

Project Report

cuDNN implementation of CNN

Table of Contents

1. Objective
2. Description
3. Project Implementation
4. Project Structure
5. Tools and Libraries to be used
6. Project/Code Walkthrough
7. Comparison between CPU and GPU
8. Steps to run the code

1. Objective

The main objective of this project is a distinct analysis of CNN implementation using CUDA and comparing the performance of a similar kind of code on a CPU.

2. Description

Convolutional Neural Network is a Deep Learning algorithm that takes in input as images and appoints weights and predispositions to significant variables to recognize the components of the image. CNN utilizes an organization of neurons to foresee the image dependent on the trained data. The bigger the CNN organization the more exact it should accomplish including some different factors excessively, for example, input information size, and the streamlining agent being utilized.

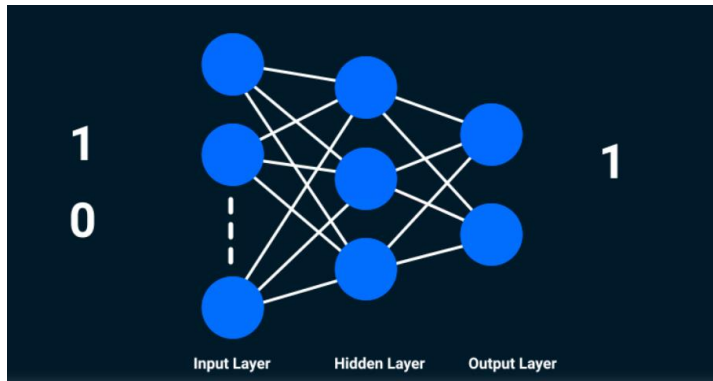
CNN's performance relies on the way the network calculations are done. If they are done in parallel, the output is faster to achieve.

This is where cuDNN implementation arrives at the big picture.

This project is on cuDNN implementation on CNN and a comparison of accuracy rate and time taken on CPU vs GPU

3. Project Implementation

CNN involves three actions – Receiving input, Processing Information and Generating output are represented in the forms of layers as input, hidden and output.



These individual units are called as neurons.

Neural Network training process involves 2 steps.

1. Forward Propagation - Images are taken care of into the input layer as numbers. These mathematical qualities denote the intensity of the pixels in the image. The neurons in the hidden layers apply a couple of numerical procedure on these values. To play out these mathematical tasks, certain parameters are randomly initialized. After these mathematical operations at the hidden layer, the result is sent to the output layer which generates the final prediction values.
2. Backward Propagation- Once the output is generated, in the next step we need to compare the output with the actual value. Based on the final values and the distance from the actual error, the values of the parameter are updated. The forward propagation process is repeated using the updated parameter values and new parameter are generated.

Convolution is often represented mathematically with an asterisk * sign. If we have an input image represented as X and a filter represented with f , then the expression would be: $Z = X * f$; These are the base of Neural Network Algorithm.

Convolutional Neural Network (CNN) Architecture

Complete Network can be broken into 2 parts-

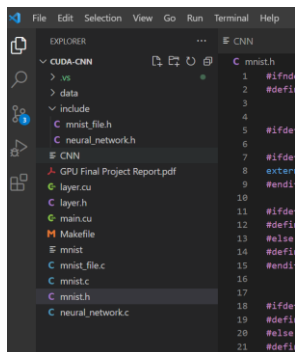
- I. The Convolution layers- Extract features from the input.
- II. Fully connected (Dense) layers- Uses data from the convolution layer to generate output.



3. Project Structure

The architecture of the project is as follows-

- mnist.c (CPU implementation file)
- mnist_file.c (CPU implementation file)
- neural_network.c (CPU implementation file)
- layer.cu (GPU implementation file)
- layer.h (GPU implementation file)
- main.cu (GPU implementation file)
- mnist.h (GPU implementation file)
- /data (MNIST dataset)
- /include(mnist and neural network header files for CPU implementation)



5. Tools and Libraries to be used

cuDNN, Visual Studio Code, C, C++

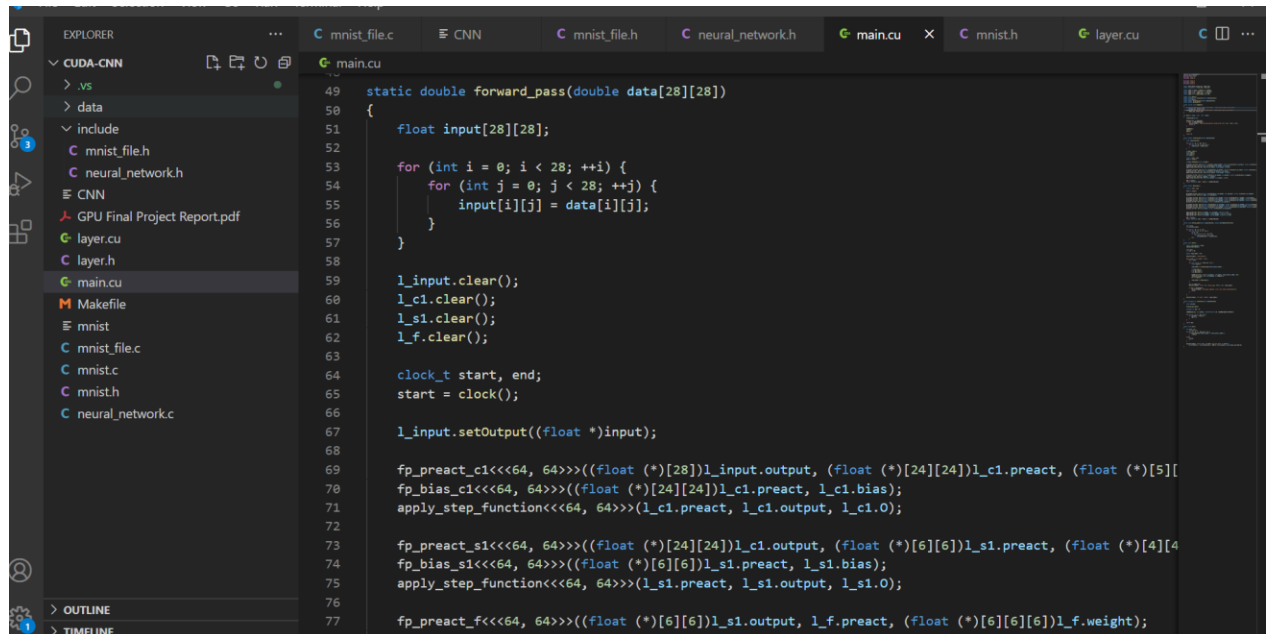
6. Project/Code Walkthrough

GPU gives an advantage of parallelizing the code, thus the parallel method is used by us to implement the forward propagation and backward propagation where the gradient calculations and weight/biases calculations are done in CUDA.

Implementation design for CNN is done in the same manner for CNN and GPU to get a comparison. The network contains 3 layers with one input layer and 2 hidden layers.

CPU implementation calculates the gradient after each layer and serially runs the code. The major difference between the GPU implementation and CPU implementation is the call to write a CUDA code in GPU and feed the weights and biases.

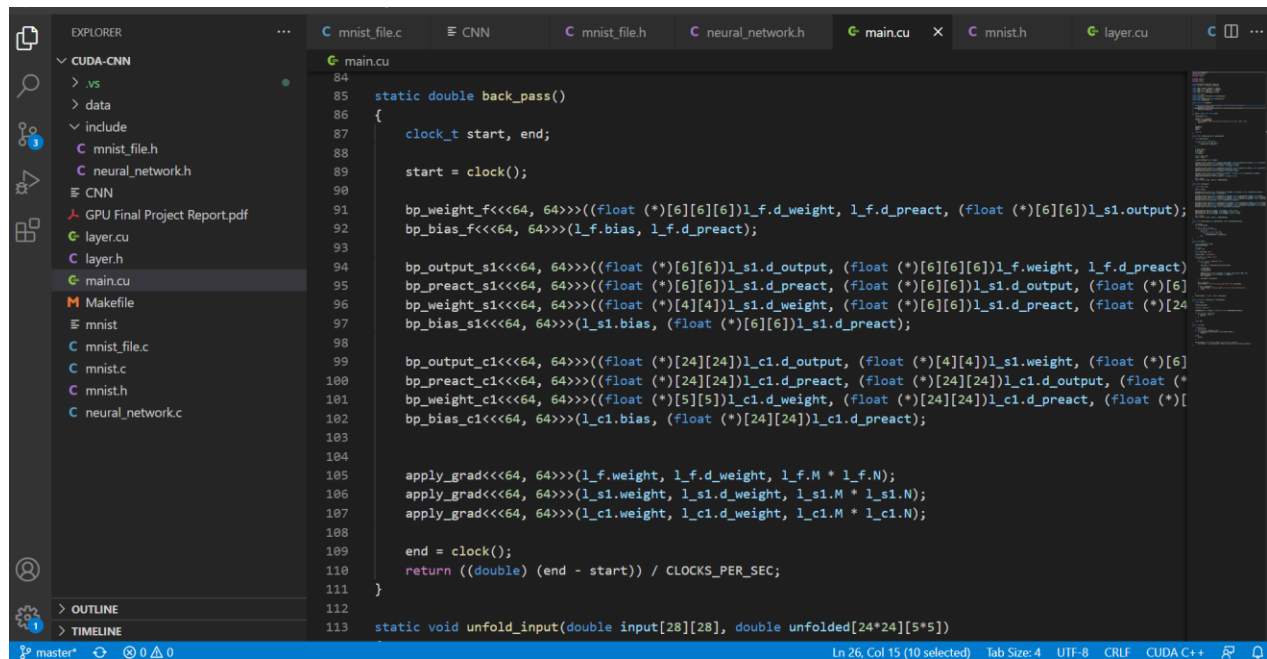
Implementation of Forward Pass



The screenshot displays the Visual Studio Code interface with the 'main.cu' file open. The left sidebar shows the Explorer view with the project structure: CUDA-CNN, .vs, data, include, mnist_file.h, neural_network.h, CNN, GPU Final Project Report.pdf, layer.cu, layer.h, main.cu, Makefile, mnist, mnist_file.c, mnist.c, mnist.h, and neural_network.c. The main editor shows the 'forward_pass' function implementation. The function takes a 28x28 double array 'data' and returns a double value. It initializes a 28x28 float array 'input' and clears the layer inputs and outputs. It then performs a forward pass through the CNN layers, calculating the output of the first layer (l1) and the second layer (l2). The function uses 'clock_t' to measure the execution time and returns the result.

```
49 static double forward_pass(double data[28][28])
50 {
51     float input[28][28];
52
53     for (int i = 0; i < 28; ++i) {
54         for (int j = 0; j < 28; ++j) {
55             input[i][j] = data[i][j];
56         }
57     }
58
59     l_input.clear();
60     l_c1.clear();
61     l_s1.clear();
62     l_f.clear();
63
64     clock_t start, end;
65     start = clock();
66
67     l_input.setOutput((float *)input);
68
69     fp_preact_c1<<<64, 64>>>((float *)l_input.output, (float *)l_c1.preact, (float *)l_s1.preact, (float *)l_c1.bias);
70     apply_step_function<<<64, 64>>>(l_c1.preact, l_c1.output, l_c1.o);
71
72     fp_preact_s1<<<64, 64>>>((float *)l_c1.output, (float *)l_s1.preact, (float *)l_f.preact, (float *)l_s1.bias);
73     apply_step_function<<<64, 64>>>(l_s1.preact, l_s1.output, l_s1.o);
74
75     fp_preact_f<<<64, 64>>>((float *)l_s1.output, l_f.preact, (float *)l_f.weight);
76
77     end = clock();
78     return ((double) (end - start)) / CLOCKS_PER_SEC;
79 }
```

Implementation of Backward Pass



The screenshot displays the Visual Studio Code interface with the 'main.cu' file open. The left sidebar shows the project structure: CUDA-CNN, .vs, data, include, mnist_file.h, neural_network.h, CNN, GPU Final Project Report.pdf, layer.cu, layer.h, main.cu, Makefile, mnist, mnist_file.c, mnist.c, mnist.h, and neural_network.c. The main editor shows the 'back_pass' function implementation. The function takes a 28x28 double array 'input' and a 24x24x5 double array 'unfolded' and returns a double value. It initializes the backward pass by clearing the layer inputs and outputs. It then performs a backward pass through the CNN layers, calculating the gradients of the first layer (l1) and the second layer (l2). The function uses 'clock_t' to measure the execution time and returns the result.

```
84 static double back_pass()
85 {
86     clock_t start, end;
87
88     start = clock();
89
90     bp_weight_f<<<64, 64>>>((float *)l_f.d_weight, l_f.d_preact, (float *)l_s1.output);
91     bp_bias_f<<<64, 64>>>(l_f.bias, l_f.d_preact);
92
93     bp_output_s1<<<64, 64>>>((float *)l_s1.d_output, (float *)l_s1.d_preact, (float *)l_s1.d_output, (float *)l_s1.d_preact);
94     bp_preact_s1<<<64, 64>>>((float *)l_s1.d_preact, (float *)l_s1.d_output, (float *)l_s1.d_preact, (float *)l_s1.d_output);
95     bp_weight_s1<<<64, 64>>>((float *)l_s1.d_weight, (float *)l_s1.d_preact, (float *)l_s1.d_preact, (float *)l_s1.d_preact);
96     bp_bias_s1<<<64, 64>>>(l_s1.bias, (float *)l_s1.d_preact);
97
98     bp_output_c1<<<64, 64>>>((float *)l_c1.d_output, (float *)l_c1.d_preact, (float *)l_c1.d_output, (float *)l_c1.d_preact);
99     bp_preact_c1<<<64, 64>>>((float *)l_c1.d_preact, (float *)l_c1.d_output, (float *)l_c1.d_preact, (float *)l_c1.d_output);
100     bp_weight_c1<<<64, 64>>>((float *)l_c1.d_weight, (float *)l_c1.d_preact, (float *)l_c1.d_preact, (float *)l_c1.d_preact);
101     bp_bias_c1<<<64, 64>>>(l_c1.bias, (float *)l_c1.d_preact);
102
103     apply_grad<<<64, 64>>>(l_f.weight, l_f.d_weight, l_f.M * l_f.N);
104     apply_grad<<<64, 64>>>(l_s1.weight, l_s1.d_weight, l_s1.M * l_s1.N);
105     apply_grad<<<64, 64>>>(l_c1.weight, l_c1.d_weight, l_c1.M * l_c1.N);
106
107     end = clock();
108     return ((double) (end - start)) / CLOCKS_PER_SEC;
109 }
110
111 static void unfold_input(double input[28][28], double unfolded[24*24][5*5])
112
113 }
```

7. Comparison between CPU and GPU

	Accuracy	Time Taken
CPU	81 %	207 seconds
GPU	97 %	172 seconds

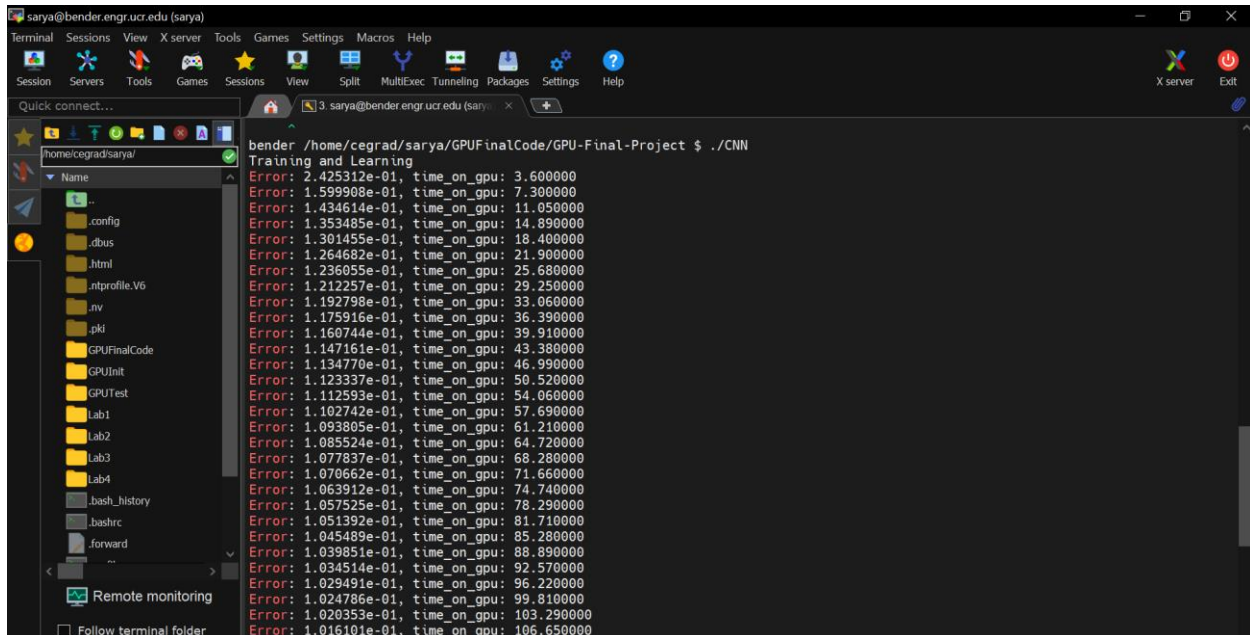
GPU takes less time and provides accuracy to a great extent as compared to CPU implementation. If we increase the number of epochs the result would be significantly different and a larger variation in time taken.

9. Steps to Run the Code

1. Extract all the files from the project folder.
2. Run **make** command inside CuDNN folder.
3. Two files, CNN and mnist will be generated.
4. Run - `./CNN` for GPU implementation.
5. Run- `./mnist` for CPU implementation.

Results produced on the bender

GPU Results



The screenshot shows a terminal window titled 'sarya@bender.engr.ucr.edu (sarya)'. The terminal displays the command `bender /home/cegrad/sarya/GPUFinalCode/GPU-Final-Project $./CNN` and its output. The output indicates 'Training and Learning' and lists 20 error messages, each followed by 'time_on_gpu' values. The values range from 3.600000 to 106.650000. The terminal window also shows a file explorer on the left with a tree view of the project directory, including files like `.config`, `.dbus`, `.html`, `.ntprofile.V6`, `.nv`, `.pkl`, `GPUFinalCode`, `GPUInit`, `GPUtest`, `Lab1`, `Lab2`, `Lab3`, `Lab4`, `.bash_history`, `.bashrc`, and `forward`. The terminal window has a menu bar with options like 'Terminal', 'Sessions', 'View', 'X server', 'Tools', 'Games', 'Settings', 'Macros', 'Help', and a toolbar with icons for 'Session', 'Servers', 'Tools', 'Games', 'Sessions', 'View', 'Split', 'MultiExec', 'Tunneling', 'Packages', 'Settings', and 'Help'.

```
bender /home/cegrad/sarya/GPUFinalCode/GPU-Final-Project $ ./CNN
Training and Learning
Error: 2.425312e-01, time_on_gpu: 3.600000
Error: 1.599908e-01, time_on_gpu: 7.300000
Error: 1.434614e-01, time_on_gpu: 11.090000
Error: 1.353405e-01, time_on_gpu: 14.890000
Error: 1.301455e-01, time_on_gpu: 18.400000
Error: 1.264682e-01, time_on_gpu: 21.900000
Error: 1.236055e-01, time_on_gpu: 25.680000
Error: 1.212257e-01, time_on_gpu: 29.250000
Error: 1.192798e-01, time_on_gpu: 33.060000
Error: 1.175916e-01, time_on_gpu: 36.390000
Error: 1.160744e-01, time_on_gpu: 39.910000
Error: 1.147161e-01, time_on_gpu: 43.380000
Error: 1.134770e-01, time_on_gpu: 46.990000
Error: 1.123337e-01, time_on_gpu: 50.520000
Error: 1.112593e-01, time_on_gpu: 54.060000
Error: 1.102742e-01, time_on_gpu: 57.690000
Error: 1.093805e-01, time_on_gpu: 61.210000
Error: 1.085524e-01, time_on_gpu: 64.720000
Error: 1.077837e-01, time_on_gpu: 68.280000
Error: 1.070662e-01, time_on_gpu: 71.660000
Error: 1.063912e-01, time_on_gpu: 74.740000
Error: 1.057525e-01, time_on_gpu: 78.290000
Error: 1.051392e-01, time_on_gpu: 81.710000
Error: 1.045489e-01, time_on_gpu: 85.280000
Error: 1.039851e-01, time_on_gpu: 88.890000
Error: 1.034514e-01, time_on_gpu: 92.570000
Error: 1.029491e-01, time_on_gpu: 96.220000
Error: 1.024786e-01, time_on_gpu: 99.810000
Error: 1.020353e-01, time_on_gpu: 103.290000
Error: 1.016101e-01, time_on_gpu: 106.650000
```

```
sarya@bender.engr.ucr.edu (sarya)
Terminal Sessions View X server Tools Games Settings Macros Help
Session Servers Tools Games Sessions View Split MultiExec Tunneling Packages Settings Help

Quick connect...
home/cegrad/sarya/
Name
..
.config
.dbus
.html
.ntprofile.V6
.nv
.pki
GPUFinalCode
GPUInit
GPUPTest
Lab1
Lab2
Lab3
Lab4
.bash_history
.bashrc
.forward
Remote monitoring
Follow terminal folder

Error: 1.057525e-01, time_on_gpu: 78.290000
Error: 1.051392e-01, time_on_gpu: 81.710000
Error: 1.045489e-01, time_on_gpu: 85.280000
Error: 1.039851e-01, time_on_gpu: 88.890000
Error: 1.034514e-01, time_on_gpu: 92.570000
Error: 1.029491e-01, time_on_gpu: 96.220000
Error: 1.024786e-01, time_on_gpu: 99.810000
Error: 1.020353e-01, time_on_gpu: 103.290000
Error: 1.016101e-01, time_on_gpu: 106.650000
Error: 1.012061e-01, time_on_gpu: 110.020000
Error: 1.008225e-01, time_on_gpu: 113.560000
Error: 1.004618e-01, time_on_gpu: 117.100000
Error: 1.001191e-01, time_on_gpu: 120.560000
Error: 9.979209e-02, time_on_gpu: 124.260000
Error: 9.947816e-02, time_on_gpu: 127.750000
Error: 9.917667e-02, time_on_gpu: 131.400000
Error: 9.888426e-02, time_on_gpu: 134.750000
Error: 9.860355e-02, time_on_gpu: 138.260000
Error: 9.833176e-02, time_on_gpu: 141.730000
Error: 9.807011e-02, time_on_gpu: 145.210000
Error: 9.781839e-02, time_on_gpu: 148.760000
Error: 9.75554e-02, time_on_gpu: 152.300000
Error: 9.734108e-02, time_on_gpu: 155.800000
Error: 9.711352e-02, time_on_gpu: 159.330000
Error: 9.689328e-02, time_on_gpu: 162.880000
Error: 9.667827e-02, time_on_gpu: 166.380000
Error: 9.646992e-02, time_on_gpu: 170.180000
Error: 9.626785e-02, time_on_gpu: 173.830000
Error: 9.607536e-02, time_on_gpu: 177.510000

Time Taken - 177.510000
Approx Error Rate: 2.88% Approx Accuracy Rate: 97.12%
bender /home/cegrad/sarya/GPUFinalCode/GPU-Final-Project $
```

CPU Results

```
sarya@bender.engr.ucr.edu (sarya)
Terminal Sessions View X server Tools Games Settings Macros Help
Session Servers Tools Games Sessions View Split MultiExec Tunneling Packages Settings Help

Quick connect...
home/cegrad/sarya/
Name
..
.config
.dbus
.html
.ntprofile.V6
.nv
.pki
GPUFinalCode
GPUInit
GPUPTest
Lab1
Lab2
Lab3
Lab4
.bash_history
.bashrc
.forward
Remote monitoring
Follow terminal folder

bender /home/cegrad/sarya/GPUFinalCode/GPU-Final-Project $ ./mnist
error: 4.355595e+00 Time taken on CPU: 4.020000
error: 3.417031e+00 Time taken on CPU: 9.040000
error: 2.970093e+00 Time taken on CPU: 13.060000
error: 2.530892e+00 Time taken on CPU: 18.090000
error: 2.186754e+00 Time taken on CPU: 23.110000
error: 2.083432e+00 Time taken on CPU: 27.130000
error: 1.725454e+00 Time taken on CPU: 31.150000
error: 1.512486e+00 Time taken on CPU: 34.170000
error: 1.570190e+00 Time taken on CPU: 37.190000
error: 1.450604e+00 Time taken on CPU: 40.220000
error: 1.775884e+00 Time taken on CPU: 44.250000
error: 1.678123e+00 Time taken on CPU: 49.270000
error: 1.213438e+00 Time taken on CPU: 52.290000
error: 1.891769e+00 Time taken on CPU: 55.320000
error: 1.431096e+00 Time taken on CPU: 60.340000
error: 1.109908e+00 Time taken on CPU: 64.360000
error: 8.951116e-01 Time taken on CPU: 67.380000
error: 6.212723e-01 Time taken on CPU: 71.400000
error: 8.720776e-01 Time taken on CPU: 75.420000
error: 9.466262e-01 Time taken on CPU: 80.440000
error: 9.007803e-01 Time taken on CPU: 84.460000
error: 6.842788e-01 Time taken on CPU: 88.480000
error: 8.234560e-01 Time taken on CPU: 93.510000
error: 7.486557e-01 Time taken on CPU: 97.530000
error: 8.676123e-01 Time taken on CPU: 102.550000
error: 7.877934e-01 Time taken on CPU: 105.570000
error: 8.052484e-01 Time taken on CPU: 108.590000
error: 8.587374e-01 Time taken on CPU: 112.620000
error: 6.209297e-01 Time taken on CPU: 116.640000
error: 5.852945e-01 Time taken on CPU: 119.660000
error: 9.446012e-01 Time taken on CPU: 122.680000
error: 6.297352e-01 Time taken on CPU: 127.700000
```

```
sarya@bender.engr.ucr.edu (sarya)
Terminal Sessions View X server Tools Games Settings Macros Help
Session Servers Tools Games Sessions View Split MultiExec Tunneling Packages Settings Help

Quick connect...
home/cegrad/sarya/
Name
..
.config
.dbus
.html
.ntprofile.V6
.nv
.pki
GPUFinalCode
GPUInit
GPUPTest
Lab1
Lab2
Lab3
Lab4
.bash_history
.bashrc
.forward
Remote monitoring
Follow terminal folder

error: 9.007803e-01 Time taken on CPU: 84.460000
error: 6.842788e-01 Time taken on CPU: 88.480000
error: 8.234560e-01 Time taken on CPU: 93.510000
error: 7.486557e-01 Time taken on CPU: 97.530000
error: 8.676123e-01 Time taken on CPU: 102.550000
error: 7.877934e-01 Time taken on CPU: 105.570000
error: 8.052484e-01 Time taken on CPU: 108.590000
error: 8.587374e-01 Time taken on CPU: 112.620000
error: 6.209297e-01 Time taken on CPU: 116.640000
error: 5.852945e-01 Time taken on CPU: 119.660000
error: 9.446012e-01 Time taken on CPU: 122.680000
error: 6.297352e-01 Time taken on CPU: 127.700000
error: 7.206808e-01 Time taken on CPU: 130.720000
error: 6.635219e-01 Time taken on CPU: 133.750000
error: 7.065411e-01 Time taken on CPU: 138.770000
error: 8.285132e-01 Time taken on CPU: 142.800000
error: 6.355624e-01 Time taken on CPU: 147.830000
error: 7.871413e-01 Time taken on CPU: 151.850000
error: 5.782230e-01 Time taken on CPU: 156.870000
error: 5.730005e-01 Time taken on CPU: 161.890000
error: 6.683274e-01 Time taken on CPU: 164.920000
error: 7.366546e-01 Time taken on CPU: 169.940000
error: 6.011595e-01 Time taken on CPU: 174.960000
error: 5.917724e-01 Time taken on CPU: 179.990000
error: 5.855922e-01 Time taken on CPU: 185.020000
error: 4.473346e-01 Time taken on CPU: 189.040000
error: 7.380794e-01 Time taken on CPU: 194.070000
error: 7.402201e-01 Time taken on CPU: 199.090000
error: 4.769174e-01 Time taken on CPU: 202.110000
error: 6.943619e-01 Time taken on CPU: 207.140000
Final Time - 207.140000
Error Rate: 69.44% Accuracy Rate: 81.53%
bender /home/cegrad/sarya/GPUFinalCode/GPU-Final-Project $
```