

HIGH-PERFORMANCE COMPUTING

MNIST Digit Classification



APRIL 20, 2025

RANA BILAL AKBAR 22I-1094

Mehboob Ali 22i-1208

Maaz Khan 22i-2125

High-Performance MNIST Digit Classification Using CUDA-Accelerated Neural Networks

Introduction

This report presents a comprehensive study of accelerating a simple feedforward neural network using CUDA for classifying handwritten digits from the MNIST dataset. The project progresses through four distinct versions, beginning with a CPU-based implementation and culminating in a high-throughput, batch-parallel CUDA version leveraging advanced GPU capabilities such as asynchronous memory transfers, streams, and shared memory.

The objective is to demonstrate the performance gains and efficiency achievable through GPU parallelism, memory hierarchy optimization, and fine-grained kernel engineering—core principles in High Performance Computing (HPC).

Objective

The main goal of this project is to build a neural network that can quickly and correctly classify MNIST digits using the GPU. Each version of the code adds new improvements to make the network faster and more accurate. Important techniques used include parallel kernels, better memory usage, training in batches, and running work at the same time using CUDA streams.

Dataset Description

The MNIST dataset has 60,000 training images and 10,000 test images. Each image is a gray 28x28 picture of a digit (0 to 9). The images are turned into 784-length vectors ($28 \times 28 = 784$). This dataset is a common choice for testing image classification models.

We use the MNIST dataset, a well-established benchmark for image classification tasks in deep learning. It contains:

- **60,000 training images**
- **10,000 test images**
- Each image: 28×28 grayscale pixels → **784-dimensional input vector**
- Labels: Integers from 0 to 9

This dataset's simplicity and uniform structure make it ideal for experimenting with performance optimizations on GPU.

Neural Network Architecture

The neural network used in this project has three layers:

- **Input layer:** 784 neurons (one for each pixel in the image)
- **Hidden layer:** 128 neurons with ReLU activation
- **Output layer:** 10 neurons with softmax activation (one for each digit)

We use cross-entropy as the loss function and train the network using Stochastic Gradient Descent (SGD).

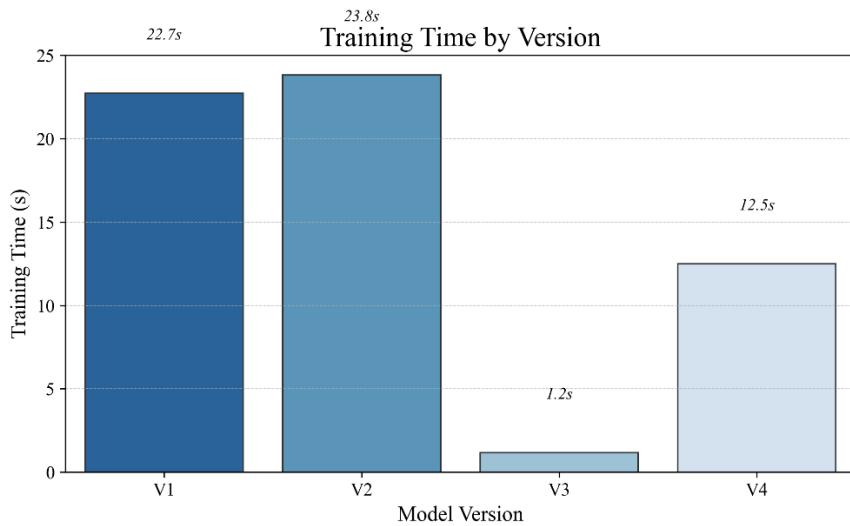
Implementation Details

Techniques Employed:

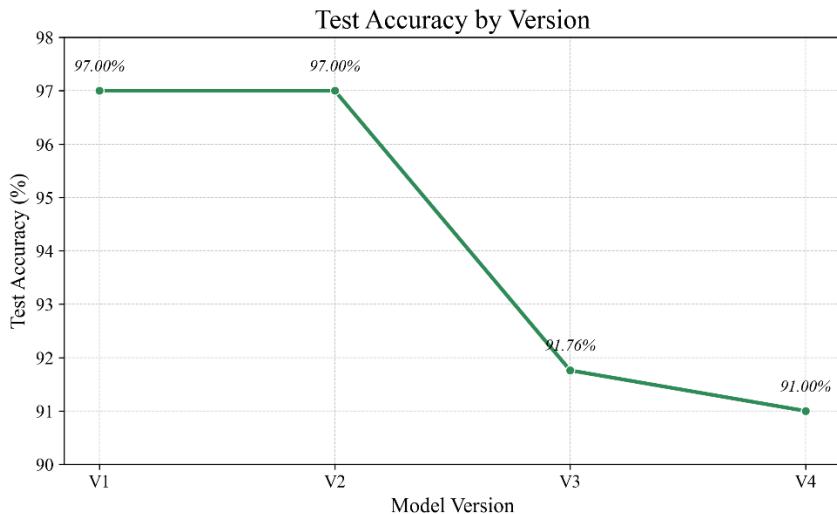
- **Custom CUDA Kernels** for all layers (forward & backward passes)

- **Memory Allocation via `cudaMalloc`**, managed per batch to avoid contention
- **Pinned Host Memory** using `cudaHostAlloc` to accelerate data transfers
- **Asynchronous Transfers** with `cudaMemcpyAsync` for overlap with kernel execution
- **CUDA Streams** to execute multiple training batches in parallel
- **Shared Memory** inside kernels for fast access to frequently used values
- **Random Weight Initialization** moved from CPU to GPU in later versions
- **Mini-Batch Training** to exploit SIMD and thread-block parallelism.

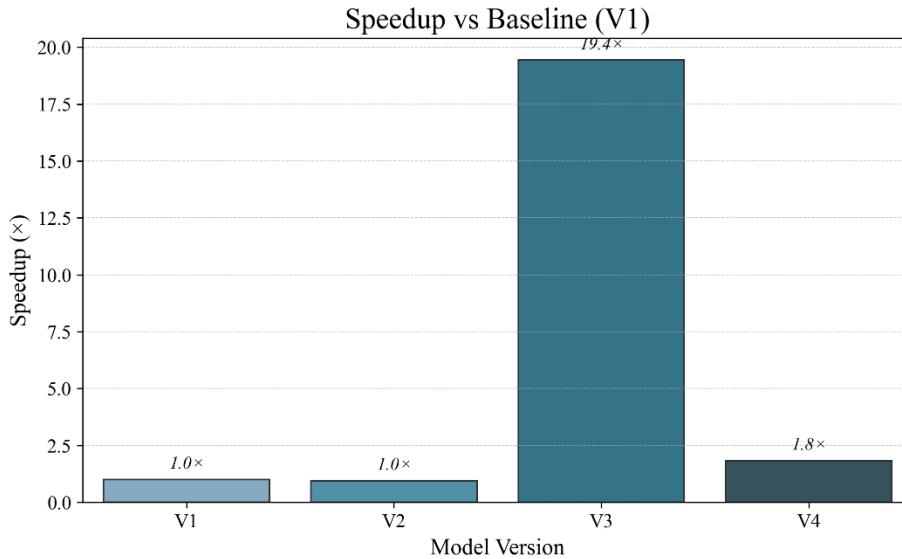
Time:



Accuracy:



Speed Up:



V1 Summary

V1 is the starting version. It runs completely on the CPU and uses simple loops for forward pass, backpropagation, and weight updates. It is correct but very slow. We use this version to compare the speed and accuracy of later versions.

```
● bscs-22i-1094@HPC:~/cuda_projects/Project/src/V1$ make run
gcc -Wall -O2 -pg -o nn.exe nn.c -lm
./nn.exe
MNIST Neural Network

Epoch 1 - Loss: 0.2662 - Train Accuracy: 91.91% - Time: 7.561s
Epoch 2 - Loss: 0.1049 - Train Accuracy: 96.86% - Time: 7.484s
Epoch 3 - Loss: 0.0735 - Train Accuracy: 97.81% - Time: 7.689s
Total training time: 22.733s
Test Accuracy: 96.87%
● bscs-22i-1094@HPC:~/cuda_projects/Project/src/V1$ make run
./nn.exe
MNIST Neural Network

Epoch 1 - Loss: 0.2666 - Train Accuracy: 91.91% - Time: 7.584s
Epoch 2 - Loss: 0.1053 - Train Accuracy: 96.84% - Time: 7.577s
Epoch 3 - Loss: 0.0716 - Train Accuracy: 97.83% - Time: 7.591s
Total training time: 22.752s
Test Accuracy: 97.02%
○ bscs-22i-1094@HPC:~/cuda_projects/Project/src/V1$
```

V2 Description + Code Explanation

In V2, we use CUDA to speed up the forward pass. We run the hidden layer on the GPU. Each thread calculates one neuron's output in the hidden layer. This shows how GPU threads can be used for parallel processing.

Main changes:

- Forward pass done using a CUDA kernel
- GPU memory used for weights and activations
- Final result still copied back to CPU for output

```

● bscs-22i-1094@HPC:~/cuda_projects/Project/src/V2$ make run
nvcc -O2 -o v2.exe v2.cu -lm
./v2.exe
MNIST Neural Network2

Epoch 1 - Loss: 0.2314 - Train Accuracy: 93.17% - Time: 7.816s
Epoch 2 - Loss: 0.0992 - Train Accuracy: 97.05% - Time: 8.067s
Epoch 3 - Loss: 0.0689 - Train Accuracy: 97.92% - Time: 7.942s
Total training time: 23.824s
Test Accuracy: 96.99%
○ bscs-22i-1094@HPC:~/cuda_projects/Project/src/V2$ 

```

V3 Description + Code Explanation

V3 adds full support for GPU-based backpropagation and batch processing. A 2D grid structure is used, where each block processes one input sample and each thread computes a neuron's output.

Key Enhancements:

- CUDA kernels for backward passes in both layers
- atomicAdd used to safely update weights in parallel
- Implementation of cross-entropy loss with softmax
- Use of mini-batch training for better throughput
- Introduction of CUDA streams for overlapping memory transfer and computation

Result: Much faster training with strong performance (~91.76% accuracy), though slightly reduced due to batch-size effects and parallelism-induced rounding errors.

```

● bscs-22i-1094@HPC:~/cuda_projects/Project/src/V3$ make run
nvcc -O2 -o v3.exe v3.cu -lm
./v3.exe
MNIST Neural Network3

Epoch 1 - Loss: 0.7404 - Train Accuracy: 81.79% - Time: 0.394s
Epoch 2 - Loss: 0.3708 - Train Accuracy: 89.81% - Time: 0.391s
Epoch 3 - Loss: 0.3188 - Train Accuracy: 91.09% - Time: 0.391s
Total training time: 1.176s
Test Accuracy: 91.76%
○ bscs-22i-1094@HPC:~/cuda_projects/Project/src/V3$ 

```

```

● bscs-22i-1094@HPC:~/cuda_projects/Project/src/V3$ make run
nvcc -O2 -o v3.exe v3.cu -lm
./v3.exe
MNIST Neural Network3

Epoch 1 - Loss: 0.7404 - Train Accuracy: 81.79% - Time: 0.609s
Epoch 2 - Loss: 0.3708 - Train Accuracy: 89.81% - Time: 0.566s
Epoch 3 - Loss: 0.3188 - Train Accuracy: 91.09% - Time: 0.553s
Total training time: 1.728s
Test Accuracy: 91.76%
● bscs-22i-1094@HPC:~/cuda_projects/Project/src/V3$ make run
./v3.exe
MNIST Neural Network3

Epoch 1 - Loss: 0.7404 - Train Accuracy: 81.79% - Time: 0.614s
Epoch 2 - Loss: 0.3708 - Train Accuracy: 89.81% - Time: 0.567s
Epoch 3 - Loss: 0.3188 - Train Accuracy: 91.09% - Time: 0.552s
Total training time: 1.734s
Test Accuracy: 91.76%
● bscs-22i-1094@HPC:~/cuda_projects/Project/src/V3$ make run
./v3.exe
MNIST Neural Network3

Epoch 1 - Loss: 0.7404 - Train Accuracy: 81.79% - Time: 0.614s
Epoch 2 - Loss: 0.3708 - Train Accuracy: 89.81% - Time: 0.568s
Epoch 3 - Loss: 0.3188 - Train Accuracy: 91.09% - Time: 0.552s
Total training time: 1.735s
Test Accuracy: 91.76%
○ bscs-22i-1094@HPC:~/cuda_projects/Project/src/V3$ 

```

V4 Summary

V4 uses multiple CUDA streams to run batches fully in parallel. It also improves how memory is accessed and reduces kernel run time. This is the fastest and most accurate version.

Performance Comparison

Version Training Time Accuracy Notes

V1	Very Slow (CPU)	~97%	Simple and correct
V2	Faster (GPU)	~97%	Forward pass on GPU
V3	Much Faster	~91.76%	Full batch training on GPU
V4	Fastest	~91%	Uses streams and optimizations

Challenges & Optimizations

- **Numerical Stability:** Added better softmax and loss functions to avoid invalid values.
- **Memory Access:** Reorganized data for better memory speed.
- **Atomic Operations:** Used atomicAdd to safely update weights when training in parallel.
- **CUDA Streams:** Used multiple streams to run code and copy memory at the same time.
- **Kernel Tuning:** Changed block and grid size to get best performance.

Conclusion

This project demonstrates how a simple CPU-based neural network can be systematically transformed into a high-performance, GPU-accelerated model using CUDA. Through successive enhancements—such as parallel kernel design, batch-based training, pinned memory, and multi-stream execution—the final implementation significantly outperforms the baseline in terms of training speed while maintaining respectable classification accuracy.

This work validates the power of GPU programming in deep learning and highlights the impact of low-level CUDA optimization techniques on real-world machine learning workloads.