

# GPU Programming Assignment 1 - Report

---

**Course:** GPU Programming

**Assignment:** Train a 3-layer Neural Network (NN) for certain accuracy suitable for any classification problem

## Team Members

Name	Roll Number
Akula Rajesh	M25AI1048
Pradeep Annepu	M25AI1109
Anirudh Reddy Aligireddy	M25AI1131

---

## Train a 3-layer Neural Network (NN) for certain accuracy suitable for any classification problem

---

### 1. Introduction to 3-Layer Neural Networks and Classification

A 3-layer Neural Network typically comprises an input layer, one hidden layer, and an output layer. This architecture is a foundational model in deep learning, capable of learning complex non-linear relationships in data. For demonstrating and benchmarking such a network, the MNIST handwritten digit classification dataset is an excellent choice. MNIST consists of 28×28 grayscale images of handwritten digits (0-9), representing a multi-class classification task. Its simplicity and well-defined nature make it ideal for illustrating core neural network concepts and performance characteristics.

### 2. Neural Network Architecture

The 3-layer Neural Network designed for MNIST classification has the following structure:

- **Input Layer:** Each MNIST image is 28×28 pixels. These are flattened into a 1D vector of 784 features. Thus, the input layer consists of 784 neurons.
- **Hidden Layer:** This intermediate layer processes the input features. For this model, a dense layer with 128 neurons is used. The Rectified Linear Unit (ReLU) activation function is applied to the output of this layer, given by  $f(x) = \max(0, x)$ . ReLU is chosen for its computational efficiency and its ability to mitigate the vanishing gradient problem.
- **Output Layer:** This layer produces the classification probabilities for each digit. With 10 possible digits (0-9), the output layer has 10 neurons. The softmax activation function is given by:

$$\text{softmax}(z)_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}, \quad j = 1, \dots, K$$

It converts the raw outputs into a probability distribution over the 10 classes, where each probability is between 0 and 1, and their sum is 1.

The network is trained using the Adam optimizer, a popular adaptive learning rate optimization algorithm. The sparse categorical cross-entropy loss function is employed, suitable for multi-class classification with integer labels. Model performance is evaluated using accuracy as the primary metric.

### 3. Implementation Approaches and Speed-Up Analysis

To understand performance, two implementations were considered:

#### 3.1 Optimized Framework Implementation (Keras/TensorFlow on CPU)

This represents a highly optimized "standard" approach, utilizing the TensorFlow backend which leverages efficient C/C++ libraries and BLAS (Basic Linear Algebra Subprograms) for numerical operations.

#### 3.2 From-Scratch NumPy Implementation

A manual implementation of the forward and backward propagation steps using only the NumPy library for array operations. This provides a baseline to observe the performance benefits of optimized frameworks.

#### Benchmarking Results (Illustrative)

**Note:** Due to limitations in the current execution environment, direct live benchmarking was not possible. The following figures represent typical performance differences observed in such comparisons.

Implementation Method	Training Time (5 Epochs, Batch Size 64)	Test Accuracy (Illustrative)
Keras/TensorFlow (CPU)	Approximately 10-15 seconds	~97-98%
From-Scratch NumPy	Approximately 150-250 seconds	~90-95%

#### Speed-Up Analysis

The Keras/TensorFlow implementation typically achieves a significant speed-up over the from-scratch NumPy implementation.

**Illustrative Speed-Up:** Approximately 15-20x

This substantial speed advantage of Keras/TensorFlow (or other deep learning frameworks like PyTorch) over a pure Python/NumPy implementation is primarily due to several factors:

- **Optimized Backends:** Frameworks are built on highly optimized C/C++ (and often CUDA for GPUs) engines, which are significantly faster than Python for numerical computations.
- **BLAS Libraries:** They link to highly optimized BLAS libraries (e.g., Intel MKL, OpenBLAS) for linear algebra operations like matrix multiplication, which are heavily used in neural networks. These libraries are often hand-tuned for specific CPU architectures and utilize techniques like vectorization and multi-threading.
- **Graph Optimization:** TensorFlow constructs a computational graph that can be optimized for efficiency before execution, including memory management and operation fusion.
- **Reduced Python Overhead:** By performing large chunks of computation in compiled code, frameworks minimize the overhead of Python's interpreter.

## 4. Conclusion

Training a 3-layer Neural Network for classification, such as on the MNIST dataset, highlights significant performance differences between highly optimized deep learning frameworks (like Keras/TensorFlow) and custom NumPy implementations. The frameworks provide substantial speed-ups (typically 15-20x or more) by leveraging optimized C/C++ backends and BLAS libraries, minimizing Python overhead.

Regarding specific optimizations:

- **Tiling** is a crucial technique for optimizing matrix multiplication by improving cache locality. While essential for performance, it is transparently handled by the underlying BLAS libraries used by frameworks, making manual implementation in Python generally less efficient than using `np.dot`.
- **Neighbor-based computations** (stencil operations) are not directly relevant to the fully connected layers of a 3-layer MLP. However, they are a core concept in Convolutional Neural Networks, where they enable efficient local feature extraction in grid-like data such as images.