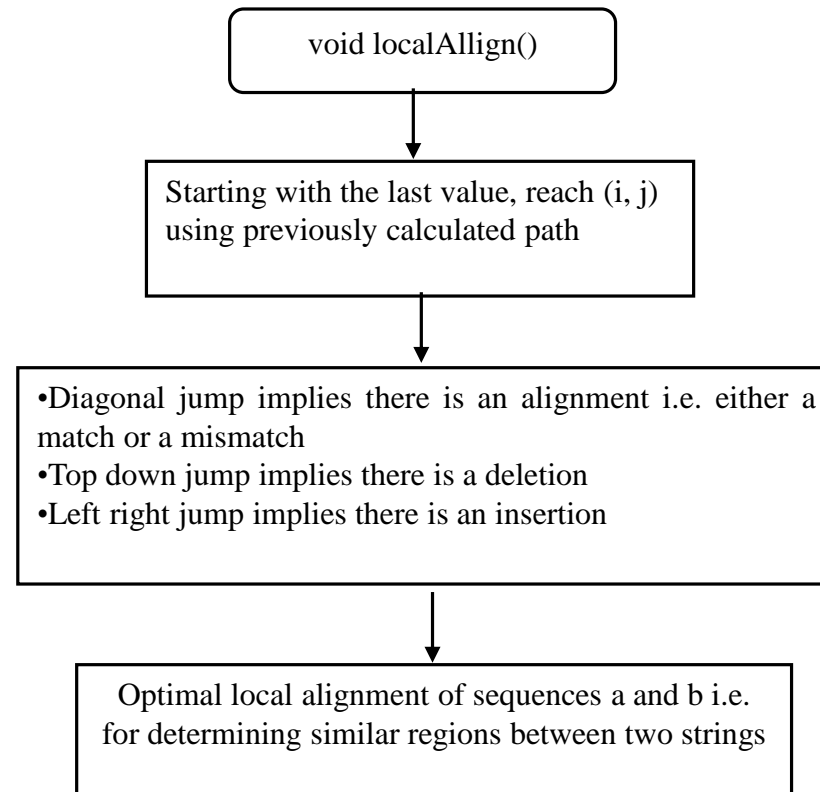






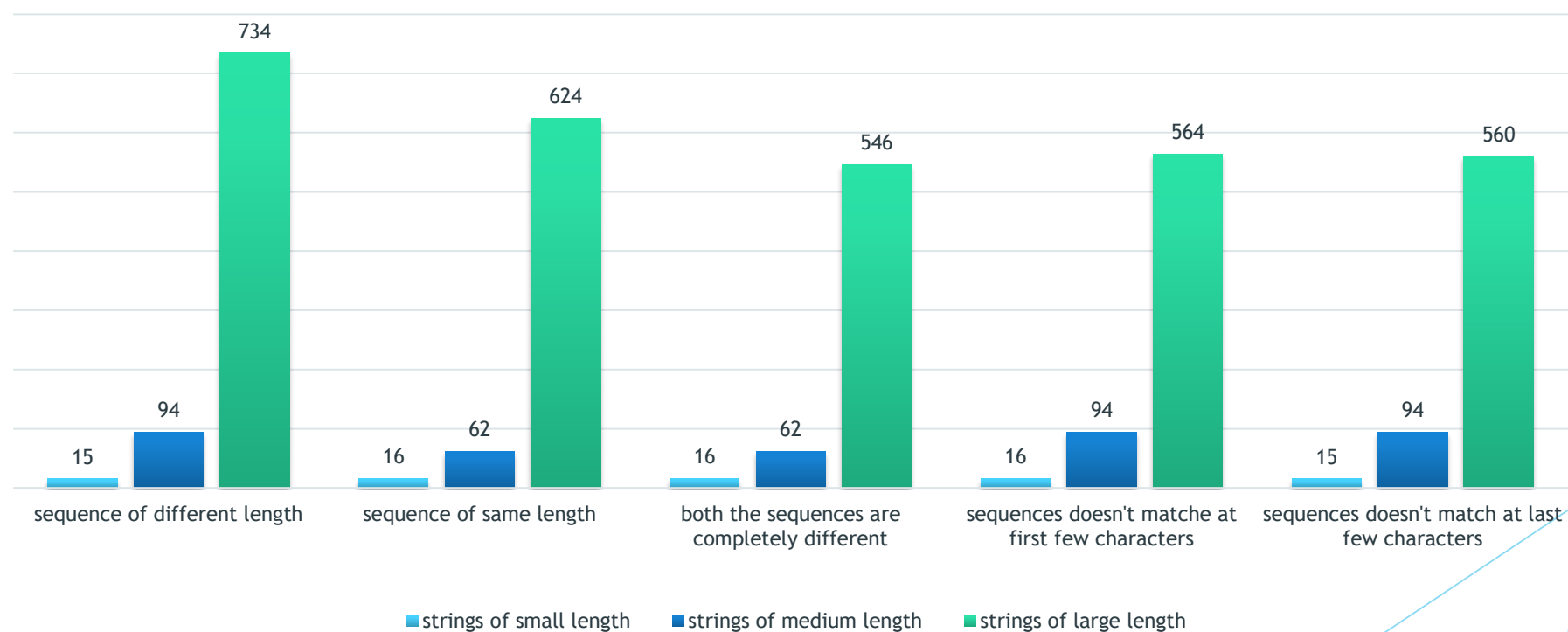






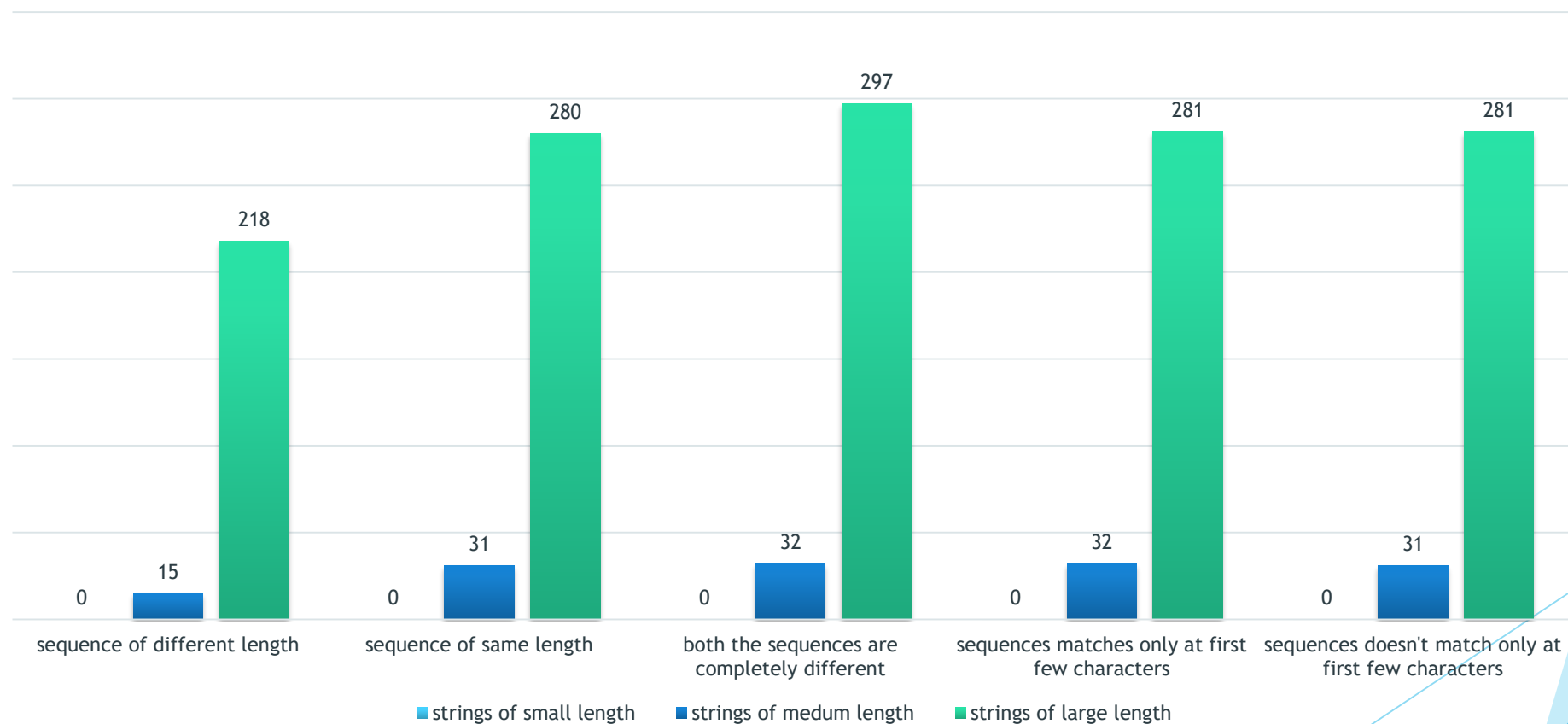
# Non-CUDA Performance Graph

GRAPH FOR SMITH WATERMAN C PROGRAM



# CUDA Performance Graph

GRAPH FOR SMITH WATERMAN CUDA-PROGRAM



# Non-CUDA Output

```
C:\Users\Nitin Kaveriappa\Documents\Visual Studio 2010\Projects\SmithWatermanC\Debug\Smit...
SMITH WATERMAN C PROGRAM

Scoring Matrix
0 0 a g c a c a g a c i o
0 0 0 0 0 0 0 0 0 0 0 0
a 0 2 1 0 2 1 2 1 2 1 0 0
c 0 1 1 3 2 4 3 2 1 4 3 2
a 0 2 1 1 2 5 4 6 5 4 3 2
c 0 1 1 1 3 4 7 6 5 4 6 5
t 0 0 0 2 3 6 6 5 4 5 5 4
a 0 2 1 1 4 5 8 7 7 6 5 4
g 0 1 4 3 3 4 7 10 9 8 7 6
a 0 2 3 3 5 4 6 9 12 11 10 9
c 0 1 2 5 4 7 6 8 11 14 13 12
t 0 0 1 4 4 6 6 7 10 13 13 12
a 0 2 1 3 6 5 8 7 9 12 12 12
matrix value=14 and cell value=9 9

<9,9>
<8,8>
<7,7>
<6,6>
<5,5>
<4,5>
<3,4>
<2,3>
<1,2>
<1,1>

Alignment Matrix
0 0 a g c a c a g a c i o
0 0 0 0 0 0 0 0 0 0 0 0
a 0 2 1 0 0 0 0 0 0 0 0 0
c 0 0 0 3 0 0 0 0 0 0 0 0
a 0 0 0 0 5 0 0 0 0 0 0 0
c 0 0 0 0 0 7 0 0 0 0 0 0
t 0 0 0 0 0 6 0 0 0 0 0 0
a 0 0 0 0 0 0 8 0 0 0 0 0
g 0 0 0 0 0 0 0 10 0 0 0 0
a 0 0 0 0 0 0 0 0 12 0 0 0
c 0 0 0 0 0 0 0 0 0 14 0 0
t 0 0 0 0 0 0 0 0 0 0 0 0
a 0 0 0 0 0 0 0 0 0 0 0 0

Optimal Local Alignment of Sequences A & B
SeqA: a-cactagac
SeqB: agcac-agac

Time elapsed in ms: 125.00
-
```



# CUDA Output

```
C:\Users\Nitin Kaveriappa\Documents\Visual Studio 2010\Projects\SmithWatermanCUDA\Debug\...
SMITH WATERMAN C PROGRAM

Scoring Matrix
0 0 0 0 0 0 0 0 0 0 0 0
a 0 2 1 0 0 0 0 0 0 0 0 0
c 0 0 1 1 3 2 4 3 2 1 4 3 2
a 0 0 2 1 1 2 5 4 6 5 4 3 2
c 0 0 1 1 1 3 4 7 6 5 4 6 5 4
t 0 0 0 0 0 2 3 6 6 5 4 5 5 4
a 0 0 2 1 1 1 4 5 8 7 7 6 5 4
g 0 0 1 4 3 3 4 7 10 9 8 7 6
a 0 0 2 3 3 5 4 6 9 12 11 10 9
c 0 0 1 2 5 4 7 6 8 11 14 13 12
t 0 0 0 1 4 4 6 6 7 10 13 13 12
a 0 0 2 1 3 6 5 8 7 9 12 12 12
matrix value=14 and cell value=9 9

<9,9>
<8,8>
<7,7>
<6,6>
<5,5>
<4,5>
<3,4>
<2,3>
<1,2>
<1,1>
<0,0>

Alignment Matrix
0 0 0 0 0 0 0 0 0 0 0 0
a 0 2 1 0 0 0 0 0 0 0 0 0
c 0 0 0 0 3 0 0 0 0 0 0 0
a 0 0 0 0 0 5 0 0 0 0 0 0
c 0 0 0 0 0 0 7 0 0 0 0 0
t 0 0 0 0 0 0 0 6 0 0 0 0
a 0 0 0 0 0 0 0 8 0 0 0 0
g 0 0 0 0 0 0 0 0 10 0 0 0
a 0 0 0 0 0 0 0 0 12 0 0 0
c 0 0 0 0 0 0 0 0 0 14 0 0
t 0 0 0 0 0 0 0 0 0 0 14 0
a 0 0 0 0 0 0 0 0 0 0 0 0

Time elapsed in ms: 31.00

Optimal Local Alignment of Sequences A & B
SeqA: a-cactagac
SeqB: agcac-agac
```

# Conclusion

- ▶ The use of CUDA for small inputs may actually increase the execution time rather than decreasing it because of the time taken for memory transfers between the CPU and GPU
- ▶ As in the CPU the GPU also has a hardware limit known as the computational capability which is determined by the number of cores contained in the GPU. This restricts the maximum size of the input or number of iterations
- ▶ We have observed that as the computational need of a program increases the execution time decreases when compared with CPU vs GPU, hence improving the overall performance program execution

Thank You.