

---

# Parallel Backpropagation for Multilayer Neural Networks

---

**Nitin Kamra**

Department of Computer Science, University of Southern California

NKAMRA@USC.EDU

**Palash Goyal**

Department of Computer Science, University of Southern California

PALASHGO@USC.EDU

**Sungyong Seo**

Department of Computer Science, University of Southern California

SUNGYONS@USC.EDU

**Vasileios Zois**

Department of Computer Science, University of Southern California

VZOIS@USC.EDU

## Abstract

We present a parallel approach to classification using neural networks as the hypothesis class. Neural networks can have millions of parameters and learning the optimum value of all parameters from huge datasets can be very time consuming, if implemented on a single processor. In this work, we have implemented parallel backpropagation to train Multi Layer Perceptrons (MLPs) for classification tasks. Specifically, we implement a serial backpropagation algorithm and its parallel counterpart with Pthread library in C++. We also compare the two implementations with CUDA implementation of backpropagation on a GPU and study the parallelization speedup obtained for various network architectures and increasing problem sizes. We perform our tests on two benchmark datasets: MNIST and KDD Cup 1999, and finally compare all our implementations with the corresponding implementations in the state-of-the-art deep learning library: Theano.

## 1. Introduction

Artificial neural networks are powerful machine learning tools used in many applications including but not limited to search engines, spam and fraud detection, image classification, diagnostic medicine applications and stock market prediction.

Prior to application, neural networks undergo a training phase which is known to be very computationally inten-

sive. This is primarily because the prevailing neural network architectures are implemented using several hidden layers, with each one consisting of thousands to millions of neurons in order to generalize well on diverse inputs. In this case, the resulting number of parameters that need to be trained are in the order of millions.

Furthermore, achieving high accuracy requires considering a large number of training examples (usually in the order of millions). For this reason training a neural network is both a data and resource intensive operation. This calls for efforts to parallelize the training process on multi-core and while emphasizing at utilizing the maximum available bandwidth.

Currently Minibatch Gradient Descent (henceforth called MGD) is the most commonly used optimization algorithm used to train neural networks in supervised settings. It is implemented in a layerwise-recursive fashion which is termed *backpropagation* in the context of neural networks.

In this paper, we have implemented parallel backpropagation to train Multi Layer Perceptrons (MLPs) for classification tasks. The rest of the paper is organized as follows: section 2 presents a description of supervised learning tasks and neural networks as classifiers. Section 3 describes the conventional serial backpropagation algorithm, and an approach to parallelization using multiple threads. Section 4 describes how to implement backpropagation on a GPU with cuda. Section 5 describes our datasets and the experiments we performed on them. Section 6 presents our results obtained for these implementations and analyzes the speedup obtained for various network architectures and increasing problem sizes. It also presents a comparison with the same algorithms implemented using a state-of-the-art neural network library Theano. Finally, we conclude in section 7 with a discussion of potential applications and future work for this project.

## 2. Problem Description

We first formally describe a classification task in the supervised learning setting. Then we describe a feedforward neural network with fully connected layers. Feedforward neural networks act as function approximators in such tasks.

### 2.1. Classification Problem

More formally, given a dataset  $\mathcal{D} = \{x_i, y_i\}_{i=1:N}$  with data points  $x_i \in \mathbb{R}^D$  and labels  $y_i \in \mathbb{R}^P$ , the classification task involves making an accurate label prediction  $\hat{y}$  on a previously unseen data point  $x$ . We approximate the label as a function of the datapoint using a classifier (a feedforward neural network here) with parameters  $\theta = \{\theta_k\}_{k=1:K}$  as follows:  $y \approx \hat{y} = f(x; \theta)$ . The classifier (neural network) learns the function  $f$  from the training data  $\mathcal{D}$  by tuning its parameters ( $\theta$ ) to minimize a pre-specified loss function, for instance, the Mean-Squared Error loss:

$$\mathcal{L}_{MSE}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i; \theta))^2 \quad (1)$$

This minimization can be carried out using optimization algorithms like Gradient Descent, Newton's method, Levenberg-Marquardt algorithm etc. Though Newton's method is a second-order optimization technique, it requires the computation of the hessian of the objective function which is very prohibitive for a large number of parameters like in a neural network. Levenberg-Marquardt algorithm also requires computing matrices of the size of hessian and can be very slow for classifiers with a large number of parameters. Currently Gradient Descent is the most successful technique to train huge neural networks with millions of parameters, since it provides a decent tradeoff between convergence speed and memory requirements. We will describe gradient descent in section 3.

### 2.2. Feedforward Neural Networks

The basic unit of feedforward neural networks is a neuron. A single neuron generally takes a vector of inputs  $x \in \mathbb{R}^n$  and outputs a single scalar  $a \in \mathbb{R}$  (called its activation). Generally the function computed by a neuron comprises of linear transformation on the input vector followed by a pointwise non-linear activation function:

$$a = f(w^T x + b) \quad (2)$$

where  $w \in \mathbb{R}^n$  is called the weight vector and  $b \in \mathbb{R}$  is the scalar bias of the neuron. Many different kinds of activation functions have been studied and used according to the task at hand e.g. linear, sigmoid, tanh, ReLU etc. We will use the sigmoid and linear( $f(z) = z$ ) activation functions in our implementation.

The sigmoid function is shown in figure 1 and is defined as

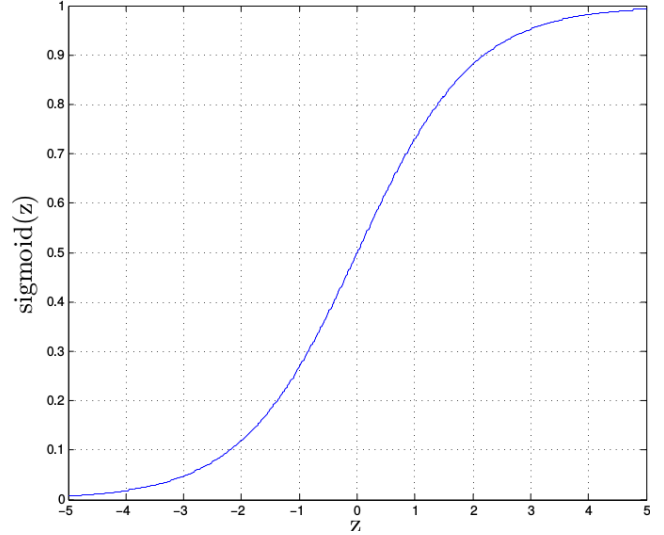


Figure 1. The sigmoid function

follows:

$$f(z) = \frac{1}{1 + e^{-z}} \quad (3)$$

A feedforward neural network is a directed acyclic graph  $G = (V, E)$  each of whose vertices  $v \in V$  is a neuron and every edge  $(u, v) \in E$  represents the output of neuron  $u$  going as an input to neuron  $v$ . In general to have a more concrete structure the graph  $G$  is organized into a layered structure and we will only use layered feedforward neural networks in our implementations.

A feedforward network with  $L$  layers comprises of a single input layer,  $L - 2$  hidden layers (their outputs are not directly observed) and a final output layer. Let the input of the neural network be  $x \in \mathbb{R}^{n_{in}}$  and the output be  $y \in \mathbb{R}^{n_{out}}$ . The  $l^{th}$  layer has  $n_l$  neurons and the full network has  $n = \sum_{l=1}^L n_l$  neurons.

The input to the first layer (denoted  $x_1$ ) is  $x$  and the output (activation) of the first layer is denoted  $a_1 \in \mathbb{R}^{n_{in}}$ . The input layer contains  $n_1 = n_{in}$  dummy neurons each of which takes a single scalar input  $n_{in}$  and passes it as it is i.e.  $a_1 = x_1 = x$ .

All subsequent layers  $l \in \{2, 3, \dots, L\}$  give an output  $a_l \in \mathbb{R}^{n_l}$  and take an input  $x_l = a_{l-1} \in \mathbb{R}^{n_{l-1}}$ . The final network output  $y$  is given by  $a_L \in \mathbb{R}^{n_L}$  ( $n_L = n_{out}$ ). Each of these layers applies a linear transformation to its input followed by a pointwise non-linear activation function i.e. for each layer:

$$z_l = (W^l)^T x_l + b_l \quad (4)$$

$$a_l = f(z_l) \quad (5)$$

where  $W^l \in \mathbb{R}^{n_{l-1} \times n_l}$  is the weight matrix of the layer and  $b_l \in \mathbb{R}^{n_l}$  is the bias vector. The function  $f(\cdot)$  is the pointwise activation function which applies to each component of the  $z_l$  vector and we will use sigmoid as our activation function for the classification task.

---

**Algorithm 1** Gradient Descent
 

---

**Input:** data  $x_i$ , size  $m$   
**repeat**  
     Initialize  $noChange = true$ .  
     **for**  $i = 1$  **to**  $m - 1$  **do**  
         **if**  $x_i > x_{i+1}$  **then**  
             Swap  $x_i$  and  $x_{i+1}$   
              $noChange = false$   
         **end if**  
     **end for**  
**until**  $noChange$  is  $true$

---

### 3. Serial and Parallel Backpropagation

As explained in section 2.1. classifiers try to choose their parameters ( $\theta$ ) in order to minimize some pre-specified loss function. In this work, we will work with the Mean-squared Error loss function as defined in equation 1.

#### 3.1. Gradient Descent

To optimize the loss function, we will use the Gradient Descent algorithm which in its most basic form takes the derivative of the cost function w.r.t. all the parameter values and updates the parameters by a value proportional to this gradient and opposite in sign. In its most basic form the Gradient Descent algorithm takes the form shown in algorithm 1.

1. Initialize all weights ( $\theta$ ) randomly with small values close to 0.
2. Repeat until convergence {

$$\theta_k := \theta_k - \alpha \frac{\partial \mathcal{L}_{MSE}}{\partial \theta_k} \quad \forall k \in \{1, 2, \dots, K\}$$

}

Note that the derivative of the loss function generally takes the following form:

$$\frac{\partial \mathcal{L}_{MSE}}{\partial \theta_k} = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i; \theta)) \frac{\partial f(x_i; \theta)}{\partial \theta_k}$$

Since this requires a summation over all the training examples, the gradient computation is the biggest bottleneck for many supervised learning tasks on huge data sets.

On the other hand, having the gradient as a sum of partial gradients with respect to individual training examples opens up the possibility of parallelizing the gradient computation efficiently. This can be done by distributing training examples across various processors/machines and letting them compute a partial gradient over their own training examples and then summing up these partial gradients to get the actual gradient. A good approximation of this process can be

obtained by using a mini-batch of training examples per iteration, instead of using the full dataset.

Describe backpropagation serial implementation in detail

Describe platform used

Describe backpropagation parallelization using pthreads

### 4. Implementing BackPropagation on GPU

Graphics Processing Units (GPUs) are massively parallel processors that were designed for efficient computer graphics rendering and image processing. In fact, GPUs are very effective when used to execute the graphics rendering pipeline which is a sequence of geometric transformations on small multidimensional vectors. For this reason, GPUs are suitable for executing very fast certain linear algebra operations including but not limited to vector-vector addition, matrix-vector and matrix-matrix multiplication. Moreover, because GPUs consist of many throughput oriented multi-processors that follow the Single Instruction Multiple Data (SIMD) execution model, they can be very useful for applications that require processing and transformation operations on large dataset.

Training neural networks is a process that can be realized through a series of matrix-matrix multiplications and additions which are applied for a large number of training examples. Every neural network can be defined and operated on by following this abstraction. Back propagation can be implemented as a series of kernel executions that are combined to ultimately compute the change in the weight values caused by the corresponding training batch. In our implementation, a neural network is viewed as a collection of layers, each one consisting of 3 distinct matrices. These matrices store the outgoing weight values  $W_i$ , the incoming activation values  $A_i$  and the delta error values  $D_i$  computed by back propagation. The activation values of the input layer (i.e.  $A_0$ ) are the actual training examples in the corresponding batch. The activation values and the delta error values are stored in column-major order. In contrast the weights for each layer are stored in row-major order. The bias values are embedded in the last column of the weight matrix.

We decompose back-propagation into 5 steps each one implemented by distinct kernel. The first step, handles the feed forward activation by implementing a tiled matrix-matrix multiplication based on Eq. 6. Here we denote with  $f$  the preferred activation function which is usually the sigmoid function. Our implementation is designed around templates and supports user defined activation functions. The feed-forward step is implemented using the tiled matrix-matrix multiplication as it is described in (). All threads participate in loading the tiles of the input matrices into shared memory. The threads use registers to accumulate partial results. After parsing the complete matrix each thread applies the activation function on the resulting value before storing it

into global memory. All read and write operations to global memory are coalesced and there are no bank conflicts because thread iterate over the secondary matrix dimension in each case. An additional step is required to initialize the thread registers with the bias values.

$$A_{i+1} = f(W_i \cdot A_i + b_i), \forall i \in [1, L-1] \quad (6)$$

$$\delta_{output\_kernel} D_L = Y - A_L \quad (7)$$

$$D_i = W_i^T \cdot D_i \circ d(W_{i-1} \cdot A_{i+1}) \quad (8)$$

$$W_i = W_i + \frac{n}{b} \cdot \sum_{j=1}^b D_{i+1}^j \cdot (A_i^j)^T \quad (9)$$

Following this previous step, the delta error values are computed for each layer using Eq. 7. This equation is computed using 3 kernels, one for computing the output layer delta and two more that compute the transpose matrix-matrix multiplication and the hadamard product of the derivative for the corresponding activation function. A variation of the tiled matrix-matrix multiplication kernel is used to perform these operations. For the first kernel, the indexing is changed to enable multiplication with the transpose of the weight matrix. Here accesses to global memory remain coalesced since we order the thread blocks vertically on the weight matrix. However, we incur few bank conflicts when accessing the data in column-major order from shared-memory. In the second case, we chose to re-compute the product of  $W_i \cdot A_i$  and multiply it with the result from the previous operation. This is because our goal was to support arbitrary defined activation functions. In case the sigmoid function is used, the derivative can be replaced with  $A_i \cdot (1.0 - A_i)$  which enables significant reduction of the required computation.

Finally, we use the computed activation and the next layer delta matrices (Eq. 8) to update the weights of the corresponding layer depending on the chosen learning rate and batch size. This operation can be considered as another variation of matrix-matrix multiplication where each column from  $D_{i+1}$  and  $A_i$  are multiplied to produce a single matrix. The number of resulting matrices are equal to the batch size and are accumulated in thread registers. The resulting summation is multiplied by  $\frac{n}{b}$  and added to the weight values of the corresponding layer. For this operation, access to global memory is again coalesced. However, shared memory access incurs many bank conflicts which are proportional to the batch size. In order to improve the performance, the corresponding tiles are loaded transposed into shared memory. This incurs only a single bank conflict during loading and avoids many bank conflicts during the multiplication phase.

## 5. Experiments

In this section, we first describe the datasets used for our experiments and the corresponding classification tasks. We then delineate the networks used for each of these datasets and the procedure to determine the values of hyperparameters. This is followed by a description of the evaluation metric of the speedup.

### 5.1. Dataset description

We test the speedups achieved by the parallel implementation on two public datasets - MNIST and KDD Cup 1999. These datasets have been widely used by the machine learning community to compare different learning algorithms. Their description is as follows:

- **MNIST** - This is an image database of handwritten digits which is commonly used for training various image processing systems. The task here is to classify the digit in the image. This dataset was derived from NIST's datasets and was formed by mixing the samples from NIST. This mixing was done because the training and testing datasets in the original NIST dataset were obtained from two different groups of people (Census Bureau employees and high school students). This dataset consists of 60,000 training images and 10,000 testing images, each of size  $28 \times 28$ . In our experiments, we further divided the training set into training and validation in 5:1 ratio in order to determine suitable values for the hyperparameters.
- **KDD Cup 1999** - This is a dataset which first appeared in the KDD Cup held in 1999. The dataset has about 4 million data points each of which has 42 features. The features are of two types - categorical and integer. The task in the competition was to build a network intrusion detector which can be used to distinguish between good and bad network connections.

### 5.2. Network details

We now describe the details of the structure of neural networks used. For our experiments, we use two different configurations of neural network -

- **Net-1h** - This network consists of 3 layers - input layer, one hidden layer and output layer. The input layer has 784 neurons for MNIST dataset and 42 neurons for KDD Cup 1999. The hidden layer is composed of 1024 neurons and the output layer has 10 neurons when used for MNIST dataset and 1 output for KDD Cup 1999.
- **Net-2h** - This network is composed of 4 layers - input layer, two hidden layers and output layer. The input and output layer are same as Net-1h and each of the hidden layers has 1024 neurons.

### 5.3. Hyperparameter selection

Hyperparameter selection, also known as model selection, is an important problem in machine learning. It refers to choosing the optimal hyperparameters for a learning algorithm. It is crucial to select the values for hyperparameters which generalize well. This enables the algorithm to perform well on test data.

In the context of neural networks, hyperparameters refer to learning rate and regularization constant. Since our implementation is free of regularization, we only have one hyperparameter - learning rate. To choose an optimal learning rate, we used the technique of grid search. As mentioned above, we divide the training set into training and validation. We use different learning rates defined in a grid and test the accuracy on the validation data. The rate which gives the highest accuracy on validation is chosen and used to evaluate the performance of the algorithm on test data.

### 5.4. Evaluation methodology

Describe datasets and classification task

Describe experimental network sizes

Describe hyperparameter selection

Describe experiments: serial vs pthreads vs cuda vs theano.

- Serial implementation of Mini-batch Gradient Descent(MGD) in C++
- Parallel implementation of MGD in C++ using Pthreads
- Comparison between serial MGD, parallel MGD with OpenMP and Theano's gradient descent (CPU)
- Parallel MGD using Cuda-C on a GPU
- Comparison between serial MGD, parallel MGD (Cuda) and Theano's gradient descent (GPU)

To test the above implementations, we will use the following datasets:

- **MNIST** - This is a database of handwritten digits that is commonly used for training various image processing systems. It has 60,000 images each of size  $28 \times 28$  pixels. The task is to classify the digit in the image.
- **KDD Cup 1999** - This is a data set designed for intrusion detection (distinguishing bad network connections from good ones). It has about 4 million data points each of which has 42 features.

Note that since Theano is a state of the art Deep Learning library which performs many other optimizations apart from just parallelizing Neural Network training, we do not expect our code to run faster than Theano. Instead the comparison will be done to see how close we can get to the current state

of the art by simply parallelizing the Neural Network training.

## 6. Results and Analysis

Provide experiment results: serial

Compare serial vs pthreads vs cuda vs theano.

Provide graphs for speedup with increasing problem size

Analyze speedup with number of threads/cores used

Analyze the level of parallelization obtained

Analyze performance bottlenecks

Do overall analysis and compare obtained results with your expectations

## 7. Conclusion and Future Work

In this paper, we have explored a technique to parallelize training of multilayer feedforward neural networks by splitting the training examples among multiple threads of execution. We have obtained significant performance gains by parallelizing the training process both by using multiple cores and also on a GPU. We have demonstrated a -x times speedup using multiple cores only and then a -x times speedup using a GPU without any significant loss of accuracy. Our parallelized backpropagation implementation can find good usage for many real-time processing tasks. Typical examples include but are not limited to:

- Real-time online control of robotic manipulators which learn online from incoming sensory data.
- Self-driving cars which need to do online object recognition in real-time using deep neural networks.
- Speech translators which deploy deep recurrent neural networks for translation in real-time.