
SIMD Optimization to RBF (with Gaussian/Cosine adaptive fusion) Neural Network Using Cuda Framework



Work Contribution

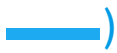
The work we have done for preparing reports and coding was mutually shared among all the project team. Here are some details of and contribution in percent of-of working

Working area	Anees Ahmed	Arif Sultan	Waqar Hameed
	Contribution %		
Problem selection and information Collection	30	40	30
Initial review of papers and papers selection	40	40	20
Conference-related working	50	25	25
Algorithm Comprehension and Adaption	20	60	20
Problem Analysis and Design	20	50	20
Design Analysis	25	20	50
Coding and Results	20	20	60
Codes review	30	30	40
Results Review	15	15	70
Report Writing	60	20	20
Report Review	40	30	30
Total Contribution Score	350	350	385



Table Of Content

Work Contribution	1
Abstract	2
Introduction	3
Artificial Neural Networks (ANN)	4
RBF Neural Networks (RBF-NN)	5
Kernel Function Design	6
Description	6
SIMD Optimization Approach	6
The Cuda Indexing Mechanism	6
Cuda C Codes	7
Network Design	8
Description	8
SIMD optimization approach.	8
Cuda C Code	8
Kernel Launch	9
Description	9
GPU optimized Helper Functions:	9
Kernel Visual Profiler	10
Kernels Kernel Time	10
Output Neuron Kernel Time	11
Multiplication kernel Time	12
AlphaUpdate Kernel Time	13
Conclusion:	16
References:	17



Abstract

Single Instruction Multiple Data (SIMD) is a good applicable choice where a grid of data available and we need to apply same computation to all data, like adjusting digital media, scaling digital media and manipulating matrices in Linear Algebra or Statistics or other computational work.

In this paper, we focus on Radial Based Function Network (Neural Network) for function approximation which is an ideal case for SIMD applicability.

We use Nvidia's Cuda framework for implementing SIMD using GPU. Most of the codes are in C/C++

Introduction

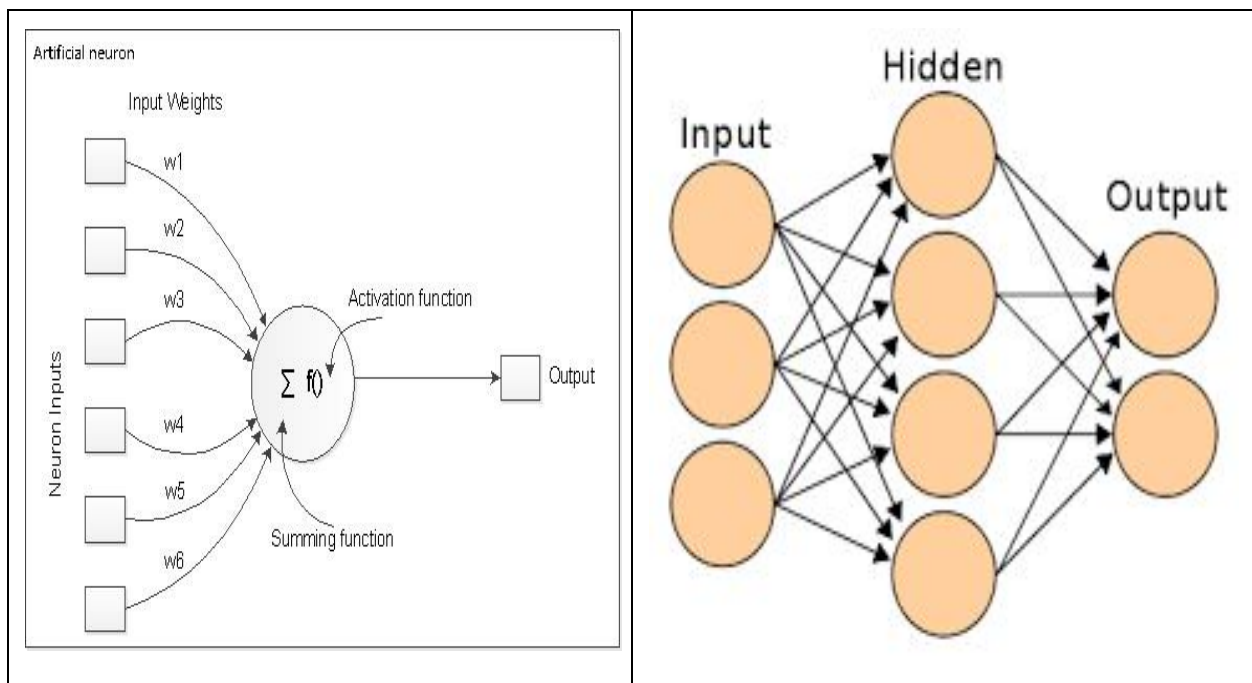
Single Instruction Multiple Data (SIMD) is an approach to parallel computation. It refers to multiple computing or processing units that perform a single operation (computing instruction) on multiple data elements simultaneously. This is also treated as data level parallelism, however, it is different if compare with concurrency. In SIMD only a single process (instruction) is available to all computation unit at a moment (1).

We should not confuse with SIMT which utilizes threads or CPU concurrency that utilizes scheduling and time slicing multiple cores of a CPU.

Using SIMD approach could bring tremendous advantages in computing by reducing computation time especially working with matrices like structure where the parallel calculation is a need and most calculation are not dependent on each other.

Artificial Neural Networks (ANN)

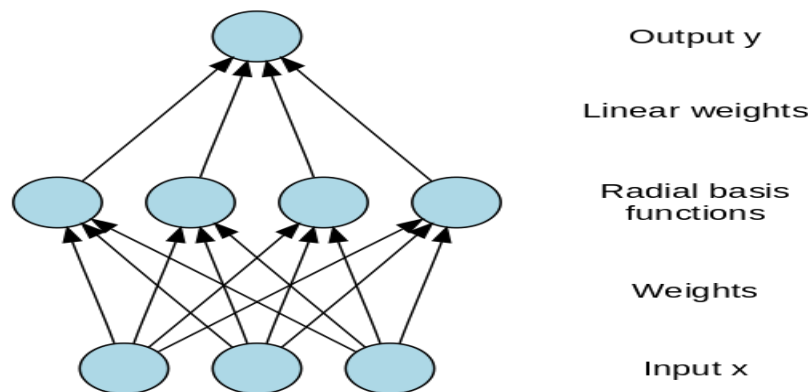
An Artificial Neural Network (ANN) is an information processing paradigm that is inspired by the way biological nervous systems, such as the brain, process information. The key element of this paradigm is the novel structure of the information processing system. It is composed of a large number of highly interconnected processing elements (neurons) working in unison to solve specific problems. ANNs, like people, learn by example. An ANN is configured for a specific application, such as pattern recognition or data classification, through a learning process. Learning in biological systems involves adjustments to the synaptic connections that exist between the neurons. This is true of ANNs as well (4).



RBF Neural Networks (RBF-NN)

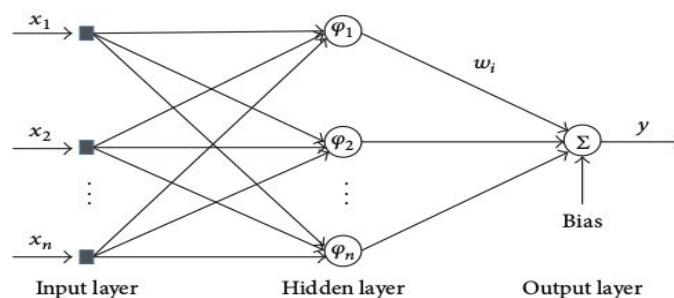
Radial Basis Functions, as a variant of Artificial Neural Network (ANN), start getting attraction in late 80 (1). They are mainly used in pattern recognition techniques but are also used for clustering, functional approximation, spline interpolation etc (2).

An RBF network has two layers of the neural network. The hidden unit implements a radial activated function while output layers of the neural network implement a weighted sum of previous layer output. The output of RBF-NN is linear, while input into the RBF is nonlinear. The nonlinear approximation properties of RBF-NN, we can model complex mappings which perceptron neural networks can only model by means of multiple intermediary layers (4).



The RBF-NN implementation is divided into different steps, like designing the kernel, data pre-processing, training, testing, and approximation.

The architecture consists of multiple layers and input layer, a nonlinear hidden layer and a linear output layer (5)





Kernel Function Design

Description

Our RBF-NN kernel is a fusion of cosine and Euclidean distances. Creating a fusion of both distance function, we get a better result as compare with the conventional approach where mostly a single function is used (5). This fusion is adaptive in nature and provides a robust result during training as an activation function (5).

$$\varphi_i(x, x_i) = \frac{\alpha_1(n) \varphi_{i1}(x, x_i) + \alpha_2(n) \varphi_{i2}(kx - x_i k)}{\alpha_1(n) + \alpha_2(n)}$$

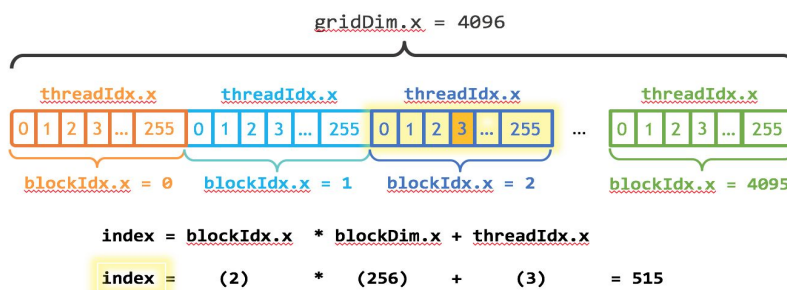
where $\varphi_{i1}(x, x_i)$ and $\varphi_{i2}(kx - x_i k)$ are the cosines and Euclidean kernels.

The kernel is implemented as sequential as following tables and can be modified with SIMD optimization. We can see the approach for optimization is straightforward. The kernel function code using an index that points to a specific thread running parallel with other threads.

SIMD Optimization Approach

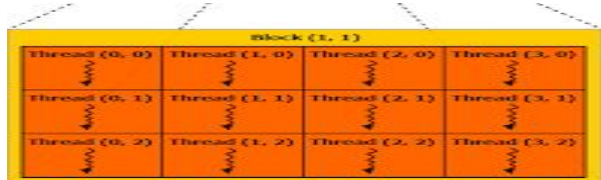
The GPU optimization is straightforward. Our optimized kernel function modified to execute instructions over a data element indexed by an indexer. The index position is calculated with the help of Cuda built-in helper variables the provides a specific pointer the particular thread inside the execution block (7).

The Cuda Indexing Mechanism



Cuda C Codes

Complete codes can be reviewed at the GitHub repository (8).

Sequential	SIMD optimized
<pre> void GaussianKernal(float x, float y, int CenterR, int CenterC, float Centers[][121], float* output) { // printf("Gauss Kernel\n\n\n"); // for (int i = 0; i < 121; i++) // { // output[i] = exp(-(pow((x - Centers[0][i]), 2) + pow((y - Centers[1][i]), 2))/0.04); // printf("%f\n", output[i]); // } // // // //float output[121]; //float sumCenter[121]; //float inputsq=x*x+y*y; //printf("\nMultiplication Kernel\n\n\n"); // // for (int i = 0; i < 121; i++) // { // float sum = 0.0; // sum = x * Centers[0][i]+ y * Centers[1][i]; // output[i] = sum; // // sumCenter[i] = sqrt((pow(Centers[0][i], 2) + pow(Centers[1][i], 2))*inputsq); // output[i] = output[i] / (sumCenter[i]+0.0000000000000001); // printf("%f\n", output[i]); // } // } </pre>	<pre> __global__ void Gauss(float* x, float* y, float* CenterX, float* CenterY, float* output, int N) { int i = blockDim.x*blockIdx.x + threadIdx.x; //printf("x= %f, y=%f \n", x[0], y[0]); if (i < N) { output[i] = exp(-(pow((x[0] - CenterX[i]), 2) + pow((y[0] - CenterY[i]), 2)) / 0.04); printf("%d: %f\n", i, output[i]); } } __global__ void Coss(float* x, float* y, float* CenterX, float* CenterY, float* output, int N) { float sumCenter; float inputsq = x[0]*x[0] + y[0]*y[0]; //printf("%f", inputsq); // Cuda Helper int i = blockDim.x*blockIdx.x + threadIdx.x; if (i < N) { output[i] = (x[0] * CenterX[i] + y[0] * CenterY[i]) / (sqrt((pow(CenterX[i], 2) + pow(CenterY[i], 2))*inputsq) + 0.0000000000000001); } } </pre> 



Network Design

Description

The RBF-NN network consists of a large set of neurons. Each neuron is responsible to apply kernel function (gaussian, cosine) with each of inputs and obtains a sum of all values for feeding the next layer

SIMD optimization approach.

Unlike a sequential approach where a nested loop is used to sum of the function output for feeding to next layer, The optimized codes call kernel parallel using Cuda helper variables. Each call is runs parallel to the sum operation performed by many threads in the time. A full set is input is collected and send for the paralleled processing so multiple data elements are calculated at the same time with different threads and asynchronously result move back to the host memory.

Cuda C Code

codes are abstracted. For complete codes please review the GitHub repository of the project (8).

```
void train2(int n, float l_rate, float error, float *RBF_out, float *w)
{
    int i = blockIdx.x*blockDim.x + threadIdx.x;
    if (i < n)
    {

        w[i] = w[i] + l_rate*error*RBF_out[i];

    }
    __syncthreads();
}
```



Kernel Launch

Description

The RBF-NN network launch requires many steps which include data generation (for testing) and loading into required structures, training using such data via RBF-NN and for training and predicting call the kernel.

Many code blocks require to run on the host in a sequential manner as their parallel implementation does not bring any optimization.

However, many helper functions like `blew` are optimized the use SIMD and run over the GPU (8).

GPU optimized Helper Functions:

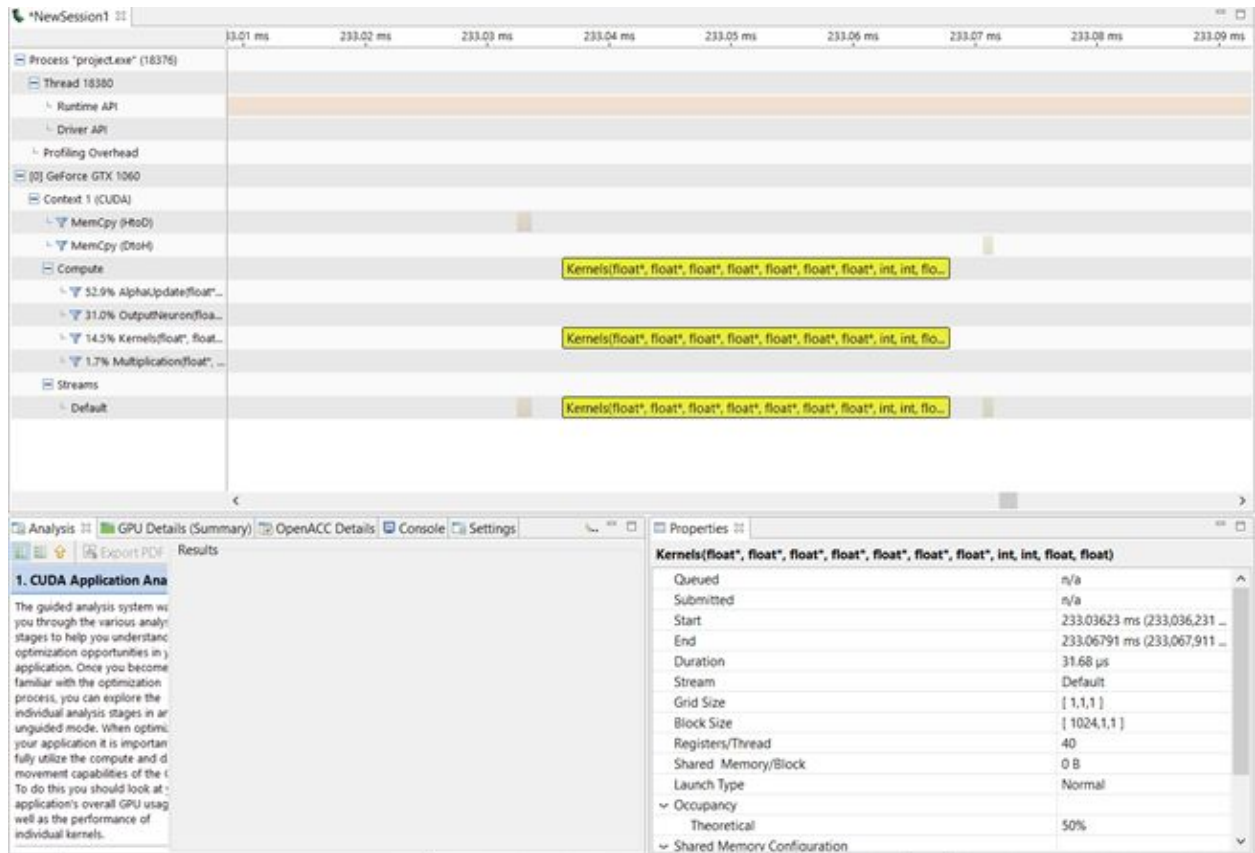
- `void OutputNeuron(float* Kernel, float* w, float* output, int N)`
- `void Multiplication(float* Kernel, float* w, float error, float learningRate, int N)`
- `void AlphaUpdate(float* KC, float* KG, float* w, float* updateAlpha1, float* updateAlpha2)`



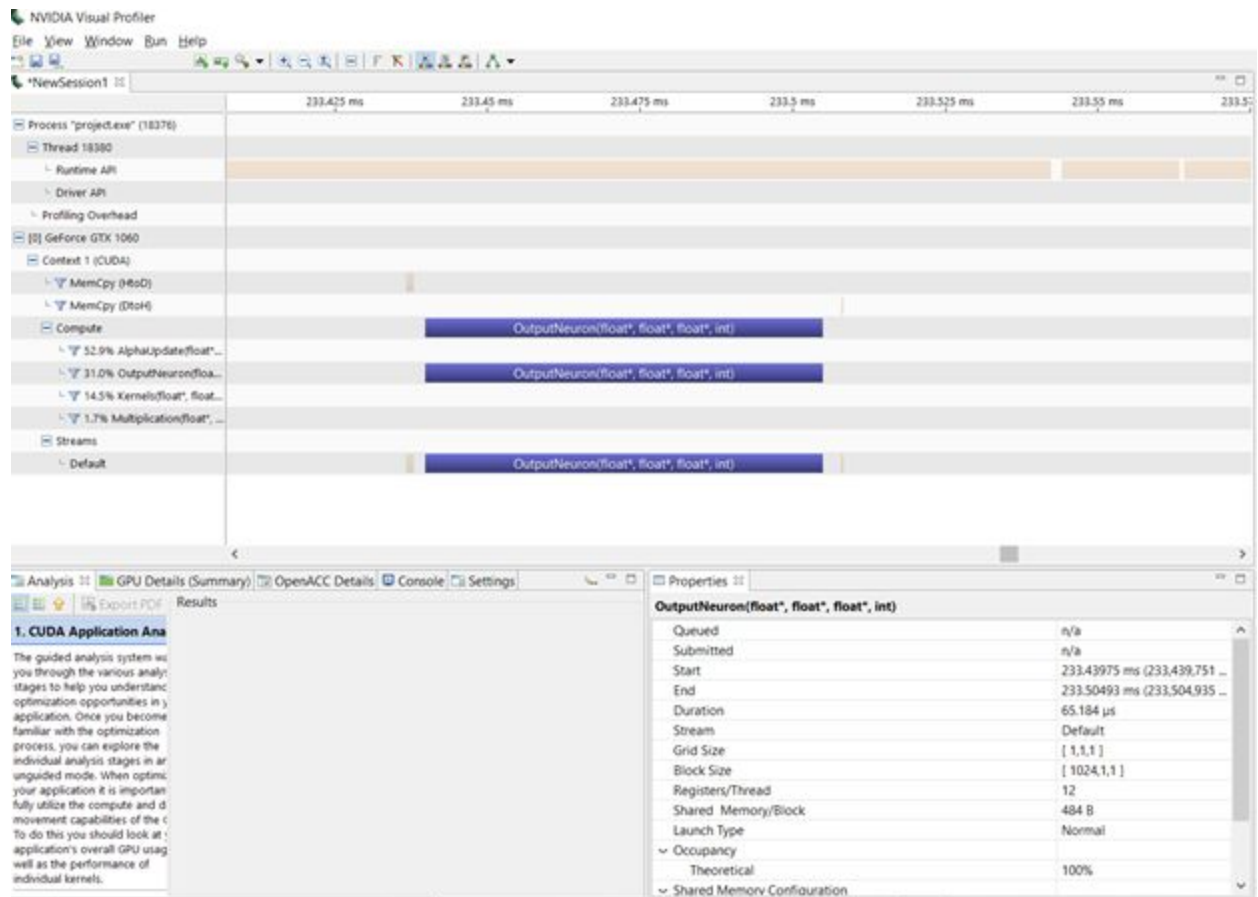
Kernel Visual Profiler

Below Images shows the profile of individual kernel execution time

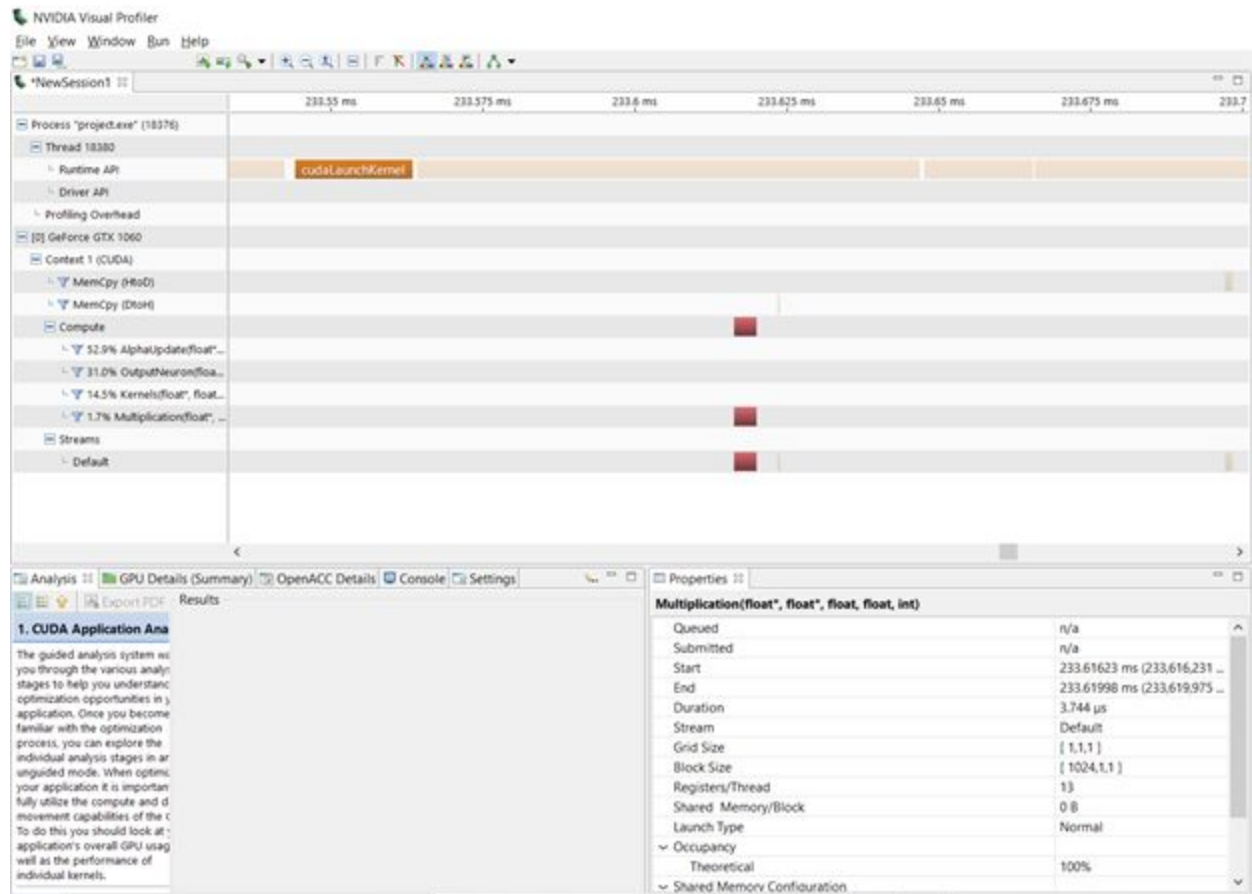
Kernels Kernel Time



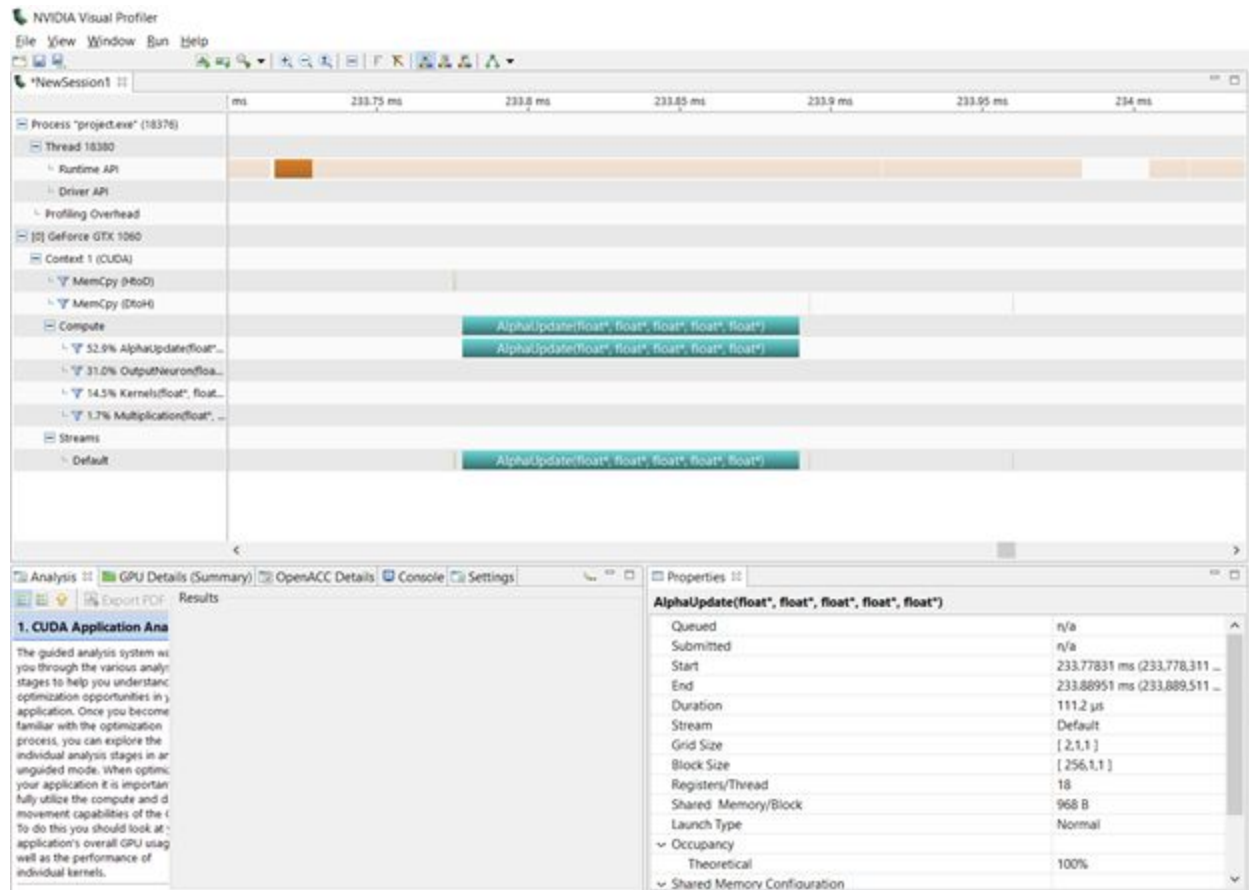
Output Neuron Kernel Time



Multiplication kernel Time



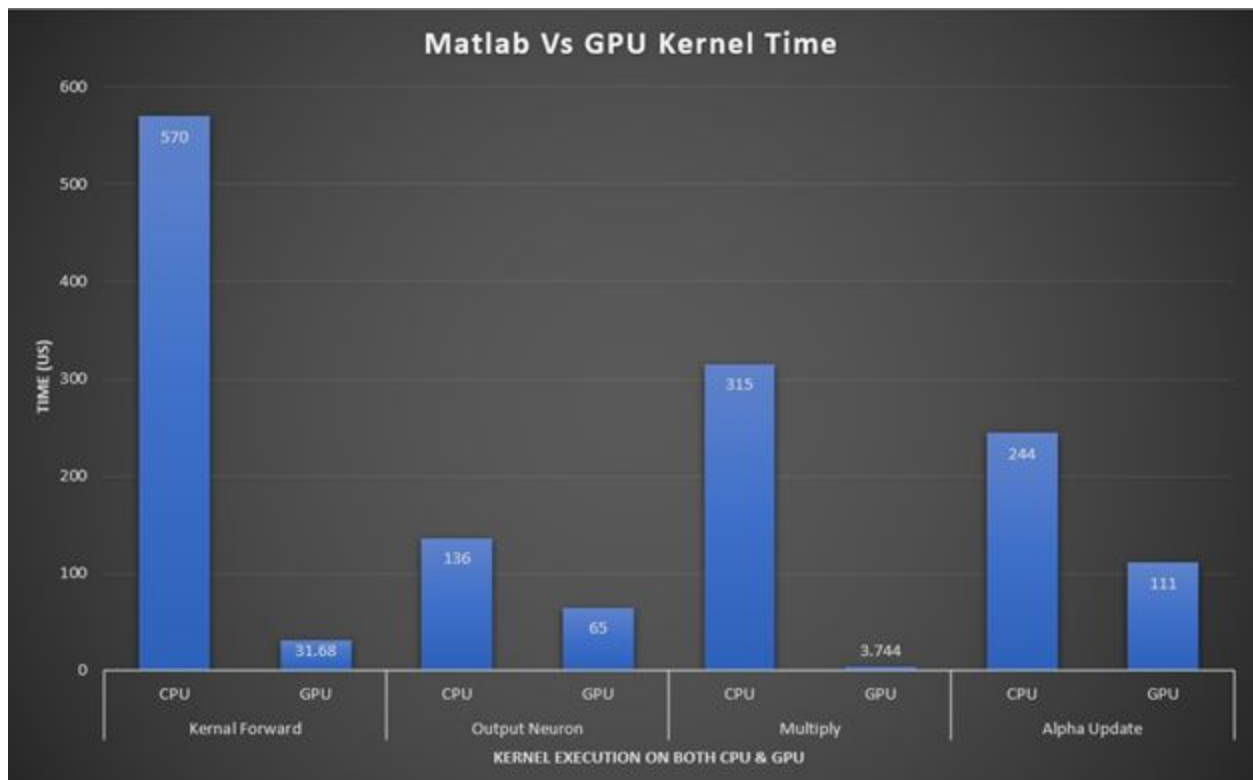
AlphaUpdate Kernel Time





Above Execution time we now able to compare it results with Matlab sequential code Since we develop sequential code in C language as well but visual studio diagnostics toolbox only give time in milliseconds. So we don't consider C language time.

Below bars show a comparison between Matlab sequential and Cuda Parallel execution task

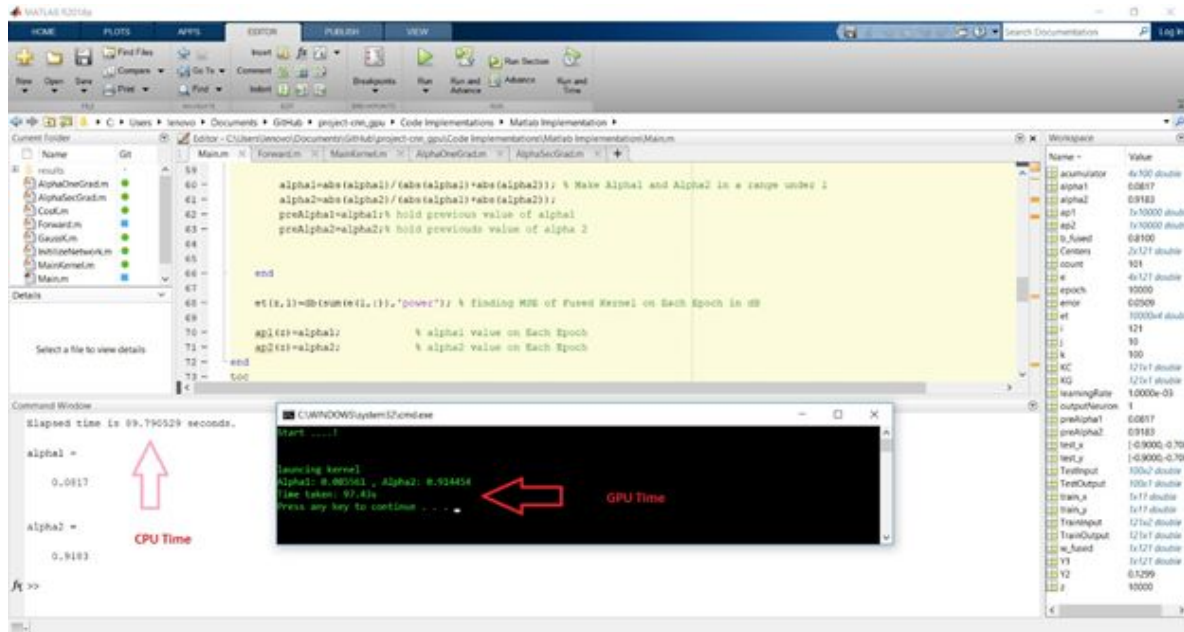


From above it is clear that we get kernel execution gain more than 2X with GPU. Following is hardware specs of the system

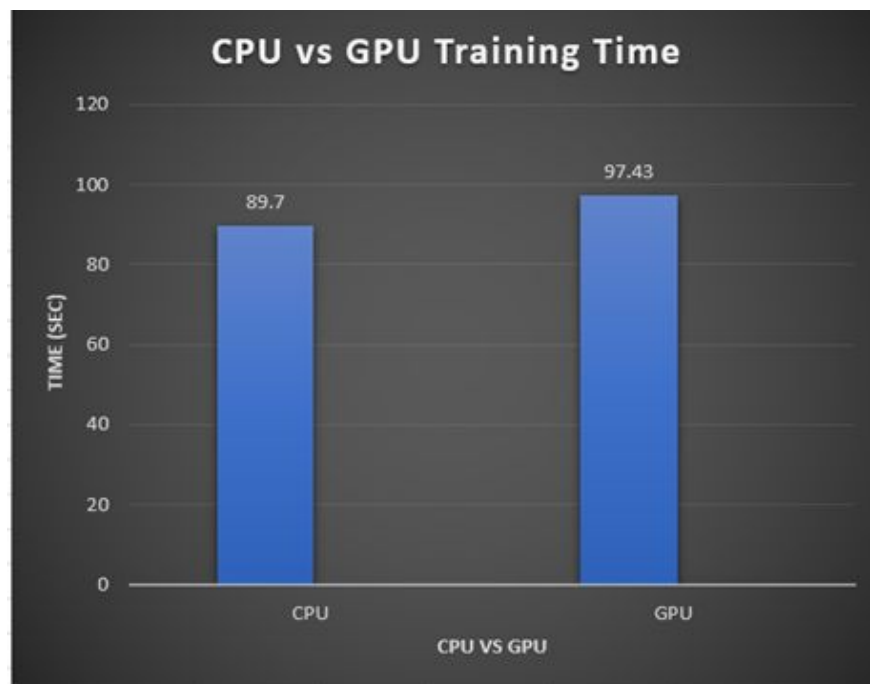
- Intel Core i7 -7700HQ (7th Generation) 3.8 GHz with 4 Core & 8 Logical Processor
- 16 Gb Ram (2400 Mhz)
- Nvidia 1060 GPU

Final Results:

After training on both the device we get result as follow



Below bar Graph shows comparison





Conclusion:

From above results, we find out that Cuda Parallel Kernel executer much faster than CPU but if we compare complete training results. We found that CPU perform better than GPU, because of too much memory transfer operation between the device to host and host to the device.

We can overcome this problem by using constant memory location. Because we have 121×2 vectors of center data. We can also improve its performance by getting burst data from global to shared memory then execute kernel.



References:

1. SIMD architectures | Ars Technica." 21 Mar. 2000,
<https://arstechnica.com/features/2000/03/simd/>. Accessed 7 Dec. 2018.
2. Introduction of the Radial Basis Function (RBF) Networks. Retrieved December 7, 2018, from
https://www.researchgate.net/profile/Adrian_Bors/publication/280445892_Introduction_of_the_Radial_Basis_Function_RBF_Networks/links/585f0e4108ae6eb871a31b01/Introduction-of-the-Radial-Basis-Function-RBF-Networks.pdf
3. Haykin, S. (1994) Neural Networks: A Comprehensive Foundation. Upper Saddle River, NJ: Prentice Hall.
4. Neural Network Definitions
5. https://www.doc.ic.ac.uk/~nd/surprise_96/journal/vol4/cs11/report.html#What%20is%20a%20Neural%20Network
6. A Novel Adaptive Kernel for the RBF Neural Networks, Shujaat Khan · Imran Naseem · Roberto Togneri · Mohammed Bennamoun
7. A Novel Kernel for RBF Based Neural Networks Wasim Aftab, Muhammad Moinuddin, 1 and Muhammad Shafique Shaikh
8. <https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#programming-model>
Accessed 8 Dec 2018
9. Project Codes
https://github.com/PAF-KIET-GPGPU-Programming-Fall-2018/project-cnn_gpu/blob/master/Code%20Implementations/Cuda%20Implementation/Parallel/project/kernel.cu