

Hardware Accelerator for the Training of Neural Networks

James Erik Groving Meade
s173476

DTU



Kongens Lyngby 2019

Technical University of Denmark
Department of Applied Mathematics and Computer Science
Richard Petersens Plads, building 324,
2800 Kongens Lyngby, Denmark
Phone +45 4525 3031
compute@compute.dtu.dk
www.compute.dtu.dk

Abstract

This thesis proposes a novel hardware architecture to accelerate the training of neural networks with small batch sizes. The accelerator uses a modular, parameterizable and computationally well-balanced design to successfully implement high-performance online training of neural networks. By using parallelism at the neuron level, the accelerator was able to achieve a speedup of 17.3 against a PyTorch CPU implementation of a specific neural network architecture. The accelerator also performs nearly as fast as a PyTorch GPU implementation of the network that used a batch size of 50 during training.

This thesis also highlights the importance of high-precision calculation for training. The highest accuracy attained by the accelerator on the MNIST dataset was 85.845%, which is a result of 18-bit fixed point precision being unable to successfully converge to a local optima due to the accumulation of precision error causing the degradation of training accuracy after a few epochs.

Preface

This thesis was prepared at DTU Compute in fulfillment of the requirements for acquiring an M.Sc. in Computer Science and Engineering.

This thesis deals with the design of a hardware accelerator for the training of neural networks. Low-level design is one of my greatest passions and thus it has been an absolute privilege to have been given the opportunity to combine hardware with the surging field of machine learning.

Lyngby, 28-June-2019

A handwritten signature in black ink, reading "Erik Meade". The script is fluid and cursive, with the first letters of the first and last names being capitalized and prominent.

James Erik Groving Meade
s173476

Acknowledgements

First and foremost, I would like to sincerely thank my advisor Jens Sparsø for his time and guidance throughout the entire duration of my thesis. Our regular meetings helped me to stay grounded and to think critically about my project.

I would also like to thank my former classmate, Cheng Fu, who is currently a PhD student at the University of California, San Diego. My conversations with him at the beginning of my foray into this thesis helped me establish my footing.

Contents

Abstract	i
Preface	iii
Acknowledgements	v
1 Introduction	1
2 Background	3
2.1 Neural Networks	3
2.1.1 The Neuron	3
2.1.2 Fully-Connected Layers	4
2.1.3 Activation Functions	5
2.1.4 Cross-Entropy Loss	7
2.1.5 Backpropagation	7
2.1.6 Hyperparameters	10
2.2 Deep-Learning Frameworks	11
2.2.1 PyTorch	11
2.3 Related Work	11
3 Software Model	13
3.1 Overview	13
3.2 Motivation	13
3.3 Design	14
3.3.1 Layers	14
3.3.2 Training	14
3.4 Source Code Structure	15
3.5 Usage	18
3.6 Testing	21

3.6.1	Test Suites	21
3.6.2	Building and Running the Test Suites	24
4	Hardware Model and Implementation	27
4.1	Hardware Setup	27
4.2	The Implemented Neural Network	28
4.3	Design Goals	28
4.4	Overall Architecture	29
4.5	Training Process	30
4.6	Computational Precision	31
4.7	Module Architecture	31
4.7.1	Fully-Connected Layers	31
4.7.2	Interlayer Architecture	39
4.7.3	Softmax Layer	40
4.8	PS – FPGA Communication	40
4.8.1	AXI Implementation for the PS	42
4.8.2	Running Programs from the PS	43
4.8.3	AXI Implementation for the FPGA	44
4.8.4	Memory Map Layout	46
4.9	PetaLinux	46
5	Hardware Model Testing and Verification	49
5.1	Simulation	49
5.1.1	Project Modifications to Simulate of the Design	50
5.1.2	Testing Environment	50
5.1.3	Simulation Output	50
5.1.4	Correctness of Simulated Outputs	51
6	Results	55
6.1	Benchmarking Models and Structure	55
6.2	Parallelism in GPU-based Training	55
6.3	Evaluation Hardware	56
6.4	Performance	56
6.4.1	Training	57
6.4.2	Inference	57
6.4.3	Active/Idle Cycles	60
6.5	Training Accuracy	60
6.5.1	Stability of Training	61
6.6	Implemented Design	62
6.6.1	Resource Usage, Power, and Timing	62

7	Analysis	65
7.1	Allocating Computational Kernels for Performance	65
7.1.1	Allocating Based on Only Forward Pass Analysis	65
7.1.2	Distribution of the Kernels	66
7.2	Cycle Analysis	68
7.3	Improving Performance	70
7.4	Granularity for Neural Network Computation	70
7.5	Ideal Learning Rate vs. Precision	71
7.6	Potential Solutions for the Lack of Precision	72
7.7	Weight Storage	73
8	Discussion	75
8.1	Overall Performance	75
8.2	Finely-Grained Parallelism	75
8.3	Limitations	76
8.3.1	Precision	76
8.3.2	Data Transfer Rate	76
8.4	Future Work	76
9	Conclusion	79
A	Stuff	81
B	Stuff2	83
	Bibliography	85

CHAPTER 1

Introduction

Neural networks have seen a surge in popularity ever since a neural network, coined AlexNet, decisively won the ImageNet Large Scale Visual Recognition Challenge in 2012, achieving 10.8% higher accuracy than the next best solution [KSH12]. It was the only neural network entry in the entire competition; and the victory stunned much of the academic world.

This moment has been generally regarded as the spark that ignited the massive surge in academic interest toward neural networks and statistical machine learning. In the 7 years since, research regarding neural networks has yet to slow down as real-world applications of neural networks continue to be found. To name a few, neural networks are currently in use for facial recognition at Facebook [TYRW14], translation for Microsoft [XHQ⁺16], spam filters for Google's Gmail [gma15] and endless more.

In order for these neural networks to have such stellar accuracy on tasks such as image classification, they must first learn from labeled data in a process known as training. The neural network training process has an incredibly high level of inherent parallelism, and thus GPUs have emerged as the device of choice to train neural networks. GPU-based training takes advantage of data-level parallelism to train networks by assigning individual inputs in a training batch to different cores, which then all perform the same computations on different data. This coarsely-grained approach to training works well for large batch

sizes. However, since these GPU models use data parallelism for speedup, GPU model training performance degrades for small batch sizes and are even slower than CPU models when it comes to online learning. Online learning is when a single labeled data sample is fed to the network at a time, or in other words, when the batch size is equal to 1 [MR17].

Today’s solutions for training neural networks with small batch sizes do not take advantage of the fine-grain parallelism available in neural networks. This thesis presents a hardware accelerator that uses fine-grain parallelism at the neuron level to achieve high training performance. The accelerator results in much faster performance compared to the current solutions available for training neural networks with small batch sizes.

This thesis proposes a novel hardware architecture for the training of neural networks. While the focus of the thesis is the architectural design of the hardware accelerator, a basic understanding of neural networks is helpful. As such, Chapter 2 reviews the basics of neural networks and surveys related work on designing hardware to optimize for neural network computation. Chapter 3 describes the software model that was implemented to verify and further understand the algorithms to be used for the hardware model. Chapter 4 covers the hardware model and implementation of the accelerator. Next, Chapter 5 documents the testing methods used to functionally verify the hardware model. Chapter 6 presents the results of the thesis and Chapter 7 provides analysis of these results. Chapter 8 discusses the project as a whole and what future work could be done to improve the project. Finally, Chapter 9 presents the conclusion of the thesis.

CHAPTER 2

Background

2.1 Neural Networks

A neural network is a machine learning tool ideal for conducting supervised learning. As a relatively recent field, the application of neural networks has rapidly extended across many domains,

2.1.1 The Neuron

The *neuron* is the basic computational unit of a neural network. A *layer* is comprised of one or more neurons. The computation performed by a neuron is shown below.

$$\text{net} = \mathbf{w} \cdot \mathbf{x} + b \tag{2.1}$$

$$y = f(\text{net}) \tag{2.2}$$

The *fan-in* to a neuron is the amount of elements in the input vector $\mathbf{x} = x_1, x_2, \dots, x_n$. For each element, there is a corresponding parameter referred to as a *weight*. The weights of a neuron form the weight vector \mathbf{w} . The neuron also has an offset b which helps with normalization. The neuron's net is first

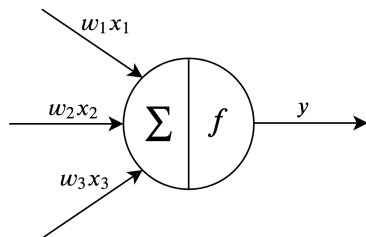


Figure 2.1: A neuron with 3 inputs; bias term omitted for simplicity.

computed as shown in equation 2.1, and then the output, or activation, is computed according to the neuron’s activation function. This is shown visually in figure 2.1.

Weight Initialization Proper weight initialization is paramount to successfully training a neural network. Firstly, weights cannot be all initialized to 0, for this will result in the same gradient for all weights, and thus all weights will be updated in the same manner. This would effectively mean that the network would become a function of a singular weight.

The most naïve way to initialize weights would to assign each weight a random value between some range. In most cases, this is good enough for the network to converge to a relatively optimal solution so long as the range is not to extreme. A recent popular and effective way to initialize the weights is through He Initialization, which randomly initializes weights using a normal distribution with a mean of 0 and a variance of $\frac{2}{f_{an_in}}$ [HZRS15].

2.1.2 Fully-Connected Layers

A fully-connected layer is a vector of neurons. All neurons in a fully-connected layer receive the same input vector. This vector is the previous layer’s output. A fully-connected layer with 3 neurons receiving input from an input layer is shown in figure 2.2. The output is a vector comprising of the outputs of each neuron. Each neuron output is calculated using the M -sized input vector as

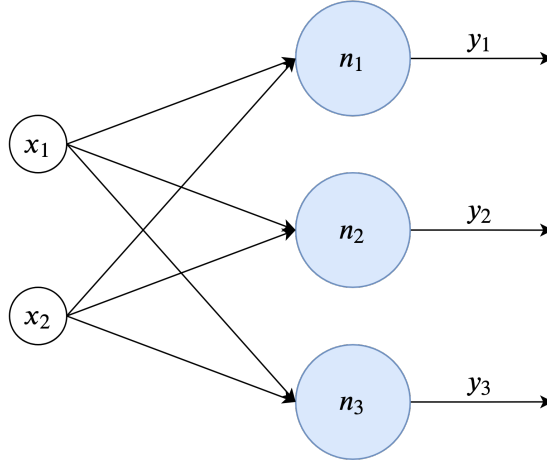


Figure 2.2: A fully-connected layer with 3 neurons, each receiving an input vector of size 2 from the input layer.

shown in equation 2.3 and added to output vector \mathbf{y} .

$$y_i = f_{\text{act}} \left(b + \sum_{j=1}^M (w_j x_j) \right) \quad (2.3)$$

$$\mathbf{y} = \{y_1, y_2, \dots, y_n\} \quad (2.4)$$

2.1.3 Activation Functions

Without activation functions, the neural network would simply devolve to a linear classifier. Activation functions provide neural networks with the non-linearity to solve complex classification problems. Two of the most common activation functions are the rectified linear unit (ReLU) and the softmax function. These are the two activation functions that were chosen to be used in the software and hardware models of this thesis.

ReLU ReLU is a powerful activation function that has found widespread use due to its mathematical simplicity. The ReLU function is shown in equation 2.5.

$$y = \max(0, x) \quad (2.5)$$

Notably, the ReLU function is much easier to compute compare to the sigmoid or hyperbolic tangent functions, which both use the exponential function. The ReLU function also quite frequently performs just as well if not better compared to other activation functions. One of the reasons is because it does not suffer as much from the vanishing gradient problem [GBB11]. The vanishing gradient problem is encountered during training using backpropagation, which uses the chain rule from calculus, briefly covered in section 2.1.5. Since gradients will always be less than 1 for most loss functions, the gradients become geometrically smaller with each layer. Since ReLU only saturates in one direction, ReLU networks will be more sparse, in the sense that many of the neurons will have an output of 0.

ReLU-based neural networks also tend to reach convergence quicker than neural networks using the sigmoid or the hyperbolic tangent functions. It also results in a sparsely activated network, in that since the neuron output is 0 if the net is negative, that many neurons in the network will have an output of 0. This is also similar to how biological neurons also follow a sparse firing model, and has shown to be effective [GBB11].

Conversely, since active neurons in ReLU network are sparse, this brings rise to another potential problem, the “Dying ReLU Problem.” This problem occurs when the sparsity increases to the point where a large majority of the neurons in the network become inactive during training and ultimately never become active again. Fortunately, this problem can be ameliorated with proper weight initialization [LSSK19].

Softmax The softmax function converts a vector of logits to a vector of probabilities. It has seen widespread use in neural networks that are used to predict the class of an input. The softmax function is shown in equation 2.6.

$$\sigma(\mathbf{x})_i = \frac{e^{x_i}}{\sum_{j=1}^C e^{x_j}} \quad (2.6)$$

In this function, x_i is the net of neuron i from the layer. Generally, the softmax function is used in the last layer to generate probabilities for multi-class problems. Each neuron in the layer represents a class, so the size of the last layer is equivalent to the number of classes, C . In much of the literature, the softmax portion of a neural network is referred to as the softmax layer as opposed to simply being the activation function of the neuron nets in the last layer.

2.1.4 Cross-Entropy Loss

Cross-entropy loss is a probabilistic loss function and as such, is frequently paired with the softmax activation function. This allows for the probabilities output from the softmax function to be used as inputs for calculating the cross-entropy loss. Cross-entropy loss is computed using probabilities and is shown in equation 2.7.

$$\mathcal{L}(\mathbf{x}) = \sum_{i=1}^N q(x_i) \log(p(x_i)) \quad (2.7)$$

In this function, $q(x_i)$ is the true probability of x belonging to class i , therefore, $q(x_i) = 1$ when x is of class i and 0 otherwise; $p(x_i)$ is equal to the predicted probability.

2.1.5 Backpropagation

Backpropagation is a method in which the weights of a network can be trained on a dataset by propagating the loss (also referred to as gradient in gradient descent) from the output layer backward through the network. There are three computational steps to be made during backpropagation: propagating loss gradients to the previous layer, using loss gradients for neurons in a layer to calculate individual weight gradients, and then finally to update the weights.

Calculating the Loss Gradients in the Output Layer For the first part of backpropagation, we must use the partial derivative of the loss function with respect to each of the neuron outputs to begin backpropagation. Note that the cross-entropy loss is calculated directly from the probabilities from the softmax function of the last layer. Therefore, the loss must derive the loss function with respect to the probabilities, and then must derive the softmax function in order to attain $\frac{\delta \mathcal{L}}{\delta \text{net}_o}$ for the neurons in the last layer. The calculus is omitted for brevity, but the final result is clean and simple, as shown in equation 2.8 [sm-].

$$\frac{\delta \mathcal{L}}{\delta \text{net}_{o,i}} = p_i - y_i \quad (2.8)$$

This equation calculates the partial derivative of the loss with respect to the net of the last layers output neuron. p_i is the probability computed from the softmax function and y_i is the true probability. Thus, if the input sample belong to class i , y_i is equal to 1, otherwise y_i is 0. Once the initial gradient for each neuron in the last layer has been calculated, backpropagation of the loss through the previous layers is possible.

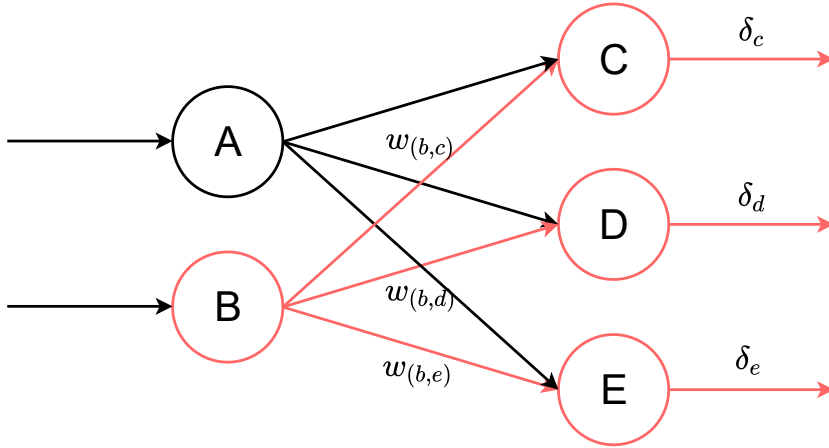


Figure 2.3: Example of backpropagating the loss gradient to the previous layer. Values used in backpropagating the loss to neuron B shown in red.

Backpropagating the Loss Gradient The strength of backpropagation is being able to use the chain rule to calculate gradients from previous layers. At a high-level, a neuron in a previous layer's output will affect the nets of neurons in the next layer. Since each activation is multiplied by a weight, the affect on the net is determined by a weight. For example, if a neuron's activation a_o increases by ϵ , then each of the next layer's neuron nets will increase by $w \times \epsilon$, where w is the weight for that connection. This connection is also sometimes referred to as a *synapse*, a term inspired from neuroscience.

An example illustrating this is shown in figure 2.3. The gradients for the nets of C , D , and E are represented by δ . The gradient of a net is commonly referred to as the *sensitivity* of a neuron. Subsequently, the weights on the synapses are also shown. With this knowledge, we can calculate $\frac{\delta \mathcal{L}}{\delta B}$ as shown in equation 2.9.

$$\frac{\delta \mathcal{L}}{\delta B} = \delta_c w_{(b,c)} + \delta_d w_{(b,d)} + \delta_e w_{(b,e)} \quad (2.9)$$

In more formal mathematical terms, if we know the $\frac{\delta \mathcal{L}}{\delta \text{net}}$, or δ , for each neuron in a layer with n neurons, then we can calculate the gradient for any neuron i 's activation in the previous layer containing m neurons as shown in equation 2.10.

$$\frac{\delta \mathcal{L}}{\delta m_i} = \sum_{j=1}^n \delta_j w_{(i,j)} \quad (2.10)$$

The sensitivity for the neurons in layer m can then be computed using the derivative of the activation function. Since this thesis only uses ReLU, the derivative is simple to calculate and shown in equation 2.11. Note that the ReLU derivative is undefined at 0, however, in practical cases using a derivative of 0 works fine.

$$f'(x) = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \\ \text{undefined} & x = 0 \end{cases} \quad (2.11)$$

Calculating Weight Gradients Once the sensitivity δ of neuron is known, calculating the gradients of individual weights and biases is possible. From a high-level, if we increase weight w by ϵ , then the product term of the net for the neuron will be $(w + \epsilon)a_i$, a net increase of $a_i \times \epsilon$. Therefore, the gradient for a weight is dependent on how large the weight's corresponding activation is. That means the weight corresponding to a large activation will have a much larger gradient than a weight corresponding to a small activation.

Returning to the previous example, the figure has now been updated to show how weight gradients for neuron C , this is shown in figure 2.4. The gradients for the 2 connecting weights are calculated as shown below. A_o and B_o are the activations of neuron A and B , respectively. As one would expect, the gradient of a weight is dependent on the magnitude of the neuron activation it is multiplied with, and the sensitivity of the neuron whose net it is summed with.

$$\begin{aligned} \frac{\delta \mathcal{L}}{\delta w_{(a,c)}} &= \delta_c A_o \\ \frac{\delta \mathcal{L}}{\delta w_{(b,c)}} &= \delta_c B_o \end{aligned}$$

Updating the Weights Once $\frac{\delta \mathcal{L}}{\delta w}$ is known for every single weight, the final step of backpropagation is to update the weights. This is performed by scaling the gradient for the weight by a value, known as the learning rate, η , and then subtracting it from the weight, since this will move the weight in the direction that lowers the loss. This is shown in equation 2.12.

$$w_{new} = w_{old} - \eta \frac{\delta \mathcal{L}}{\delta w} \quad (2.12)$$

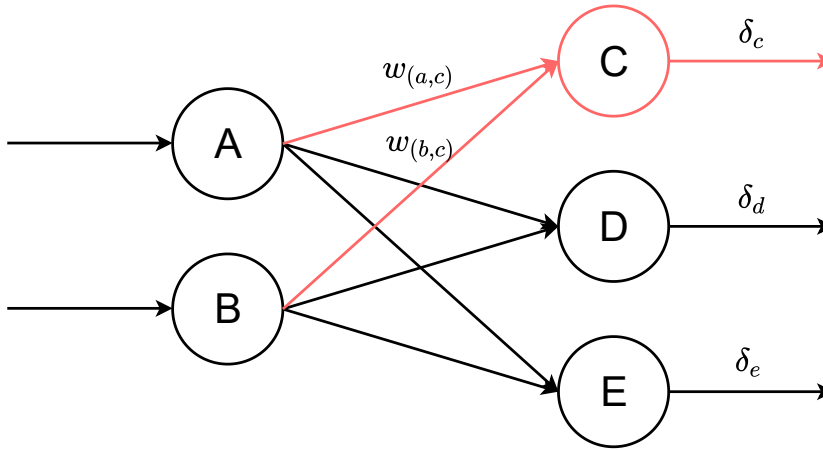


Figure 2.4: Example of computing weight gradients. Relevant values shown in red.

2.1.6 Hyperparameters

There are many hyperparameters to consider when designing a neural network. As described in section 2.1.5, the learning rate determines how of an impact the loss gradient has when updating the weight. Related to the learning rate is another hyperparameter known as momentum. The term is inspired from physics and has the effect of incorporating past updates in a geometrically decreasing fashion. We first define a few parameters:

m	—	the momentum parameter
v	—	‘velocity’
η	—	the learning rate
dw	—	the loss gradient for some weight or bias w .

The momentum-based update can then be mathematically represented in the following manner:

$$v = (m \times v) - (\eta \times w)$$

$$w = w + v$$

Each time each time we update w , previous updates update will also have an effect. A typical value for momentum is 0.9.

The final hyperparameter to be discussed in this section is batch size. Batch size determines the amount of forward and backward passes should be computed

before performing a weight update. There is no typical value for the batch size and the optimal batch sizes varies largely by dataset and problem type.

2.2 Deep-Learning Frameworks

Deep-learning has come into the spotlight in the past few years and as such, many popular and robust frameworks have been developed. Some of the most popular frameworks are TensorFlow which is developed by Google [AAB⁺15], and PyTorch which is developed by Facebook [PGC⁺17]. Keras is another popular framework that has introduced the most popular syntax style for describing neural networks [C⁺15]. Caffe was one of the earliest public frameworks available and remains in widespread use today [JSD⁺14]. These frameworks are generally relatively simple to use and deliver high performance.

2.2.1 PyTorch

For this thesis, PyTorch has been chosen as the framework to construct a model against which to benchmark my results. PyTorch offers a simplistic interface to build highly customizable neural networks. In addition, it also has support for GPU-training, thus both CPU and GPU benchmarks can be obtained.

2.3 Related Work

With the surge in popularity of neural networks, there has been a lot of research focusing on improving the performance of inference and training. Zhao et al. developed a data-streaming solution (F-CNN) using an FPGA to perform 32-bit floating point training and inference. They used a CPU to communicate weights, addresses, and training data over a PCI-E bus, ultimately obtaining roughly a 4 times speedup over a CPU implementation and a 7.5 times more power-efficient compared to a GPU implementation [ZFL⁺16].

The Alternative Computing Technologies lab at the Georgia Institute of Technology has begun promising work on a framework called TABLA, that generates FPGA code for a multitude of machine learning models. The framework is based on creating individual processing engines inside processing units and thus creating a more generalized design by having schedulers assign work to these processing units. Training and inference have been tested with promising results

on a Xilinx Zynq-7000 SoC and the group are aiming to make it public in the near future [MPA⁺16].

The majority of work based on accelerators for neural networks has been focused on inference. Google’s Tensor Processing Unit (TPU) has found real-world success by performing inference on networks using Google’s TensorFlow Lite framework [JYP⁺17]. This framework sacrifices precision for speed and obtains great performance with little punishment to overall accuracy, due to inference being less accuracy reliant than training. The Eyeriss is another hardware accelerator chip that has been developed to improve inference. It is designed by the Eyeriss Project team at the Massachusetts Institute of Technology. Its primary contribution is a novel dataflow optimization for convolutional neural networks called RS, for row stationary. This optimization allows for a high percentage of data reuse during inference [CES16].

The overwhelming majority of research regarding hardware acceleration for neural networks is focused on inference. Consequently, aside from the aforementioned F-CNN FPGA-based framework for training neural network, there has not been much research investigating neural network training on an FPGA. This work differs from the F-CNN as it investigates the effectiveness of using fixed-point arithmetic instead of floating point during training, allowing for extra speed at the sacrifice of accuracy. Furthermore, this work proposes that the main advantage of using a training accelerator is for datasets that train optimally with smaller batch sizes.

CHAPTER 3

Software Model

3.1 Overview

This section documents the general-purpose neural network framework that was written in C++ for this thesis. There is an example program that trains on the MNIST dataset and documents epoch-by-epoch training statistics. MNIST is a dataset of handwritten digits, containing 60,000 training images and 10,000 test images. The source code for the software model can be found in the appendix as well as online on GitHub.¹

3.2 Motivation

The software neural network framework was written so that the FPGA hardware model could be benchmarked against a CPU-based model that performs neural network inference and backward passes using the same method as the hardware model. This benchmark could be used to evaluate the performance

¹<https://github.com/erikgroving/NeuralNetworkHardwareAccelerator/tree/master/SWModel>.

of the hardware model. In addition, it could be benchmarked against professional open-source deep-learning frameworks that make use of advanced algebraic methods to perform computation such as matrix multiplication that inherently offer more efficiency. Furthermore, by developing a software model, the algorithmic integrity of the proposed network was able to be verified and tested in an expedient manner by using a well-known testing framework, Google Test [?]. Finally, if high floating-point precision were needed for training a network, then the software model could be used to learn the weights and parameters, and then subsequently be loaded into the weight BRAM of the FPGA hardware model.

3.3 Design

3.3.1 Layers

The software model was designed to be flexible such that any neural network architecture may be constructed so long as the layer types were implemented. The model currently supports 2D convolutional, fully connected, and pooling layers.

All layers are derived from a base class, `Layer`. Certain methods such as `forward()` and `backward()` must be implemented by all derived classes. There is then a `Net` class that contains a `vector` of `Layer` objects. This allows for a flexible design, as one only needs to add layers to the `Net` object. Furthermore, the model can easily be extended to other layer types so long as the layer type derives from `Layer`.

The non-linear activation function used in the model is ReLU because the derivative is trivial to compute. Compared to the sigmoid function, ReLU is much more computationally feasible for an FPGA hardware implementation, and therefore, ReLU was used in the software model so that both models would use the same activation function.

3.3.2 Training

The Softmax Function and Computing Loss Gradients The network uses an implicit Softmax function for the last layer since this converts the logits in the last layer to numbers that can be interpreted as probabilities, ideal for image classification.

The loss gradients for the neurons in the last layer are computed using multi-class cross entropy loss. Therefore, only one probability will account for loss, however, since each probability is an output from the softmax function which takes in all neuron outputs as input, all neurons in the last layer will have a loss gradient.

The derivative of this loss function is needed to perform backpropagation. We define \mathcal{L}_i as the loss for neuron i in the last layer and z_i as the output of neuron i . We also introduce y_i , which is 1 if x is an instance of class i and 0 otherwise. We can then compute the loss gradient for neuron i in the last layer quite simply as follows:

$$\frac{\delta \mathcal{L}_i}{\delta z_i} = z_i - y_i$$

Batch Size The software model supports batch training and thus a batch size is to be specified when creating an instance of a new network.

Learning Rate and Momentum The software model learns using stochastic gradient descent. As such, the network is configured with a learning rate and momentum. The learning rate may be manually readjusted during training epochs.

3.4 Source Code Structure

The software model contains a Makefile and three folders: *data*, *src* and *test*. The *data* folder contains the MNIST binary data files, and is loaded by the example program that trains on the MNIST dataset. The *src* folder contains the source code of the neural network framework. The *test* folder contains test made using the Google Test C++ testing framework. The Makefile is used to build the source as well as tests. This section will detail the source files in the *src* folder that are core to the software model framework. The files *main.cpp* and *parse_data{.cpp, .h}* will be described in section 3.5 that focuses on usage.

net{.cpp, .h} These files contain the definition of the **Net** class, the highest-level class of the network. After initializing a **Net** object, layers can be added to the neural network by calling the **addLayer()** method which will add a **Layer** object to a **vector**. The **Net** class also stores intermediate activations from the current inference, which are required when performing backward pass to

Name	Type	Description
<code>in</code>	<code>uint32_t</code>	Size of the input to the neural network.
<code>out</code>	<code>uint32_t</code>	Size of the output of the neural network.
<code>bs</code>	<code>uint32_t</code>	Size of the batch size to be used when training the net.
<code>lr</code>	<code>double</code>	The learning rate to be used during training of the network. Can be set and read using the functions <code>setLearningRate()</code> and <code>getLearningRate()</code> .
<code>momentum</code>	<code>double</code>	The momentum to be used when performing updates to the weights and biases of the network.

Table 3.1: Description of parameters for the constructor `Net` class.

calculate loss gradients. The key parameters to the `Net` object are set in its constructor, and are defined in table 3.1.

The `Net` class has a method `inference()` that computes the forward pass for a batch of inputs, thus the argument is a 2-d `vector`, with each outer index corresponding to an input. The `()` operator has also been overloaded to call `inference()`. This is all that is needed to compute a forward pass.

To compute the backward pass, `computeLossAndGradients()` should be called first. This method takes in the label data as a `vector` for the inputs as an argument and computes the loss gradients for the outer layer of the network. Next, a call to `backpropLoss()` should be made; this method propagates the outer layer loss gradients back through the neural network. After the loss has been back-propagated, weights of each `Neuron` in the network should be updated by calling `update()`. Previously cached forward pass activation data should then be cleared with a call to `clearSavedData()`.

layer.h This file contains the `Layer` class, which serves as the base class for all the different types of layer classes in the framework. It contains virtual methods `forward()` and `backward()`, representing the forward and backward pass functionality that must be implemented. All layer classes must also contain a `getType()` method to identify the layer type, as well as methods for `updateWeights()`, `clearData()`, and `getOutput()`.

convolutional{.cpp, .h} These files contain the definition of the `ConvLayer` class, which implements a 2D-convolutional layer, and derives from the `Layer` class. A unique method to the `ConvLayer` class is the `getWindowPixels()` method, which returns the pixels inside the filter window, and is used when

Name	Type	Description
<code>dim</code>	<code>uint32_t</code>	Dimensions of the input. The dimension is assumed square, meaning that <code>rows = dim</code> and <code>columns = dim</code> .
<code>filt_size</code>	<code>uint32_t</code>	Dimension of the filter used for the convolution, dimension also assumed square.
<code>stride</code>	<code>uint32_t</code>	Size of the stride
<code>padding</code>	<code>uint32_t</code>	Padding used for convolution.
<code>in_channels</code>	<code>uint32_t</code>	Amount of channels in the input.
<code>out_channels</code>	<code>uint32_t</code>	Amount of channels in the output.

Table 3.2: Description of parameters for the `ConvLayer` class.

computing both the forward and backward passes. The class' constructor and key parameters are described in table 3.2.

fullyconnected{.cpp,.h} These files define the `FullyConnected` class. The class only has two defining parameters in its constructor: `in` and `out`, which are of type `uint32_t` and specify the input and output size to the layer, respectively. It derives from the base `Layer` class, so methods such as `forward()` and `backward()` are also implemented.

pooling{.cpp,.h} These files define the `PoolingLayer` class. The class derives from `Layer` and performs a 2D 2×2 max pooling operation. There are three main parameters for the class: `dim_i`, `dim_o`, and `channels`. The parameters `dim_i` and `dim_o` specify the dimension of the input and output feature vectors. Since the layer currently only performs 2×2 max pooling, `dim_o` will always be half of `dim_i`, though if different types of pooling filters were to be supported, then `dim_o` would be necessary. The `channels` parameter is used to specify the number of channels of size $\text{dim}_i \times \text{dim}_i$ present in the input.

neuron{.cpp, .h} These files define the `Neuron` class. The `Neuron` class is the computational building block of the fully connected and convolutional layers. The fan-in of the neuron is specified in the constructor. Weights should be initialized using the `initWeights()` method, which implements He initialization [HZRS15]. He initialization randomly initializes weights using a normal distribution with a mean of 0 and a variance of $\frac{2}{\text{fan_in}}$.

The class implements all necessary computational elements for a neuron in a

neural network. During a forward pass, a neuron's net and activation are computed with `computeNet()` and `computeActivation()` respectively. When computing the backward pass, the gradients for the neuron's weights are computed using `calculateGradient()`. Weights can be subsequently updated using the `updateWeights()` function. Finally, all gradient data can be cleared using `clearBackwardData()`.

3.5 Usage

This section will show how the software model may be used for image classification. In the following example, the software model will be trained to classify handwritten digits from the MNIST database. Each image is a handwritten digit of size 28×28. The relevant files specific to this example are *main.cpp* and *parse_data.cpp*.

Load the Training and Testing Data The first step to any neural network problem is to load the training and testing dataset. The MNIST dataset is provided as binary files and helper functions to load the data have been provided in *parse_data.cpp*. Training and testing data can be loaded as shown below.

```
1  std::vector< std::vector<double> > trainX;
2  std::vector<int> trainY;
3  std::vector< std::vector<double> > testX;
4  std::vector<int> testY;
5  trainX = readImages("data/train-images.idx3-ubyte");
6  trainY = readLabels("data/train-labels.idx1-ubyte");
7  testX = readImages("data/t10k-images.idx3-ubyte");
8  testY = readLabels("data/t10k-labels.idx1-ubyte");
```

Create a Net Instance The next step is to create a `Net` object with the relevant hyperparameters to be used for the neural network. The below code accomplishes this.

```
1  int      input_size   = 28*28;
2  int      output_size  = 10;
3  int      batch_size   = 200;
4  double   momentum     = 0.9;
5  double   lr           = 0.01;
6  Net net(input_size, output_size, batch_size, lr, momentum);
```

Create Layer Objects and Add them to the Net Object After the `Net` object has been created, layers need to be added to the network. Two configuration options are present in *main.cpp*; one implements a 7-layer convolutional neural network, and the other implements a 4-layer fully connected neural network. The below code snippet shows how the 7-layer convolutional neural network is implemented. The software model was designed with simplicity in mind, so the below code is relatively straightforward to follow.

```
1 Layer* conv1 = new ConvLayer(28, 3, 1, 1, 1, 8);
2 Layer* pool1 = new PoolingLayer(28, 14, 8);
3 Layer* conv2 = new ConvLayer(14, 3, 1, 1, 1, 8, 16);
4 Layer* pool2 = new PoolingLayer(14, 7, 16);
5 Layer* fc1 = new FullyConnected(16*7*7, 64);
6 Layer* fc2 = new FullyConnected(64, 10);
7
8 net.addLayer(conv1);
9 net.addLayer(pool1);
10 net.addLayer(conv2);
11 net.addLayer(pool2);
12 net.addLayer(fc1);
13 net.addLayer(fc2);
```

Train the Net In *main.cpp*, a function `trainNet()` has been implemented, which trains the net using batch training. The actual training for a given batch only requires 5 lines of code, and is shown below.

```
1 net(in_batch);
2 net.computeLossAndGradients(out_batch);
3 net.backpropLoss();
4 net.update();
5 net.clearSavedData();
```

Build and Run the Model Compile the code by running `make` in the *SW-Model* directory. The model will then train for the amount of epochs specified in the call to the `trainNet()` function in `main()`. Since the model is initialized with random weights, the final result of training is non-deterministic. Output similar to the output shown in figure 3.1 can be expected. In this case, the fully connected model was used, and train to a maximum accuracy of 97.62%. it is also worth noting the expected differences in loss and accuracy between the training and test datasets. This discrepancy is expected as the network never learns from the test dataset. The difference between test and training dataset accuracy is normally used to quantify how well the network is able to generalize from the training dataset.

```
1 Running software model...
2 Starting Accuracy
3 Total correct: 1022 / 10000
4 Accuracy: 0.1022
5
6 Epoch: 0
7 --- Training Stats ---
8 Total correct: 54914 / 60000
9 Accuracy: 0.915233
10 Loss: 0.290908
11 --- Test Stats ---
12 Total correct: 9183 / 10000
13 Accuracy: 0.9183
14 Loss: 0.280574
15
16 Epoch: 1
17 --- Training Stats ---
18 Total correct: 56213 / 60000
19 Accuracy: 0.936883
20 Loss: 0.218062
21 --- Test Stats ---
22 Total correct: 9390 / 10000
23 Accuracy: 0.939
24 Loss: 0.214584
25
26 ...
27
28 Epoch: 36
29 --- Training Stats ---
30 Total correct: 59168 / 60000
31 Accuracy: 0.986133
32 Loss: 0.0516957
33 --- Test Stats ---
34 Total correct: 9762 / 10000
35 Accuracy: 0.9762
36 Loss: 0.0845137
```

Figure 3.1: An expected output from using the software model on the provided MNIST dataset. Epochs 2-35 omitted for brevity. In this training run, the network reached a maximum test set accuracy of 97.62%.

3.6 Testing

To ensure the correctness of the software model, several test suites were created during development. Source code for the test suites can be found in the *test* folder as well as in the appendix.

source code in appendix

3.6.1 Test Suites

Four test suites were created during the development of the software model. The test cases were written to test features as they were developed. As such, the tests include neuron functionality, forward pass for fully connected and convolutional layers, and finally a gradient checking test suite to verify the backward pass. This section elaborates on the test suites that were used during development.

Neuron Testing The neuron test suite, found in *neuron_test.cpp*, contains one primary test case that sets the weights of a neuron, computes the activation, and verifies that the activation is correct.

Fully Connected Forward Pass The test case for a fully connected layer's forward pass is located in *fullyconnected_test.cpp*. The test case creates a `FullyConnected` layer that has 3 inputs and 4 outputs. The weights are then set and an input is sent forward through the layer. Each of the 4 outputs are then verified to be correct.

Convolutional Forward Pass There is a test case to verify the convolutional forward pass located in *conv_test.cpp*. The test creates a convolutional layer that takes a 2×2 feature vector with 2 channels, uses a 3×3 filter for convolution, uses a stride and padding of 1, and produces 2 output channels. Weights and inputs were the arbitrarily assigned and the forward pass was computed and verified against the output that had been previously calculated manually.

Gradient Checking It would be very tedious and error-prone to debug the backward pass of a neural network using manual calculations, thus the general standard method of testing the gradients computed during a backward pass

is to use gradient checking. Note that during the backward pass, all the loss gradients for every single weight and bias are calculated. For every weight (and bias), the partial derivative $\frac{\delta \mathcal{L}}{\delta w_i}$ is computed. Gradient checking verifies that the mathematically computed analytic derivative aligns with a numerically estimated derivative [Kar]. The numerical gradient can be computed as follows:

$$\frac{\delta \mathcal{L}(w_i)}{\delta w_i} = \frac{\mathcal{L}(w_i + \epsilon) - \mathcal{L}(w_i - \epsilon)}{2\epsilon}$$

The partial derivative of the loss with respect to a certain weight w_i can thus be estimated by calculating the loss after incrementing w_i by a small ϵ , calculating the loss after decrementing w_i by ϵ , and then dividing the difference by 2ϵ . As long as ϵ is rather small, the derivatives should be near exact. In these test cases, $\epsilon = 10^{-4}$. Once we have the analytic and numerical gradient, we can compute the relative error as shown below:

$$\text{Relative gradient error} = \frac{|\mathcal{L}'(w_i)_a - \mathcal{L}'(w_i)_n|}{\max(|\mathcal{L}'(w_i)_a|, |\mathcal{L}'(w_i)_n|)}$$

If the relative error is below a certain threshold, then it is safe to assume the gradient has been calculated correctly. In this test suite, the relative error threshold must be lower than 10^{-7} .

The two test cases in *gradient_check_test.cpp* perform gradient checks for a fully connected network and for a convolutional neural network. The fully connected network gradient check test creates a neural network with an architecture shown in figure 3.2.

The test then creates 10 random inputs, each having a random label. Each input sample is fed forward through the network and analytic gradients are computed for each weight. The numerical gradient is then subsequently computed for a random weight. The random weight can belong to any neuron and any layer. This process of choosing a random weight, calculating the numerical gradient, comparing it to the analytic gradient is then repeated 100 times. The test asserts that the relative error is less than 10^{-7} each time. A portion of the computed analytic and numerical gradients are shown in figure 3.3.

The convolutional gradient checking test is set up in the same manner as the fully connected gradient checking test, except that the network structure is different. The network is now a **convolutional layer — pooling layer — convolutional layer — fully connected layer**. The input is randomized 8x8 data, and convolutional layers use 3x3 filters with a padding and stride set to 1. The first convolutional layer has 3 output channels and the second convolutional layer has 3 input channels and 6 output channels. The code used to create the network is shown in figure 3.4.

```
1  int      input_size   = 100;
2  int      output_size  = 2;
3  int      batch_size   = 1;
4  double   momentum     = 0.9;
5  double   lr            = 0.001;
6  Net net(input_size, output_size, batch_size, lr, momentum);
7
8
9  Layer* fc1 = new FullyConnected(input_size, 98);
10 Layer* fc2 = new FullyConnected(98, 64);
11 Layer* fc3 = new FullyConnected(64, output_size);
12
13 net.addLayer(fc1);
14 net.addLayer(fc2);
15 net.addLayer(fc3);
```

Figure 3.2: Layer created for the fully connected gradient check test.

```
1  Layer: 2, Neuron: 0, Weight: 31
2  Analytic Gradient: -0.0638284 Numerical Gradient: -0.0638284
3
4  Layer: 0, Neuron: 93, Weight: 71
5  Analytic Gradient: -0.156235 Numerical Gradient: -0.156235
6
7  Layer: 1, Neuron: 34, Weight: 29
8  Analytic Gradient: -1.22615 Numerical Gradient: -1.22615
9
10 Layer: 1, Neuron: 12, Weight: 43
11 Analytic Gradient: 0.376021 Numerical Gradient: 0.376021
```

Figure 3.3: Results from the fully connected test using randomly sampled weights to perform gradient checking

```
1  int    input_size    = 8*8;
2  int    output_size   = 2;
3  int    batch_size    = 1;
4  double momentum      = 0.9;
5  double lr             = 0.001;
6  Net net(input_size, output_size, batch_size, lr, momentum);
7
8  Layer* conv1 = new ConvLayer(8, 3, 1, 1, 1, 3);
9  Layer* pool1 = new PoolingLayer(8, 4, 3);
10 Layer* conv2 = new ConvLayer(4, 3, 1, 1, 3, 6);
11 Layer* fc1   = new FullyConnected(4*4*6, output_size);
12
13 net.addLayer(conv1);
14 net.addLayer(pool1);
15 net.addLayer(conv2);
16 net.addLayer(fc1);
```

Figure 3.4: Layer created for the convolutional layer gradient check test.

3.6.2 Building and Running the Test Suites

The test suites requires Google Test to compile. Google Test can be downloaded online at GitHub ². The *googletest* directory should then be placed under the *SWModel* folder. The test suite can then be compiled using the provided Makefile and the following command:

```
1 > make all_tests
```

This will produce an executable in the *SWModel* directory called **all_tests**. The test suites can be run by invoking the executable. The output is shown in figure 3.5

²<https://github.com/google/googletest>

```
> ./all_tests
Running main() from ./googletest/src/gtest_main.cc
[=====] Running 6 tests from 4 test cases.
[-----] Global test environment set-up.
[-----] 1 test from ConvTest
[ RUN      ] ConvTest.TestForward
[          OK ] ConvTest.TestForward (1 ms)
[-----] 1 test from ConvTest (11 ms total)

[-----] 1 test from FCTest
[ RUN      ] FCTest.TestForward
[          OK ] FCTest.TestForward (0 ms)
[-----] 1 test from FCTest (10 ms total)

[-----] 2 tests from NeuronTest
[ RUN      ] NeuronTest.InitWeights
[          OK ] NeuronTest.InitWeights (0 ms)
[ RUN      ] NeuronTest.SetWeightsAndGetOutput
[          OK ] NeuronTest.SetWeightsAndGetOutput (0 ms)
[-----] 2 tests from NeuronTest (29 ms total)

[-----] 2 tests from GradientTest
[ RUN      ] GradientTest.FCGradientCheck
[          OK ] GradientTest.FCGradientCheck (950 ms)
[ RUN      ] GradientTest.ConvGradientCheck
[          OK ] GradientTest.ConvGradientCheck (2260 ms)
[-----] 2 tests from GradientTest (3223 ms total)

[-----] Global test environment tear-down
[=====] 6 tests from 4 test cases ran. (3329 ms total)
[ PASSED  ] 6 tests.
```

Figure 3.5: Test coverage output using the Google Test C++ testing framework to verify the correctness of the software model for both forward and backward passes.

CHAPTER 4

Hardware Model and Implementation

This chapter details the hardware designed during this Master’s thesis to accelerate neural network training. The current hardware implements both training and inference acceleration for the neural network architecture described in section 4.2.

Refer to [github](#), [appendix](#), and [project link](#)

4.1 Hardware Setup

The hardware model was implemented using a ZedBoard. The ZedBoard is a development board equipped with a Zynq-7000 XC7Z020 SoC. The Zynq series has both a processing system and programmable logic, where the processing system is a ARM Cortex-A9 based processor (hereafter referred to as the “PS”) and the programmable logic is an Artix-7 series FPGA. Bitstreams for the FPGA were generated using Vivado 2018.3 and PetaLinux boot images for the PS were created using Xilinx SDK. The hardware description language (HDL) code for the project was primarily written in SystemVerilog. The programs run on the PS were written in C.

Layer Name	Input Size	Output Size
FC0	784 (28×28)	98
FC1	98	64
FC2	64	10
Softmax	10	10

Table 4.1: The hidden and output layers in the implemented neural network

4.2 The Implemented Neural Network

The classical MNIST handwritten digit dataset was chosen as the problem setting for the hardware model as a proof-of-concept. This problem has been chosen to verify the value in designing accelerators that take advantage of the finer-grained parallelism present in neural networks. The network consists of an input layer, 3 fully-connected layers, and a softmax output layer. The input layer is a 28×28 grayscale image of a handwritten digit. The dimensions of the rest of the layers in the network are shown in table 4.1. Layers whose name starts with FC are fully-connected layers.

Note that in this implementation, while biases are supported for forward computation, they are not used as the MNIST dataset is already fairly normalized. As such, the biases read are always 0 during the forward pass, and during the backward pass, no updates or gradients are calculated for the bias. Note that the gradient of the bias would just be the gradient of the neuron, unless the neuron had a ReLU activation function with negative net, so implementing this update would be trivial as all neuron gradients are already calculated. In addition, the current implementation only supports online training (training using 1 labeled data item at a time, a batch size of 1); offline training using larger batches is not supported by this hardware model.

4.3 Design Goals

There were a few key principles that guided the overall design process throughout the development of the hardware accelerator. A core tenet was to maintain the project such that in the future HDL could be generated for training a network of any architecture so long as the desired layer types had an implementation. As a result, all layers have been modularized and internal components are parameterized. Designing in a modular and parameterizable fashion also allows for quick and easy readjustments to the neural network architecture if needed.

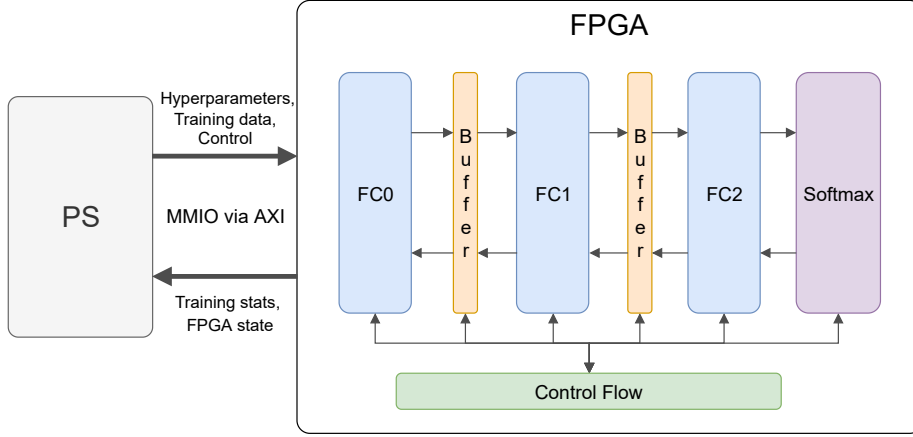


Figure 4.1: Architecture of the hardware accelerator

In addition, optimal usage of resources available was prioritized. For example, the limiting FPGA resource was the amount of digital signal processing slices (DSPs). Therefore, the FPGA design optimized the distribution of DSPs over other resources as opposed to saving an extra Block RAM (BRAM) block.

4.4 Overall Architecture

In the hardware model, both the Zynq’s PS and the FPGA were used to facilitate a cohesive and efficient architecture to accelerate neural network computation. The overall system architecture can be seen in figure 4.1.

Through memory-mapped I/O, the PS transfers neural net hyperparameters, training data, and control signals to the FPGA. The FPGA transfers training statistics and state data back to the PS. The interface is further described in sections 4.8 and 4.8.4.

Inside the FPGA, the neural network described in section 4.2 is implemented. Layers are connected in both forward and backward directions in order to support training. There are three types of primary modules in the top-level of the FPGA: fully-connected layers, interlayer activation buffers, and the softmax layer. In addition, there is a general control flow in the top-level with which all the primary modules interact.

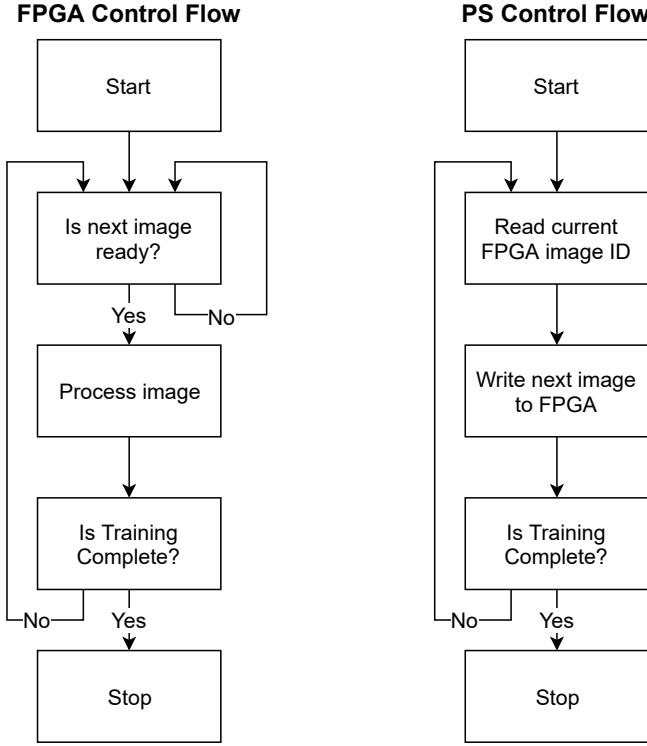


Figure 4.2: The high-level control flow of the training process

4.5 Training Process

The training process begins with the PS writing a 1 to the `start` register. This signals to the FPGA to start training whenever data becomes available. From the FPGA side, if the image in MMIO has ID equal to the FPGA's current image ID plus 1 (modulo image set size), then the image is ready. Otherwise, the FPGA will wait for the next image to be written. This process is summed up in Figure 4.2. The FPGA loops back around to 0 once the image set size is reached. The image set size is assigned via MMIO.

When the FPGA has processed the last image in the set during an epoch, the epoch counter is incremented and training stats will be available for the PS to read. If the epoch counter has reached the set number of training epochs, then training will stop, otherwise, the next training epoch will begin.

4.6 Computational Precision

In this implementation, a bit-width of 18 was chosen for all weight gradients and activations. This value was chosen because the multiplication portion of DSP slices have an input multiplicands with bit widths of 25 and 18 [Xila]. This thesis uses the Q number format to define precision types. For example, Q10.6 would mean that a 16-bit value has 10 integer bits and 6 fractional bits [cen01]. For this accelerator, activations have a precision of Q6.12. Weights and weight gradients both have a precision of Q1.17. These values were chosen through experimental analysis of minimum and maximum activation, weight, and gradients values using the software model described in chapter 3.

4.7 Module Architecture

As mentioned in the design goal section, one of the tenets of this design was to allow for modularity and parameterization, such that changing a network architecture would not require too much work. As such, there are a few global parameters defined, such as the amount of bits specified for the fixed-point precision. There are also parameters defined for each of the fully-connected layers. These parameters can all be found in the *sys_defs.vh* file in the Appendix

fpga appendix

, or on GitHub.

4.7.1 Fully-Connected Layers

The fully-connected layer modules implement both forward and backward passes. The general architecture is shown in figure 4.3. As DSP slices are limited, both the forward and backward computational units make use of the same resources to compute multiplications, known as the kernel pool. There are 4 modes of computation in the fully-connected layer: forward pass, backpropagating neuron gradients, computing weight gradients, and updating the weights. Of these 4 modes of computation, all except updating the weights make use of the kernel pool. This is because updating the weights makes uses of bit shifting instead of multiplication to multiply gradients by the learning rate.

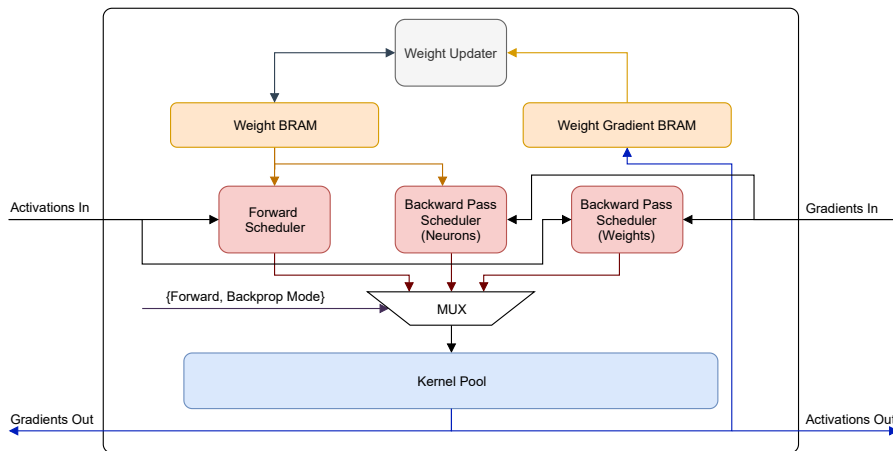


Figure 4.3: Architecture of the fully connected layer

The forward pass multiplies weights and input activations to produce output activations. Backpropagating neuron gradients multiplies weights by current layer input gradients to produce previous layer gradients as output. The weight gradient computation multiplies input activation from the forward pass by the current layer gradient, and then writing the computed gradient to the weight gradient BRAM.

Since being flexible and modular was one of the design goals, all the fully-connected layers use the same kernel and scheduler modules, with different parameters in the instantiation.

Scheduling Each of the computational modes needs to have a scheduler to generate addresses to be read and guide the computation. For this, the a generalized scheduler module was implemented. The scheduler uses two pointers starting from the head and middle of the BRAM, and iterates through the entirety during the forward pass. Since the weight BRAMs of each layer are different, certain parameters are assigned for instantiations of the scheduler are also different.

```

1 fc_scheduler #( .ADDR(`FC1_ADDR),
2   .BIAS_ADDR(`FC1_BIAS_ADDR),
3   .MID_PTR_OFFSET(`FC1_MID_PTR_OFFSET),
4   .FAN_IN(`FC1_FAN_IN))
5   fc1_scheduler_i (
6     // inputs
7     .clk(clk),

```

```

8      .rst(rst),
9      .forward(forward),
10     .valid_i(sch_valid_i),
11     //outputs
12     .head_ptr(head_ptr),
13     .mid_ptr(mid_ptr),
14     .bias_ptr(bias_ptr),
15     .has_bias(sch_has_bias)
16 );

```

Listing 4.1: The instantiation of the scheduler for the FC1 layer

The instantiation of the scheduler shown in listing 4.1 is similar across all the 3 fully connected layers, with only the parameters in the instantiations differing. The outputs are the pointers whose starting addresses are at the head and middle of the weight BRAM. In addition, there is a bias pointer and a signal to indicate if there is a bias.

Kernel The same kernel module is used in all the fully-connected layers. A high-level architecture of the computational kernel is provided in figure 4.4. Note that saturation checking is not shown in the figure for simplicity, though it is implemented and verified. The scarcest computational resource in this FPGA architecture are the DSP slices, thus the kernel has been designed so that all required forms of multiplication are supported in the network. This is why there are 3 different outputs. During both forward and backward passes, a kernel works on a single specific neuron until computation for that neuron has been completed, after which if there is still more work to do, will become computing for another neuron.

From the figure, the top output is the neuron gradient. This multiplies a weight with a gradient. Multiplying two Q1.17 values results in a Q2.34 product, which must be checked for saturation in the top 2 bits and the bottom 17 bits must be truncated to obtain a Q1.17 output.

The middle output is the weight gradient. The weight gradients computation multiplies a gradient with an activation. As gradients are Q1.17 and activations are Q6.12, the output is Q7.29. To convert the resultant Q7.29 result to the desired Q1.17 format required for a weight gradient, the top 7 bits must be checked for saturation and the bottom 12 bits truncated.

Finally, the bottom output is the net output calculated during the forward pass. This value becomes valid after performing n MACs, where n is the fan-in of a

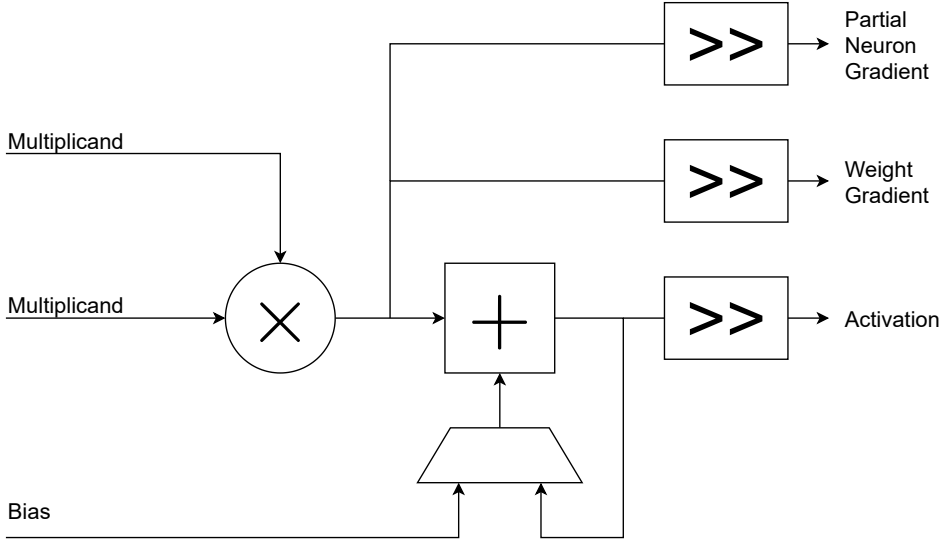


Figure 4.4: Architecture of the kernel, saturation checking not shown.

neuron. An input activation and corresponding weight are multiplied and added to either a bias or the running sum for the current in-progress net calculation.

Since the forward pass multiplies Q6.12 activations by Q1.17, the multiplied result is 36-bits, Q7.29. Since for some layers fan-in can be quite large, extra precision is used during the accumulation phase. The accumulated sum uses 32 bits: Q6.26, thus the internal sum must be checked for saturation and truncation of the bottom 3 bits for each MAC. The conversion from the internal sum precision of Q6.26 to the net output of Q6.12 is a simple truncation of the bottom 14 bits.

With this amount of internal precision during the forward pass, truncation error during the internal summation is sufficiently minimized. This is because the largest fan-in in this network is 784 in the FC0 layer. For each internal MAC, the 27th fractional bit onward is truncated when converted from the product with precision Q7.29 to the internal sum precision of Q6.26. Assuming that the 27th bit is equally probable to be 0 or 1, then the 27th bit is 1 approximately half of the time. This means that the average truncation per MAC is $0.5 \times 2^{-27} = 2^{-28}$. Given the 784 MACs of layer FC0, the expected truncation error is $784 \times 2^{-28} = 2.92 \times 10^{-6}$. The final truncation error from Q6.26 to Q6.12 truncates from the 13th fractional bit, resulting in an expected truncation error of $2^{-14} = 6.1 \times 10^{-5}$. Note that worst-case error is simply double the expected error, since it would assume that a 1-bit is truncated every time. Thus, truncation error from internal

summation is expected to be more than 20 times less than the truncation to the 18-bit net output. Therefore, truncation error is successfully minimized during the internal summation of net computation during the forward pass.

The kernel is parameterized to support reuse in all the layers. The kernel has 2 parameters that are specified upon instantiation: neuron fan-in and the amount of bits needed to represent a neuron ID in the layer.

Weight Updates During the weight update phase, scaled weight gradients are added to the corresponding weight. This phase is implemented by iterating over the weight BRAM and weight gradient BRAM in the fully-connected layers. Scaling down the gradient by a learning rate is accomplished by using right bitshifts, which allows for efficient computation without any real sacrifice in training accuracy, as learning rates are generally arbitrary; frequently chosen as some power of 10.

One update phase takes two cycles. In the first cycle, the weights and their corresponding gradients for an address are read from the BRAMs. In the second cycle, the gradient is scaled down using bitshifting and added to the weight. The result is then written to the weight BRAM in the same cycle. The address will then be incremented and the process continues until all weights have been updated. The control logic is implemented by simply using a counter. The address to the BRAM contains all the bits except the bottom bit, so the address is incremented once every 2 cycles. The bottom bit of the counter is used to indicate the update phase and determine read and write enables on the BRAMs.

Backpropagation Priority When the backward pass is in progress and the input neuron gradient's to a fully-connected layer become valid, there are two tasks that become ready to be performed. The first task is to use the valid input gradients to backpropagate neuron gradients to the previous layer. The second task is to calculate the weight gradients for the current layer. In this case, backpropagating the neuron gradients is given priority, because that way, once the gradients are backpropagated, the previous layer can also start performing its backward pass. Note that for the first layer, it is not possible to further backpropagate to previous neurons, so when the first layer receives its valid gradients, then it simply calculated its own weight gradients and then finishes.

Layer	# Kernels
FC0	196
FC1	16
FC2	2
Total	214

Table 4.2: Kernel allocation for the fully-connected layers in this implementation

4.7.1.1 Individual Fully-Connected Layer Implementation

While the scheduler and kernels are the same across fully-connected layers, the weight BRAM specifications are not, since number of neurons, kernels and fan-in are different for each layer. For this reason, all fully-connected layers needed to have separate files defining them. If the weights and gradients were loaded and stored to from DRAM instead of BRAM, then the fully-connected layer could be parameterized.

Kernels Per Layer Given that the layers all have different amounts of MACs, the amount of computational kernels to allocate to each layer should be balanced to roughly even out the amount of cycles the computational phases require. There are 220 DSP slices available on the FPGA, and each kernel uses 1 DSP slice. In this design, 215 kernels were instantiated and the distribution is shown in table 4.2. The mathematical reasoning for this allocation is discussed in Chapter 7.

Weight BRAM Initialization The BRAMs cannot be initialized with all 0 values as that devolve to being a linear classifier since all weights would have the same gradients, as explained in Chapter 2. Therefore, the weight BRAMs have been pre-initialized with values generated using He Initialization [HZRS15]. The BRAMs are configured using Xilinx coefficient (COE) files. The initialization is performed using floating point, converted to Q1.17 binary format, and then written to a file in COE format using a python script. This script, `weight_coeff.py` may be found in Appendix

appendix for this

or in the *misc* folder of the GitHub repository.

Address	Word Content
0	$w_{(0,0)} w_{(0,1)} \cdots w_{(0,96)} w_{(0,97)}$
1	$w_{(1,0)} w_{(1,1)} \cdots w_{(1,96)} w_{(1,97)}$
2	$w_{(2,0)} w_{(2,1)} \cdots w_{(2,96)} w_{(2,97)}$
...	...
782	$w_{(782,0)} w_{(782,1)} \cdots w_{(782,96)} w_{(782,97)}$
783	$w_{(783,0)} w_{(783,1)} \cdots w_{(783,96)} w_{(783,97)}$

Table 4.3: FC0 BRAM layout

BRAM Structure The memory storage and throughput requirements differed between the layers. As such, the weight and gradient BRAMs are organized differently. All the BRAMs use the true-dual port RAM configuration of the Xilinx Block Memory Generator IP Core, version 8.4. Each layer must be able to read 1 weight per kernel per cycle during computational steps to prevent kernels from idling. Note that the weight and gradient BRAMs are organized in the same way.

The FC0 layer has 196 kernels and 98 neurons, thus 196 weights need to be read per cycle. This is accomplished by having two ports of width 98 weights wide. A weight is 18 bits, so the word length for each port is 1,764 bits wide. Furthermore, since the fan-in of each neuron is 784 (28×28), there are 784 words in this BRAM. In total, the FC0 layer requires 49 36K BRAMs for the weights and the gradients each, or 98 total. The BRAM layout is shown in table 4.3. The format for the weights listed in the word content is $w_{(i,j)}$, which meaning weight i of neuron j .

The FC1 layer has 16 kernels and 64 neurons. 16 weights must be read each cycle to supply the neurons, so two ports of width 8 words are used. This means that the bitwidth of each port is 144 bits. The fan-in of each of the 64 neurons is 98, so there are 784 ($\frac{64 \times 98}{8}$) words in each BRAM for the weights and gradients. This results in needing eight 36K BRAMs for the FC1 layer. The first 98 words contain the weights for neurons 0-7. The subsequent 98 weights contain the weights for neurons 8-15. This continues through the entire contents of the BRAM, concluding with the last 98 words containing the weights for neurons 56-63. Since not every neuron is represented in every word, the memory layout is slightly different and shown in table 4.4.

The FC2 layer is the smallest fully-connected layer in this hardware model, containing 10 neurons each with a fan-in of 64. There are 2 kernels, so each port on the BRAM has a word width equal to 1 weight, or 18 bits. The depth of the BRAM is 640, and can be entirely contained within one 36K BRAM. The layout is shown in table 4.5.

Address	Word Content
0	$w_{(0,0)}w_{(0,1)} \cdots w_{(0,6)}w_{(0,7)}$
1	$w_{(1,0)}w_{(1,1)} \cdots w_{(1,6)}w_{(1,7)}$
...	...
97	$w_{(97,0)}w_{(97,1)} \cdots w_{(97,6)}w_{(97,7)}$
98	$w_{(0,8)}w_{(0,9)} \cdots w_{(0,14)}w_{(0,15)}$
99	$w_{(1,8)}w_{(1,9)} \cdots w_{(1,14)}w_{(1,15)}$
...	...
195	$w_{(97,8)}w_{(97,9)} \cdots w_{(97,14)}w_{(97,15)}$
196	$w_{(0,16)}w_{(0,17)} \cdots w_{(0,22)}w_{(0,23)}$
197	$w_{(1,16)}w_{(1,17)} \cdots w_{(1,22)}w_{(1,23)}$
...	...
293	$w_{(97,16)}w_{(97,17)} \cdots w_{(97,22)}w_{(97,23)}$
...	...
686	$w_{(0,56)}w_{(0,57)} \cdots w_{(0,62)}w_{(0,63)}$
687	$w_{(1,56)}w_{(1,57)} \cdots w_{(1,62)}w_{(1,63)}$
...	...
783	$w_{(97,56)}w_{(97,57)} \cdots w_{(97,62)}w_{(97,63)}$

Table 4.4: FC1 BRAM layout

Address	Word Content
0	$w_{(0,0)}$
1	$w_{(0,1)}$
...	...
63	$w_{(0,63)}$
64	$w_{(1,0)}$
65	$w_{(1,1)}$
...	...
127	$w_{(1,63)}$
...	...
576	$w_{(9,0)}$
577	$w_{(9,1)}$
...	...
639	$w_{(9,63)}$

Table 4.5: FC2 BRAM layout

Parameter Name	Brief Description
N_KERNELS_I	The width of the input write port. This is equivalent to the number of kernels of the previous layer.
N_KERNELS_O	The width of the output read port. This is equivalent to the number of kernels in the next layer after the buffer.
ID_WIDTH	The amount of bits needed to represent a neuron of the previous layer.
BUFF_SIZE	The amount of entries in the buffer.
LOOPS	Amount of times the buffer needs to be looped through for the next layer after the buffer to finish its computation.

Table 4.6: Parameters required for instantiation of the interlayer activation buffer.



Figure 4.5: The interlayer activation buffer

4.7.2 Interlayer Architecture

In this implemented hardware model, activations stored in interlayer buffers are stored directly in flipflops. This is because there are only 98 activations from FC0 to FC1 and 64 from FC1 to FC2. This results in 162 18-bit activations being stored in interlayer activation buffers, far within the resource limitations of the FPGA. The interlayer activation buffer module is also parameterized, so both buffers use the same SystemVerilog file. There are 5 parameters to be provided upon instantiation of the module, shown in table 4.6.

The architecture of the interlayer activation buffer is shown in figure 4.5. There is one write port which has a write width of N_KERNELS_I words. There are two read ports. The top one has a read width of N_KERNELS_O words and is used during the forward pass. The bottom one has a read width of 1 word and is used during the backward pass, particularly during the weight-gradient calculation phase.

4.7.3 Softmax Layer

To implement training of the neural network, meaningful gradients needed to be calculated for the output layer neurons. Cross-entropy loss, one of the most popular loss functions in deep learning, was chosen for this network. As such, the softmax function (also described in Chapter 2) needed to be implemented. The softmax function is shown again in equation 4.1 for convenience.

$$\sigma(\mathbf{x})_i = \frac{e^{x_i}}{\sum_{j=1}^C e^{x_j}} \quad (4.1)$$

The dataflow of the softmax layer is shown in figure 4.6. The softmax layer has 10 inputs in this network for MNIST digit classification. This layer is fully pipelined, so while there is a relatively long latency, the computation can still be performed quickly.

The implemented softmax function is also referred to as a numerically stable softmax. By subtracting a constant from the exponents, the final probabilities will not be affected. This is why the first step of the softmax layer is to subtract the maximum value input from all inputs. This then results in low, stable exponents being fed to the exponential function.

Since there is no support for the exponential function using fixed-point inputs in the Xilinx IP core repository, logits are first converted to 32-bit floating point numbers. After this, the exponential function for each input is calculated. The e^x core uses 1 DSP slice. The exponential function output is then converted from floating point back to fixed point. At this point, e^x is known for all the inputs, so all the numerators required for $\sigma(\mathbf{x})_i$ are known. To calculate the denominator, these values must also be summed up, so this occurs in the next stage of the layer. Finally, the numerators are divided by this denominator to finish the softmax process of converting the outputs from logits to probabilities.

4.8 PS – FPGA Communication

The processing system and the FPGA communicate via an AXI4 bus. In this AXI bus, the PS is the master and the FPGA is the slave.

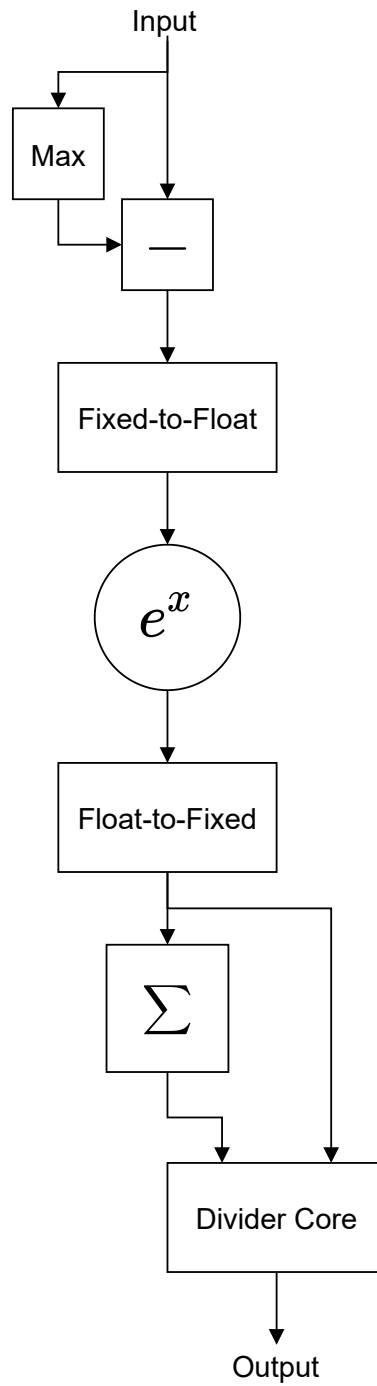


Figure 4.6: Architecture of the softmax layer

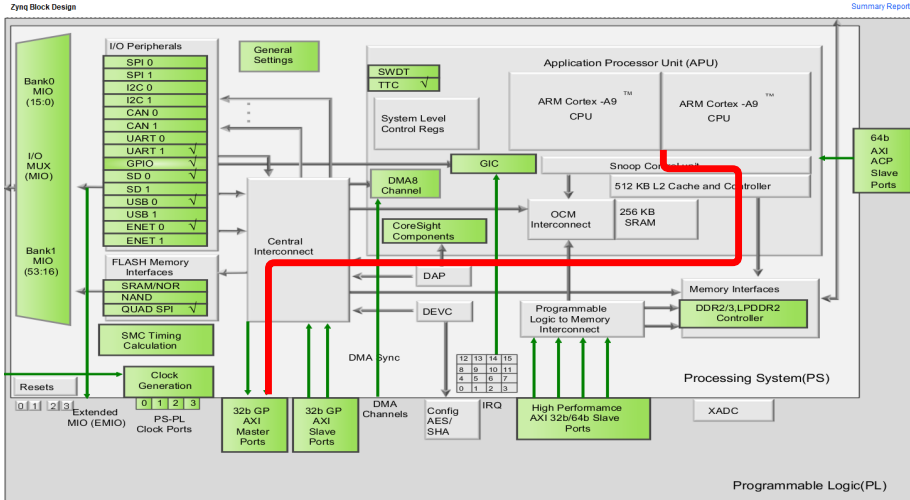


Figure 4.7: Communication from the PS to the FPGA

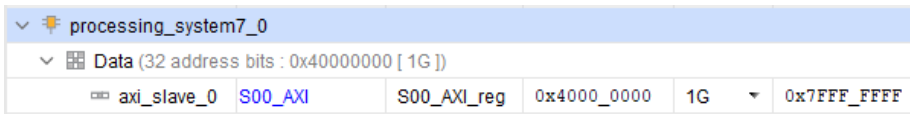


Figure 4.8: Specifying the Address Range for the AXI Bus

4.8.1 AXI Implementation for the PS

From the PS side, communication is performed as shown by the highlighted red line in figure 4.7. The AXI communication base address is 0x40000000 and spans until 0x7FFFFFFF. Having done this, by mapping a pointer to this location on `/dev/mem`, data that is written to or read from addresses within this region of memory will invoke an AXI bus transaction. This was set up by adding the *Zynq7 Processing System Version 5.5* IP core to a block diagram in Vivado, and defining an address range in the address editor as shown in Figure 4.8.

Once the address range and Zynq IP core has been added to the block diagram, C code to run on the PS must be written. There are only 2 things that need to be done to be able to start performing AXI bus transactions. The first step is to create a file handle by opening `/dev/mem` with the proper flags set. Then one only needs to memory map a pointer to this region. In the example shown in Listing 4.2, the pointer is of type `volatile`, because the order of reads and writes is critical during the training process when images are transferred to the

FPGA.

```

1 int handle = open("/dev/mem", O_RDWR | O_SYNC);
2 volatile ddr_data_t* ddr_ptr = mmap(NULL, 134217728,
    PROT_READ | PROT_WRITE, MAP_SHARED, handle, 0x40000000);

```

Listing 4.2: Code to allow the PS to perform AXI transactions with the FPGA

4.8.2 Running Programs from the PS

Several programs have been created for this project. Currently on the SD-Card, there are 3 programs: `train`, `train_small`, and `inference_only`. The source code, written in C, for these programs can be found in the

appendix

or in the `ps_code` folder of the project repository. The `train` program runs full-scale training on the entire MNIST dataset, splitting the data into training and testing sets of size according to a user-define. The `train_small` program uses a very small subset of the MNIST dataset to demonstrate the model's ability to have the network successfully learn an entire dataset. Finally, `inference_only` performs inference on every image in the MNIST dataset. A portion of terminal output from 5 epochs of training when running the `train` program is shown in Listing 4.3.

```

1 000 Loading MNIST images...
2 000 Loading complete!
3
4 000 EPOCH 1
5 000 Training Images: 2950/5000
6 Accuracy: 58.999997%
7 000Test Images: 43120/65000
8 Accuracy: 66.338462%
9 Active Cycles: 89107985 Idle Cycles: 228978233
10 Active Cycle Percentage: 28.013784%
11 Elapsed time: 6.36157 seconds
12
13 000 EPOCH 2
14 000 Training Images: 3623/5000
15 Accuracy: 72.460002%
16 000Test Images: 45904/65000
17 Accuracy: 70.621538%
18 Active Cycles: 178215969 Idle Cycles: 458029353
19 Active Cycle Percentage: 28.010575%

```

```

20 Elapsed time: 6.36297 seconds
21
22   *** EPOCH 3
23   *** Training Images: 4003/5000
24   Accuracy: 80.059999%
25   ***Test Images: 45345/65000
26   Accuracy: 69.761539%
27   Active Cycles: 267323953           Idle Cycles: 687061998
28   Active Cycle Percentage: 28.010046%
29   Elapsed time: 6.36260 seconds
30
31   *** EPOCH 4
32   *** Training Images: 4172/5000
33   Accuracy: 83.440000%
34   ***Test Images: 55452/65000
35   Accuracy: 85.310769%
36   Active Cycles: 356431937           Idle Cycles: 916126394
37   Active Cycle Percentage: 28.009084%
38   Elapsed time: 6.36324 seconds
39
40   *** EPOCH 5
41   *** Training Images: 4260/5000
42   Accuracy: 85.200000%
43   ***Test Images: 53002/65000
44   Accuracy: 81.541538%
45   Active Cycles: 445539921           Idle Cycles: 1145184941
46   Active Cycle Percentage: 28.008611%
47   Elapsed time: 6.36312 seconds

```

Listing 4.3: Training output from the `train` program

4.8.3 AXI Implementation for the FPGA

To implement the FPGA side of AXI communication, a block diagram to interface with the *Zynq7 Processing System* was created. Next, a custom IP core was created using Vivado's base AXI4 slave and then customized to meet the design needs of the project. The block diagram was then completed as shown in Figure 4.9. As can be seen, the Zynq's master AXI output is connected to an AXI interconnect which is then connected to the AXI slave port of the custom AXI module. Once the block diagram is complete, a wrapper for the block-diagram was generated and instantiated inside the top module of the design.

The final part of the FPGA AXI implementation was to modify the generated AXI module. The module was generated in Verilog, so it differed slightly from

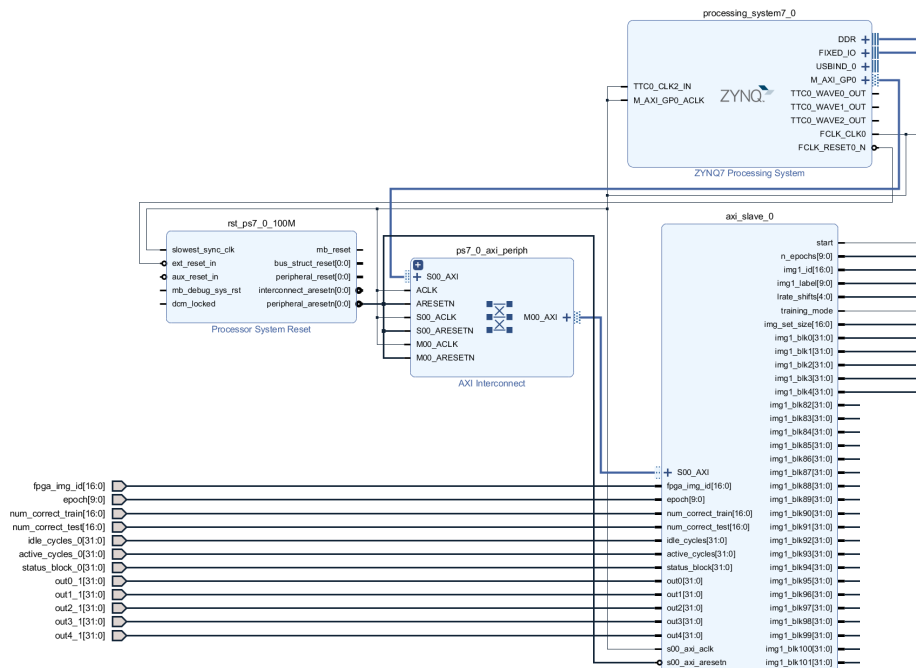


Figure 4.9: Block diagram to establish connection from the FPGA side to the AXI bus from the PS

the rest of the SystemVerilog project. It is also for this reason that there are 196 separate 32-bit registers to contain image data, since Verilog does not support 2D packed arrays as ports to modules. The Verilog code for implementing the memory map described in section 4.8.4 is included in Appendix

verilog appendix

4.8.4 Memory Map Layout

The memory map may be extended easily by adding or modifying address definitions to the AXI slave connected to the PS in the FPGA. In the PS code, one only need modify the `ddr_data` struct definition found in the C files of the PS source code. The memory map layout is defined in table 4.7. All addresses have a base address of 0x40000000. Note that values with addresses starting from 0x0 to 0x18 are registers written to by the FPGA. Values with addresses starting from 0x1C until 0x343 are written to by the PS. Output data from the FPGA is provided in the address range of 0x344 – 0x358.

4.9 PetaLinux

The boot image on the SD card has been modified to run Xilinx’s PetaLinux. This was done by using the 2016.4 prebuilt PetaLinux image as a base from Xilinx’s PetaLinux website [Xilb]. The image has been slightly modified by changing the `/etc/init.d/rcS` file to have the PS acquire a certain IP address and to mount the SD card to the filesystem. The *BOOT.bin* for booting Petalinux on the PS is created by writing a first stage bootloader *elf* file created by the Xilinx SDK for the Vivado project, the bitstream generated by Vivado, and the *u-boot.elf* file from the 2016.4 PetaLinux image. The boot image is created by using Xilinx SDK’s “Create Boot Image” utility. Aside from these changes to the 2016.4 SD card image, the other files in the image are not changed.

Cross-Compiling for PetaLinux Since the PS is an ARM-based processor, all C code for this project must be cross-compiled before it can be run on the PS. This is done by using the Linaro cross-compiling toolchain. A C file may be compiled to run on the PS using the below command.

PS – FPGA Memory Map

Offset	Name	Brief Description
0x0	fpga_img_id	The ID of the image that the FPGA is currently processing.
0x4	epoch	The current epoch in the FPGA
0x8	num_correct_train	The amount of correctly classified training images during the current epoch.
0xC	num_correct_test	The amount of correct classified test images during the current epoch.
0x10	idle_cycles	The amount of idle cycles in the FPGA since the start signal was received from the PS. An idle cycle is one in which none of the layers are performing any form of computation.
0x14	active_cycles	The amount of active cycles in the FPGA since the start signal was received from the PS. An active cycle is one in which at least one of the layers is performing computations.
0x18	status	A 32-bit register with many different flags from the FPGA, such as the layer states, for example.
0x1C	start	The start signal for training.
0x20	n_epochs	The number of epochs to train for in the FPGA
0x24	learning_rate	Value set to specify the amount of right shifts weight gradient should incur before updating a weight.
0x28	training_mode	Specifies whether the backward pass should be performed or not during computation.
0x2C	img_set_size	The size of the image set used during computation.
0x30	img_label	The label of the current image being computed.
0x34 - 0x343	img	Image data for the FPGA.
0x344 - 0x358	output	Output data from the last layer in the FPGA, before the softmax function is performed, so they are still logits in this case.

Table 4.7: Current memory map for communication between the PS and the FPGA.

```
1 > arm-linux-gnueabi-gcc <source_file> -o <executable>
```

CHAPTER 5

Hardware Model Testing and Verification

Verification is a vital part of hardware design. For this project, all relevant functionality in the FPGA was verified by simulation.

5.1 Simulation

During FPGA development, four different testbenches were created to test the functionality of the design. The first three were module level testbenches to test the scheduler, fully-connected layer, and softmax layers. The fourth testbench tests the entire FPGA design, thus verification of the fourth testbench means that all modules are functional. Therefore, this section will discuss verification of the fourth testbench, which is a full test of the network: `neural_net_top_tb.sv`, found in the directory *FPGA/FPGA.srscs/sim_1/new* of the GitHub repository as well as in Appendix

appendix for testbenches

5.1.1 Project Modifications to Simulate of the Design

To conduct the full-scale test of the network, an input needed to be provided to the network on which to perform training. This was done by using a BRAM to store a random input generated by the same Python script (`weight_coeff.py`) used to generate the weights for the weight BRAMS. When the simulation runs, the input to the network comes from this input BRAM rather than from the PS.

5.1.2 Testing Environment

The Vivado Simulator was used to perform simulation of the hardware. The testbench was run through Vivado's Tcl shell. During the testing process, a simulation would be ran and then diagnostic data in a test file could be. The below commands run from the Vivado Tcl show how to open the project, run the testbench for 50000 ns, and have all output written to a file.

```

1 Vivado% open_proj FPGA.xpr
2 Scanning sources...
3 Finished scanning sources
4 open_project: Time (s): cpu = 00:00:11 ; elapsed = 00:00:13
  . Memory (MB): peak = 322.016 ; gain = 71.828
5 Vivado% launch_simulation > sim_out
6 Vivado% run 50000 >> sim_out

```

5.1.3 Simulation Output

Verification and debugging was simple through the use of informative simulation output files. The project formatted signal data to be easy to read through the use of `$display` statements. An example of this is shown in listing 5.1. This example is from `neural_net_top.sv` and prints the current cycle number and FC2 output and gradient data; `sf` and `sf2` are scaling factors for activations and gradients, respectively. These scaling factors allow the fixed-point Q format values to be displayed as their floating point equivalent. These types of display statements are ubiquitous in the modules of the project. Once verified, most of these display statements were commented out to prevent clutter of the simulation output file.

```

1 `ifdef DEBUG
2 integer clk_cycle;

```

```

3 integer it;
4
5 always_ff @(posedge clk) begin
6     if (reset) begin
7         clk_cycle    <= 0;
8     end
9     else begin
10        clk_cycle    <=  clk_cycle + 1'b1;
11    end
12    $display("\n\n----- CYCLE %04d -----", clk_cycle);
13
14    $display("---FC2 GRADIENTS---");
15    $display("img_label: %d", img_label);
16    for (it = 0; it < `FC2_NEURONS; it = it + 1) begin
17        $display("%02d:\t%f", it,
18            $itor($signed(fc2_gradients[it])) * sf2);
19    end
20
21    $display("--- FC2 OUT ---");
22    $display("fc2_buf_valid: %01b" , fc2_buf_valid);
23    for (it= 0; it < `FC2_NEURONS; it=it+1) begin
24        $display("%02d: %f", it,
25            $itor($signed(fc2_act_o_buf[it])) * sf);
26    end
27 end
28 `endif

```

Listing 5.1: Example debug code for simulating the functionality of the hardware model

With this output redirected to a text file, signal data per cycle can be easily found. Furthermore, jumping to a previous or next cycle is quick. For example, in Vim, this can be done with by pressing ‘N’ or ‘n’, respectively. This method of debugging with Vim is shown in Figure 5.1, displaying the output generated from the code in Listing 5.1.

5.1.4 Correctness of Simulated Outputs

A Python script was written to aid in verification of the hardware simulation. This script, `fpga_forward_backward_pass_test.py`, is located in Appendix

```

erik@erik: /mnt/c/Users/Erik/Desktop/NeuralNetworkHardwareAccelerator/FPGA
100111 ----- CYCLE 1240 -----
100112 --FC2 GRADIENTS---
100113 img_label:      0
100114 00: -0.897644
100115 01: 0.096283
100116 02: 0.089554
100117 03: 0.100243
100118 04: 0.086243
100119 05: 0.113922
100120 06: 0.093460
100121 07: 0.100021
100122 08: 0.111229
100123 09: 0.106659
100124 --- FC2 OUT ---
100125 fc2_buf_valid: 1
100126 00: -0.021729
100127 01: -0.082764
100128 02: -0.155273
100129 03: -0.042480
100130 04: -0.192871
100131 05: 0.085449
100132 06: -0.112549
100133 07: -0.044678
100134 08: 0.061523
100135 09: 0.019531

100123,8 24%

```

Figure 5.1: Jumping from cycle to cycle to view debugging data using Vim

as well as in the *misc* folder of the project on GitHub.

This script parses the Xilinx coefficient files for the input to the network, as well as for all the neurons in all the layers and converts them to floating point numbers. The script then performs the forward pass using the parsed weights and prints the outputs. The script then computes the backward pass and also prints out all neuron and weight gradients. The hardware model can then be verified by checking that the outputs at each stage of computation align with the Python script. Note that values are not compared for equality, but for relative correctness, since the Python script uses floating point and the hardware model uses 18-bit fixed point.

Furthermore, to ensure that the script's computed outputs and gradients are correct, the script also implements gradient checking tests for itself. With this assurance, the script's computed output and gradients were successfully verified as correct and thus could be used as a baseline against which to compare the hardware model. Note that the gradient check testing for the testing script was based on the gradient checks implemented for the software model in Chapter 3, and example gradient check tests are shown in Listing 5.2.

```

1 > python3 fpga_forward_backward_pass_test.py
2 ../FPGA/FPGA.srscs/sources_1/ip/fc0_weights_1.17.coe
3 ../FPGA/FPGA.srscs/sources_1/ip/fc1_weights2_1.17.coe
4 ../FPGA/FPGA.srscs/sources_1/ip/fc2_weights_1.17.coe

```



```

5 Calculated gradient:      -0.003531676695546401
6 Numerical gradient:      -0.003531676693313557
7
8 Calculated gradient:      -0.006374946805618298
9 Numerical gradient:      -0.0063749433798498956
10
11 Calculated gradient:      0.0006677415295515585
12 Numerical gradient:      0.0006677441533042838

```

Listing 5.2: Gradient checks for randomly chosen weights in the Python verification script that uses inputs and weights read from the Xilinx coefficient files of the BRAMs in the hardware model. Only three non-zero gradients shown for brevity.

Forward Pass Verification The Python script was used to verify the correctness of the forward pass of the FPGA layer by layer. Forward pass layer outputs for softmax layer from the simulation and script are compared side-by-side in Listing 5.3. Since the softmax output depends on the outputs from FC0, FC1, and FC2, these layer outputs are not shown for the sake of space. From these tests, the simulated forward pass outputs of the hardware model are shown to be correct. The full outputs for every layer may be seen in the *HW_Verification* folder of the GitHub repository.

SIMULATION		PYTHON SCRIPT	
Neuron	Activation	Neuron	Activation
00	0.102348	0	0.10235213231346099
01	0.096283	1	0.09627338472902423
02	0.089554	2	0.08953177512141873
03	0.100243	3	0.1002532810000227
04	0.086243	4	0.08622264587243636
05	0.113922	5	0.11397991662882098
06	0.093460	6	0.09343919203092571
07	0.100021	7	0.10005157131620508
08	0.111229	8	0.11123918032286678
09	0.106659	9	0.10665692066481845

Listing 5.3: Softmax output. All 10 Neuron outputs shown.

Backward Pass Verification The backward pass was verified in the same way as the forward pass, though there are many more gradients than outputs. There is 1 gradient for each neuron and weight, totalling over 80,000 gradients. The forward pass outputs and backward pass gradients can be seen in their entirety in the *hardware_verification* folder of the GitHub repository.

The gradients in the backward pass all stem from the output layer gradients which come from the softmax function. The steps for deriving the gradients of the output layer is a softmax based neural network are explained in more detail in Chapter 2, though the gradients are essentially the softmax output visible in Listing 5.3 except that the neuron representing the inputs class label is subtracted by 1.

Listing 5.4 shows randomly selected weight gradients from each of the fully-connected layers. The weight gradients depend on the neuron gradients, thus the neuron gradients for that layer must be correct for the weight gradients to be correct; because of this, only weight gradients are shown in the figure, though neuron gradients are also available for viewing in the *hardware_verification* folder. As can be seen, the gradients are calculated to relatively high accuracy. This level of accuracy is directly correlated to the fact that the gradients are all Q1.17, maximizing the amount of fractional bits. Note that the 1 integer bit is required to represent the output layer gradient (since the input class label neuron is subtracted by 1), so the radix cannot be moved any further.

SIMULATION			PYTHON SCRIPT		
FC0					
Neuron	Weight	Gradient	Neuron	Weight	Gradient
09	593	0.000763	09	593	0.00076932259
19	711	0.006874	19	711	0.00688892029
37	412	-0.006149	37	412	-0.00613842723
57	128	0.000610	57	128	0.00061567956
74	485	-0.000282	74	485	-0.00027281649
FC1					
Neuron	Weight	Gradient	Neuron	Weight	Gradient
02	051	-0.003815	02	051	-0.00380934976
19	097	-0.019463	19	097	-0.01948172921
24	035	-0.013214	24	035	-0.01325251269
37	094	0.016045	37	094	0.01610831241
51	030	0.016563	51	030	0.01659535729
FC2					
Neuron	Weight	Gradient	Neuron	Weight	Gradient
01	043	0.015907	01	043	0.01595727359
03	002	0.016861	03	002	0.01688975169
04	057	0.005745	04	057	0.00578264697
08	023	0.024437	08	023	0.02451471064
09	024	0.000542	09	024	0.00055094484

Listing 5.4: 5 randomly selected weight gradients from each of the fully connected layers

Results

6.1 Benchmarking Models and Structure

Some of the results in this chapter are based on evaluating the hardware model (HWM) against other models implementing the same neural network. The other models include my software model (SWM), PyTorch running on the CPU (PyCPU), and PyTorch running on the GPU (PyGPU). My software model performs the same computations as the hardware model, so this provides insight to speedup over CPU without computational optimizations. The PyTorch CPU and GPU models then compare my hardware accelerator against heavily-optimized neural network frameworks. A training epoch in the following experiments is defined as performing learning on the 60,000 training images of the MNIST dataset. Inference experiments measure time to perform inference on all 70,000 images in the MNIST dataset.

6.2 Parallelism in GPU-based Training

Figure 6.1 shows PyGPU speedup for 1 training epoch over the 60,000 training images for varying batch sizes. As one might notice, speedup grows at

approximately the same rate as batch size. This is a textbook display of data parallelism; where it is clear that the images in the batch are processed individually by separate CUDA cores. The slope is able to stay linear even at massive batch sizes because a larger batch size means fewer updates to the weights. This results in the overhead from a larger batch size being negated by the decrease in amount of weight updates. The amount of forward and backward passes, which can be completely parallelized, remains the same. The weight update, the serial portion of training, is only performed once per batch. This means that a batch size of 1 performs 5 times more weight updates than a batch size of 5, which performs 20 times more weight updates than a batch size of 100. Therefore it is for this reason to see such a speedup.

From the figures, it is clear that the GPU model takes advantage of data-level parallelism to achieve performance, as epoch time is a near linear function of batch size. As a result, since the GPU-based implementation uses a coarser form of parallelism compared to the HWM, it would be illogical to benchmark speedup against the GPU with a batch size of 1 since the GPU's parallelism depends on batch size. Therefore, in the following experiments, the PyGPU model has been benchmarked using a batch size of 50 unless otherwise specified. It should be noted that the PyGPU model also performs 49 fewer weight updates compared to the other models as a result of this. Moreover, each weight update on a GPU would require reductions of partial gradient results from the CUDA kernels, so this should be taken into consideration when observing the following performance benchmarks.

6.3 Evaluation Hardware

The hardware model is evaluated using a ZedBoard equipped with a Zynq-7000 XC7Z020 SoC. The SWM and PyCPU both run on a Intel Core i7-4720HQ CPU. The GPU is an Nvidia GeForce GTX 970M equipped with 6 GB of GDDR5 RAM.

6.4 Performance

One of the most important metrics for an accelerator is runtime performance. While this hardware model is primarily focused on training, experiments to determine performance for both training and inference modes have both been conducted and are shown in this section.

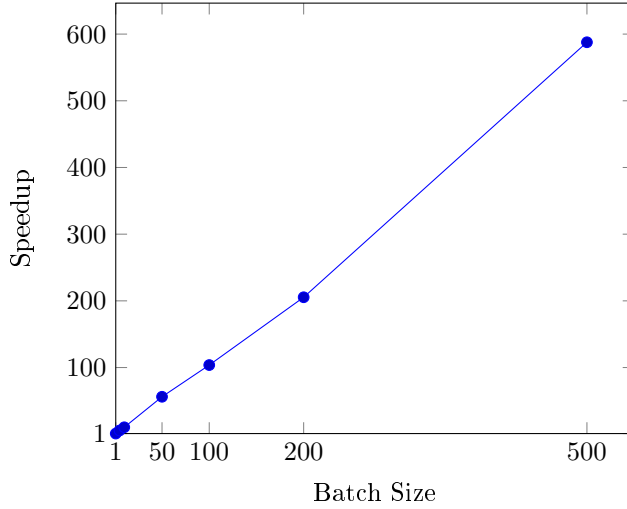


Figure 6.1: Speedup for 1 training epoch when training using different batch sizes for PyGPU. Speedup is calculated using the PyGPU model with a batch size of 1 as the reference.

6.4.1 Training

The average time for 1 training epoch has been recorded for each of the neural network models. The result is shown in Figure 6.2. This graph shows that the accelerator massively outperforms CPU models. Figure 6.3 shows the speedup of the models, using PyCPU as a baseline. Notably, the HWM achieves a speedup of of nearly equal to that of the PyGPU model.

6.4.2 Inference

Inference performance was also measured for each of the models. The result is shown in Figure 6.4. This graph shows that the accelerator also outperforms CPU models for inference, though falls short of the GPU model. Figure 6.5 shows the speedup of the models, using PyCPU as a baseline. The HWM achieves a speedup of 2.282 compared to the PyCPU model.

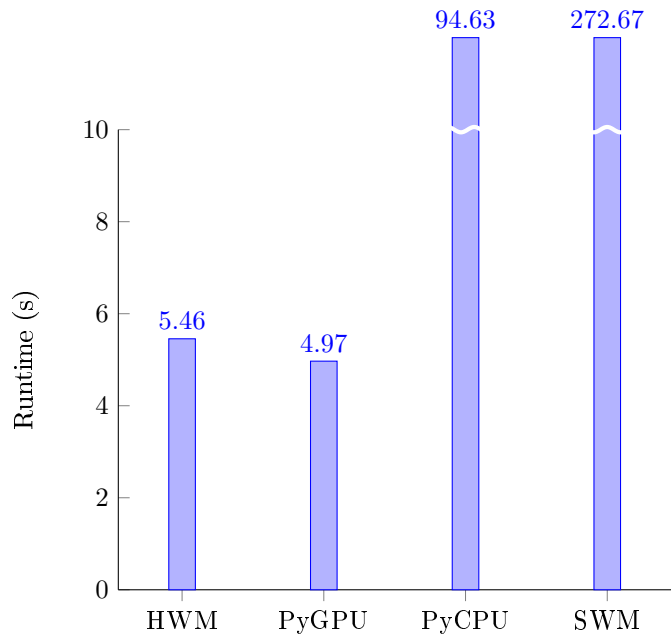


Figure 6.2: Training runtime for various network models

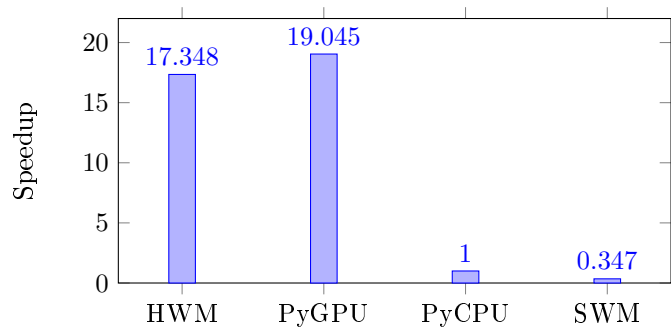


Figure 6.3: Training speedup using the PyCPU as a baseline

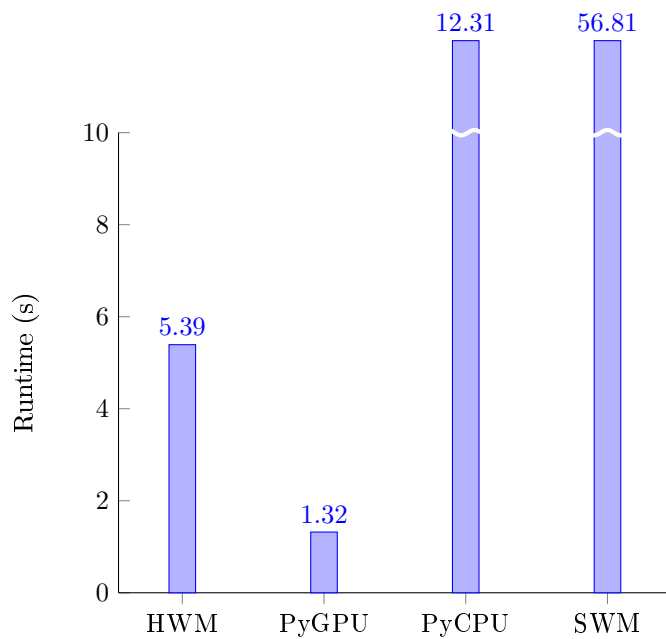


Figure 6.4: Inference runtime for various network models

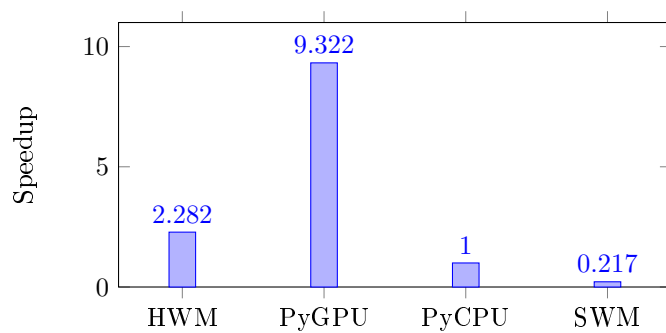


Figure 6.5: Inference speedup using the PyCPU as a baseline

	Active Cycle Percentage
Inference	25.13%
Training	69.20%

Table 6.1: Active Cycle Percentages for inference and training.

6.4.3 Active/Idle Cycles

To determine the impact of using MMIO via AXI bus to transfer image data between the PS and the FPGA, active and idle cycles were measured during training and inference. An active cycle is defined as a cycle on the FPGA during which at least one of the layers was computationally active. An idle cycle is thus defined as a non-active cycle.

An experiment was performed to measure active cycle percentages for the HWM during inference and training. The dataset was the entire MNIST dataset in both cases. The active cycle percentage for inference and training are shown in Table 6.1.

This experiment was performed to evaluate if the sending of input over MMIO was the bottleneck of the system. As can be clearly seen from the table, the MMIO transfer of training data was indeed the bottleneck. Furthermore, since backpropagation requires roughly double the amount of work compared to inference, it makes sense that training (which is both inference and backpropagation) is roughly a factor of 3 more active.

6.5 Training Accuracy

This section details the accuracy of the training process using the hardware accelerator. Varying training dataset sizes were chosen during the training process as the reduced precision training resulted in non-convergent training. As such, the training accuracy experiment conducted modified two variables: the learning rate and the training dataset size.

The tested learning rates were 2^{-7} , 2^{-8} , 2^{-9} , and 2^{-10} (0.0078, 0.0039, 0.00195, and 0.000977). This is because the hardware model performs the learning rate multiplication by using bitshifts. The experiments recorded the peak test data set accuracy during the training process. Note that the test dataset size for each run is 70000 minus the size of the training dataset. The results are shown in Figure 6.6. In this experiment, the highest test set accuracy, 85.845%, was

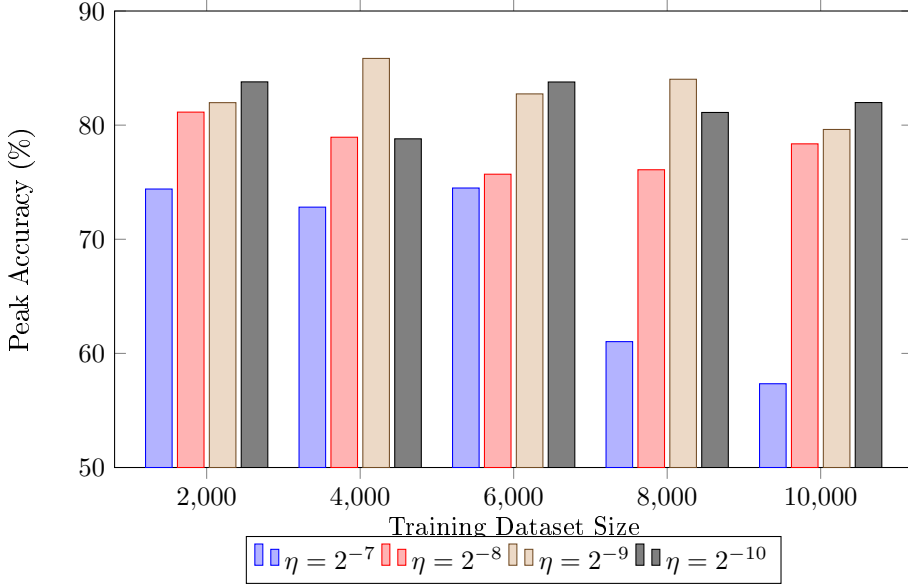


Figure 6.6: Maximum training accuracy reached for various learning rate and training set sizes.

achieved with a learning rate of $\eta = 2^{-9}$ and with a dataset of size 4,000. When using the same neural network architecture as the HWM, the SWM converges to 97.6% test set accuracy.

6.5.1 Stability of Training

As previously mentioned, due to the relatively low training precision of 18-bit fixed-point, the training process does not converge to a maximum training accuracy, but rather it will reach a maximum training accuracy, and then accuracy will degrade as precision errors accumulate over the training process. Training statistics for the first 10 epochs of the most optimal training configuration from Figure 6.6 illustrate this phenomenon and are shown in Figure 6.7.

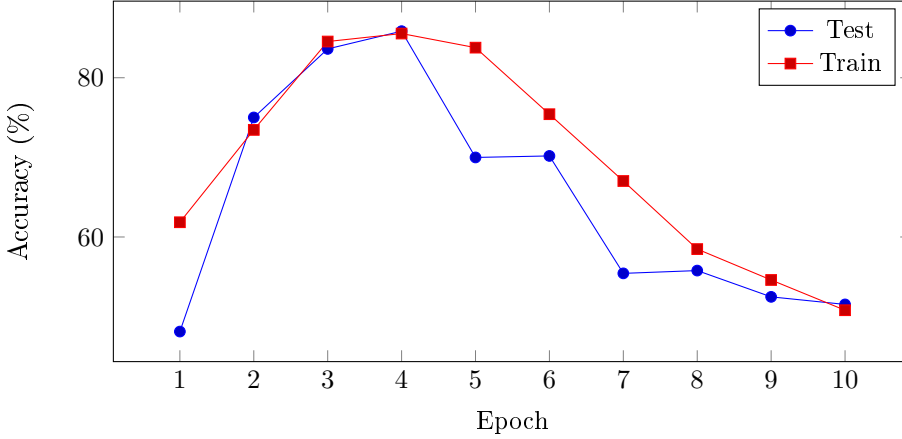


Figure 6.7: Epoch-by-epoch training data for an HWM configuration. Clearly visible degradation of accuracy instead of convergence after epoch 4.

6.6 Implemented Design

The design implemented for the FPGA is shown in Figure 6.8. As expected, the FC1 layer is by and large the most resource intensive, as it utilizes 196 kernels. It is interesting to observe the clustering of individual layer modules, while the interlayer activation buffer for FC0 and FC1 is widely spread out through the FPGA. This would indicate that this interlayer activation buffer was frequently routed to as a midpoint between FC0 and FC1.

It should be noted that implementation is a non-deterministic process and every design run should result in a slightly different implemented design. However, general trends for routing of the design tend to persist throughout multiple runs, despite the non-determinism of the placing and routing algorithms.

6.6.1 Resource Usage, Power, and Timing

The resource usage of the hardware model is shown in Table 6.2. As can be seen, the DSP slice is the scarcest resource, with LUTs and BRAMs also heavily being used. Overall, high utilization of the FPGA resources were made to optimize the performance of the accelerator as much as possible.

According to the Vivado report, the total on-chip power of the design is 2.798

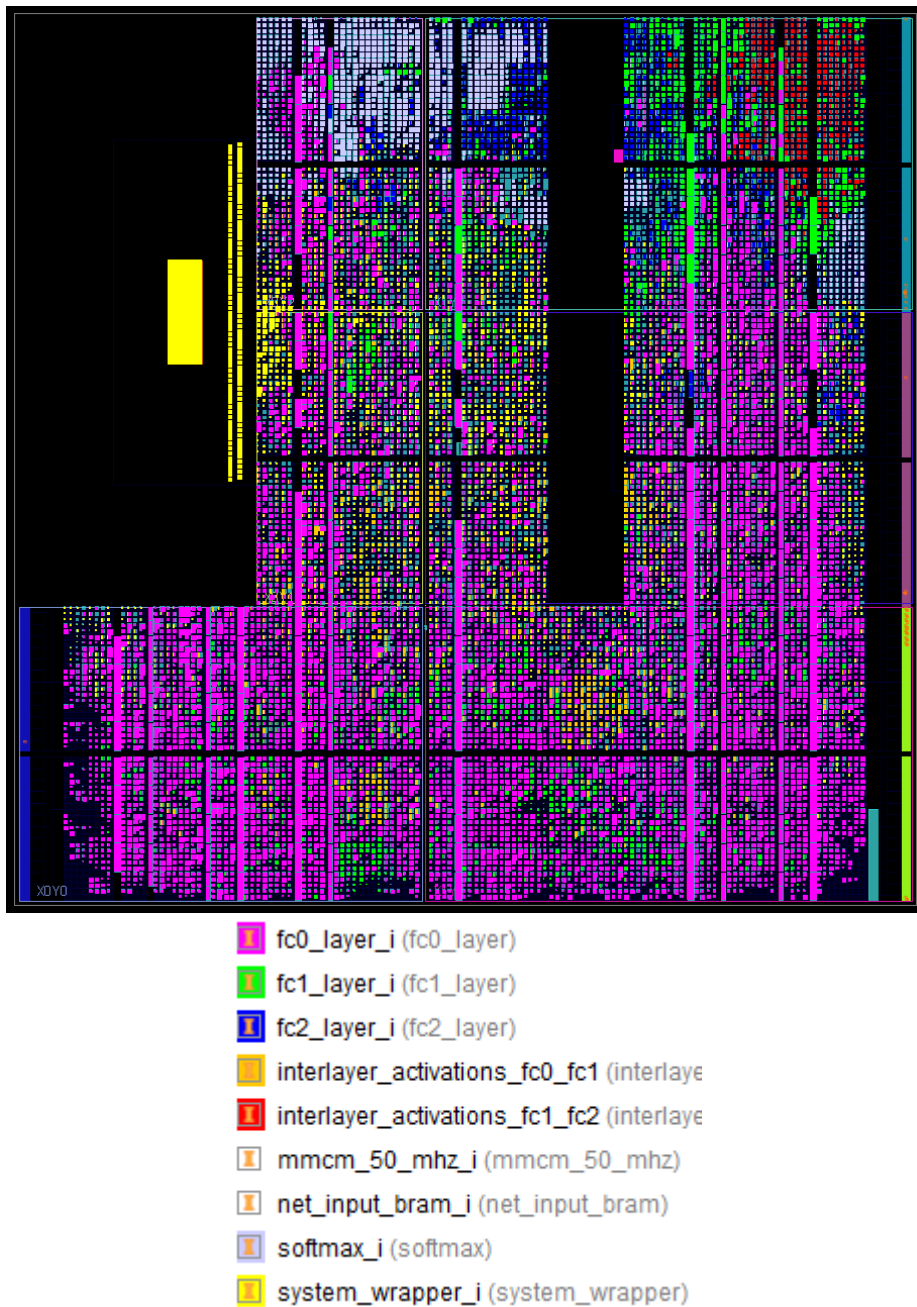


Figure 6.8: The implemented design of the hardware model

Resource	Utilization	Available	Utilization %
LUT	41132	53200	77.32
FF	54097	106400	50.84
BRAM	107.5	140	76.79
DSP	215	220	97.73

Table 6.2: Resource usage of the implemented design

Watts. While this number is reported with ‘Low’ confidence by Vivado, this wattage is far lower than typical GPU power consumptions. Power measurements have not been made for the internal GPU during these experiments, though measurements using the Torch framework were made and reported GPU average power at 94.19 Watts while performing training using the AlexNet architecture. The power measurements were performed with an Nvidia Tesla K20M GPU with 5 GB of GDDR5 SDRAM, a top-of-the-line GPU [Che16]. This is several magnitudes higher than the FPGA based solution in this thesis.

The implemented design for the hardware model is clocked at a frequency of 50 MHz. Placement and routing are both able to successfully complete with 0 timing violations. There was no need to improve frequency because the current performance bottleneck of this design stems from data transfer over the AXI bus and not from FPGA computational speed.

7.1 Allocating Computational Kernels for Performance

7.1.1 Allocating Based on Only Forward Pass Analysis

When computing an optimal allocation of kernels to the fully-connected layers, it was enough to account only for the forward pass. This is because in the backward pass, there are 2 times the kernel is used, during previous layer neuron gradient calculation, and during weight gradient calculation. In each case, the amount of multiplications is equal to the sum of the fan-ins of all the neurons. For the forward pass, every neuron receives an input from every neuron in the previous layer, so the amount of MACs will be:

$$MACs = \# \text{previous layer neurons} \times \# \text{current layer neurons} \quad (7.1)$$

For the backward pass, each backpropagated neuron gradient to a previous layer requires an MAC on all the neurons in the current layer. This must be done for each neuron in the previous layer, thus the total amount of MACs for backpropagating neuron gradients is also equivalent to Equation 7.1.

Layer	Fan-in per Neuron	# Neurons	MACs
FC0	784	98	76832
FC1	98	64	6272
FC2	64	10	640

Table 7.1: MACs per layer during the forward pass

For computing all the weight gradients in a layer, every weight for every neuron must be multiplied by a gradient. Each neuron in the current layer will have a # previous layer neurons weights, as that is the fan-in for each neuron. Thus, there same amount of multiplications is also equal to the expresion in Equation 7.1.

The backward pass also has a weight updating step, however, this uses bitshifts and not DSP slices to multiply the weight gradient by the learning rate. As such, the backward pass in this model uses exactly twice the amount of multiplications as the forward pass, so the optimal allocation of kernels is optimal for both the forward and backward pass. Note however, that even if the update step used multiplication, the amount of extra multiplications would be 1 for every weight, which would still be a multiple of the amount of multiplciations in the forward pass.

7.1.2 Distribution of the Kernels

Table 7.1 shows the fan-in, number of neurons, and thus the number of MACs per layer during the forward pass.

Furthermore, recall that a kernel operates on only 1 neuron at a time. Therefore, the amount of kernels allocated should be either a multiple or a factor of the amount of neurons such that work can be evenly distributed across the kernels. Given the 220 DSP slices on the FPGA, this becomes an optimization problem such that the amount of time for each layer to finish computing their outputs should be roughly the same. This would allow for pipelined input processing if offline training for large batch-sizes were to be implemented.

FC0 has 120.05 times more MACs than FC2. FC1 has 9.8 times more MACs than FC2. To balance runtime per layer, FC0 should then have roughly 120.05 times more DSPs than FC2 and so on. Using this information, we can write an equation for the amount of MACs and use substitution to come up with an ideal allocation scheme if allocate partial DSPs and ignore the fact that kernels

work on 1 neuron at a time.

$$DSP_{FC0} = 120.05 DSP_{FC2} \quad (7.2)$$

$$DSP_{FC1} = 9.8 DSP_{FC2} \quad (7.3)$$

$$\# \text{ DSPs} = 220 = DSP_{FC0} + DSP_{FC1} + DSP_{FC2} \quad (7.4)$$

Substituting into equation 7.4 using equations 7.2 and 7.3:

$$(7.5)$$

$$220 = 120.05 DSP_{FC2} + 9.8 DSP_{FC2} + DSP_{FC2} \quad (7.6)$$

$$220 = 129.85 DSP_{FC2} \quad (7.7)$$

$$DSP_{FC2} = 1.69 \quad (7.8)$$

Substituting the result from 7.8 into equations 7.2 and 7.3:

$$DSP_{FC0} = 120.05 \times 1.69 = 203.36 \quad (7.9)$$

$$DSP_{FC1} = 9.8 \times 1.69 = 16.60 \quad (7.10)$$

Thus if the 220 DSPs could be divided up ignoring all previous restrictions, the DSPs should be allocated according to Equations 7.8, 7.9 and 7.10. However, this is not possible, and DSPs are indivisible and the amount of kernels per layer should be a factor or multiple of the number of neurons, but it provides an maximum upper bound for performance.

Starting with layer FC0, which has 98 neurons, we should delegate 196 kernels. This is quite close to the optimal 203.36 computed above, thus layer FC0 should be allocated 196 kernels, which is $\frac{196}{203.36} = 96.38\%$ of the maximum upper bound. Continuing to layer FC1, which has 64 neurons, the optimal allocation is 16.6. The closest factor of 64 is thus 16 so 16 kernels are allocated, resulting in the same computation time as FC0 with $\frac{16}{16.6} = 96.38\%$ of the upper bound. Finally, layer FC2, with 10 neurons, rounding down and allocating 1 kernel would result in only $\frac{1}{1.69} = 59.17\%$ of the upper bound performance. Since a few kernels were freed up a from rounding down in FC0 and FC1, FC2 could be allocated 2 kernels, which allows it to finish faster than the optimally balanced latency for the 3 layers.

The final allocation of kernels is shown in Table 7.2. A pipelined solution is only as fast as its slowest step. Since FC0 and FC1 is the farthest away from the optimal upper bound, this solution performs at 96.38% of the theoretical upper bound for performance. It is worth re-mentioning that this upper bound is not actually possible since it assumes DSPs as divisible and that kernels can arbitrarily switched from neuron to neuron mid-computation, which would require finer-parallelism than what is supported in this architecture. Also note that since 214 of the 220 DSPs are used, the softmax layer was also able to use a DSP for calculation of the exponential function.

Layer	# Kernels
FC0	196
FC1	16
FC2	2
Total	214

Table 7.2: Kernel allocation between the fully-connected layers.

7.2 Cycle Analysis

This section calculates the computational cycles required for computing the forward and backward passes. Cycles spent pipelining or performing non-computational work are not included in this analysis. FP and BP are the amount of cycles needed to compute the forward pass and backward pass respectively. The amount of cycles for fully-connected layers in the forward pass can be roughly represented by $\frac{\#MACs}{\#kernels}$.

$$FP_{FC0} = \frac{76832}{196} = 392 \text{ cycles} \quad (7.11)$$

$$FP_{FC1} = \frac{6272}{16} = 392 \text{ cycles} \quad (7.12)$$

$$FP_{FC2} = \frac{640}{2} = 320 \text{ cycles} \quad (7.13)$$

For the softmax layer, there are several computational steps. For simplicity, conversions from fixed to floating point and vice versa are not included nor is the max circuit and the subsequent subtraction of the max. The softmax forward pass perform the exponential function, a sum of 10 values, and then a fixed point division. The exponential function requires 20 cycles, the summation requires 10 cycles and the divisor takes 46 cycles. This results in $FP_{softmax} = 20 + 10 + 46 = 76$ cycles. FP is computed in Equation 7.16.

$$FP = FP_{FC0} + FP_{FC1} + FP_{FC2} + FP_{softmax} \quad (7.14)$$

$$FP = 392 + 392 + 320 + 76 \quad (7.15)$$

$$FP = 1180 \text{ cycles} \quad (7.16)$$

Computing the backward pass is a bit more involved. Note that backpropagation of from the softmax layer to FC2 takes 1 cycle so it is not included. Each fully-connected layer first backpropagates neuron gradients for the previous layer and then computes the weight gradients, however, the previous layer can start computing the backward pass as soon as the neuron gradients are ready. Thus, the time for backpropagation to finish is not based on layer computation time,

but instead BP can be calculated by using neuron gradient computation time. Recall that the amount of computational cycles required to perform backpropagation of neuron gradients (NG) and computation of weight gradients (WG) is the same. As described in Chapter 4, updating the weights (UW) requires 2 cycles per-weight rather than 1-cycle per weight. Thus the UW step requires twice as many cycles as the others. The longest path during backpropagation is thus backpropagating neuron gradients from FC2 to FC1, from FC1 to FC0 and then performing WG in FC0 followed by UW .

$$BP = NG_{FC2} + NG_{FC1} + WG_{FC0} + UW_{FC0} \quad (7.17)$$

Substituting the forward pass cycles for the layers, as they are equivalent:

$$BP = 320 + 392 + 392 + (2 \times 392) \quad (7.18)$$

$$BP = 1888 \quad (7.19)$$

Thus, the amount of cycles spent performing computation during 1 training cycle is:

$$\text{Cycles} = FP + BP = 1180 + 1888 \quad (7.20)$$

$$\text{Cycles} = 3068 \text{ cycles} \quad (7.21)$$

This analysis has shown that the forward pass has 1,180 cycles of computation, and that one iteration of training contains 3,068 cycles of computation. Indeed, by viewing the output simulation file, this is confirmed as the forward pass finishes at cycle 1,235 and the training iteration concludes at cycle 3,145. The discrepancy in simulated cycles to computational cycles comes from the overhead of other actions in the design such as a max circuit, a pipelined addition reduction for neuron gradients, buffering of data between layers, non-computational data pipelining in layers, and floating-fixed type conversions. With a clock period of 20 nanoseconds (50 MHz) and ignoring data transfer overhead and delay, this would mean that 1 training iteration should take: $3145 \times 20ns = 62900ns$ or approximately 62.9 microseconds.

This result can be compared with the experimental results achieved. One training epoch was measured to take 5.455 seconds, as shown in Figure 6.2. Using the above calculated cycle results, at 62.9 microseconds per training iteration, training over a 60,000 image dataset we have:

$$t = 62.9e^{-6} \times 60000 = 3.774 \text{ seconds} \quad (7.22)$$

Recall however, how the active cycle percentage during training was only 69.2%. Thus, the analytically derived computation time should account for this:

$$t = 3.774 \times \frac{1.0}{0.692} \quad (7.23)$$

$$t = 5.454 \text{ seconds} \quad (7.24)$$

Compared to the experimentally measured epoch time of 5.455 seconds, the analytically computed training time of 5.454 seconds using cycle analysis is nearly exactly the same. The experimental data has thus successfully validated the cycle analysis.

7.3 Improving Performance

As was shown in Table 6.1 of the results section, the active cycle percentage for training is only 69.20%. Thus the first step to improve runtime performance would be to make data available for the FPGA to process faster. A suggested approach would be to use DRAM to stream data to the FPGA as done in other projects such as the neural network inference hardware accelerator proposed by Qiao et. al [QSX⁺16].

If an active cycle percentage of near 100% can be achieved from doing this, the next step to would be to optimize performance on the FPGA. The quickest route to doing this would be to improve the clock frequency. The design was clocked at a relatively low-frequency since it was not the bottleneck for performance. As such, the clock frequency was not investigated heavily during the design of the FPGA architecture, since improving this clock frequency would only increase the amount of time that the layers in the FPGA spend idling.

7.4 Granularity for Neural Network Computation

A key difference between training using this accelerator and training using GPUs is that this accelerator uses a much more fine-grained level of parallelism. While GPUs use data-level parallelism, this design uses neuron-level parallelism. Some attempts have been made to implement finer-level parallelism training on GPUs by Jiang et. al, though only yielded modest improvements of 1.58 to 2.19 times the speedup [JZL⁺18].

As a result, if one were to use the PyGPU solution to perform online training, then the GPU is 2.95 times slower than the CPU solution, which was 17.35 times slower than the hardware model (results in [?, inline]). Therefore, for online training, using fine-grained parallelism at the neuron level is the only place to find speedup, as data-level parallelism is not possible during online training where the batch size is only 1.

7.5 Ideal Learning Rate vs. Precision

One of the key intricacies in optimizing hyperparameters for the training process was balancing an ideal learning rate against increased precision error. Normally one need not think about precision when choosing a learning rate. However, in this accelerator, choosing what perhaps would be a more ideal learning rate might actually induce higher training error if it is too small since a smaller learning rate means fewer bits of gradient data are kept. For this project, this is compounded even further by the general notion that online training performs best with smaller learning rates.

For example, a learning rate of 0.001 results in a much better solution than a learning rate of 0.016 when a batch size of 1 is used for the implemented network architecture. However, the smaller in the implemented accelerator, this would mean an additional 4 bitshifts to the right when updating the weights. This results in a weight update that will have 4 fewer bits of information. In this project, the best training solution was found using 9 bitshifts to the right, or a learning rate of 0.00195. The weight gradients in this project are of number format Q1.17. Shifting this to the right 9 bits results in a Q1.8 number. This means that each and every weight update in the network only have 8 bits of information.

Vanishing Gradient Problem To add on even further to the aforementioned loss of information is that many of the gradients are already quite small. This is largely in part to a phenomenon referred to as the vanishing gradient problem. The vanishing gradient problem in neural networks refers to the fact that as one backpropagates further and further through the network, the gradients become smaller and smaller. To illustrate this phenomenon, the distributions of non-zero gradients in the implemented network has been plotted in Figure 7.1. As can be seen from the figure, gradients become smaller as backpropagation progresses from FC2 to FC0. Thus, as the gradients become smaller, they also become harder to represent using the limited precision of the hardware model. In an architecture where many bits of information are already being lost due to the learning rate, the vanishing gradient problem exacerbates the precision-induced error during training even more.

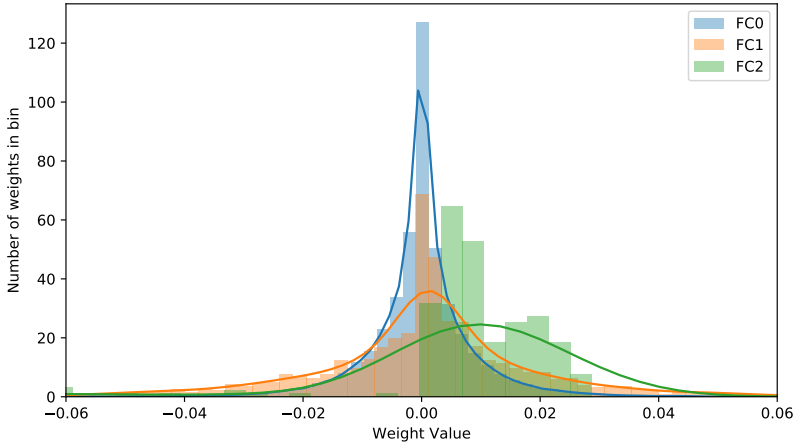


Figure 7.1: Weight gradient distribution for the 3 fully-connected layers after performing a single backward pass on the network using the input used to verify the hardware model from Chapter 5.

7.6 Potential Solutions for the Lack of Precision

As 18-bit fixed-point computation has been shown to be too imprecise to perform training on this network, potential solutions to this problem should be observed. It is the author's opinion that future accelerators for the training of neural networks would be best implemented by using 32-bit floating point as was done in the F-CNN accelerator by Zhao et. al [ZFL⁺16].

However, if 32-bit floating point is infeasible, or if accuracy is to be traded off for improved speed, area, and storage of weights, then investigations into designing accelerators using 16-bit or 24-bit floating-point computation could be made.

If the architecture must use fixed point, then the author would suggest first investigating 32-bit fixed-point computation with varying radices. If storage restrictions permit and the design performs multiplications using DSP slices, then a maximum of 36 bits would be supported for completing a multiplication using 2 DSPs in 1 cycle. This stems from the fact that 18 bits is the maximum width for one of the ports in a DSP-multiply. Otherwise, multiplication could also be a multi-cycle computation to trade off time for improved precision.

7.7 Weight Storage

The implemented neural network was specifically designed such that the weights and weight gradients could fit in the BRAM. Since 76.79% of the BRAM was utilized, the implemented network is representative of what architectures may be supported in BRAM as it comes close to hitting the upper limit of network size that can be supported entirely using BRAM. For networks larger than the implemented neural network for this thesis, other solutions such as a streaming weight and weight gradient datapath to DRAM would be required.

In addition, since the precision of the weight and weight gradients in this project proved to be inadequate for convergence to a local optimum during training, it should be noted that increasing precision would also increase BRAM utilization. This network would be able to use a maximum of 23-bits of precision for the weights while still fitting into BRAM at roughly 98.12% utilization. If more bits are needed for successful training then the network architecture would have to be made smaller or the hardware architecture would need to use a streaming datapath solution.

Discussion

8.1 Overall Performance

Regarding the performance, the accelerator has outperformed all compared CPU benchmarks. It performs online training with a speedup of 17.35 compared to the PyTorch CPU model. Considering how the PyTorch GPU achieved a 19.05 speedup using a batch size of 50, the accelerator was nearly able to keep pace. Furthermore, the GPU model does not use fine-level parallelism, so the accelerator achieves the highest speedup of all models for training with a batch size of 1.

8.2 Finely-Grained Parallelism

Training of neural networks in today's world is done almost exclusively using GPUs and occasionally using CPUs. This is a stark contrast compared to inference, for which many different chips such as Google's TPU have been developed [JYP⁺17]. However, as this thesis has shown, for neural network training problems that do not have vast amount of data parallelism available, there is no highly optimized solution. As such, the accelerator developed during the process of this thesis shows a massive potential for this side of training since it

takes advantage of the finely-grained parallelism available at the neuron-level, something not done by options available in today's world.

8.3 Limitations

8.3.1 Precision

Precision is a major limitation of training for the current design. It is the reason why the training process is not able to smoothly converge to a local optimum. This results in contradicting desires to have more bits of information available in weights gradients while at the same time having a low learning rate.

8.3.2 Data Transfer Rate

Another major limitation of this work is the method of transferring training data by using a memory-mapped interface between the PS to the FPGA. This approach was used for convenience, however, as the FPGA active cycle results from Table 6.1 showed, this approach is inefficient and became the largest bottleneck of performance for the design.

8.4 Future Work

While the potential for application-specific hardware accelerators training has been demonstrated in this thesis, there is a lot of potential for future work to improve the project.

Increased Precision As was demonstrated in the results section, training a neural network requires high precision computation. This is especially true for deeper neural networks as a result of the vanishing gradient problem. Therefore, increasing the precision, either via changing to floating point or using more bits in fixed point would be a great improvement.

Larger Batch Sizes Online training is only applicable to certain datasets. While the usefulness of an accelerator for online training has been shown, there

are also many datasets that converge faster by using a larger batch size and off-line training. In addition, a larger batch size provides a more accurate gradient of the actual loss function of the training set.

Since the amount of data-level parallelism increases with the batch size, it becomes increasingly harder to compete with the performance of GPUs. Furthermore, a solution to storing activations in memory to compute the backward must be designed. That being said, using a larger batch size would also open up the possibility to taking advantage of data-level parallelism and using an array of training accelerators. In such a setup, both data-level and neuron-level parallelism would be working together.

Additional Layer Types This design only implemented the fully-connected and softmax layer types. There are many other types of layers for neural networks, and this project could be expanded by implementing other layer types such as convolutional or pooling layers, which are frequently used in image recognition.

Backward Pass for Biases In the interest of time and since the input data is already fairly normalized, only the backward pass for weights was computed. A rather quick improvement to the project would be to implement the backward pass for biases, so that the network architecture could be applied to non-normalized datasets as well. The gradient for a bias is simply the gradient of the net, as it is added directly to the net. Therefore, the bias gradients are already known in the hardware, and all that would be done would be to add BRAMs for the biases and slightly modify the update phase to update the biases.

Additional Activation Functions In both the software and hardware models for this project, the ReLU function was chosen specifically due to its computational simplicity, quick convergence during training, and its ability to converge to strong local optima. That being said, there are still many other activation functions in the realm of neural networks that also achieve strong training results. As the dataset and network architecture changes, so may the the most optimal activation function. Other activations functions such as the sigmoid function, leaky ReLU, hyperbolic tangent, and many others may be preferred to ReLU under certain circumstances. These functions would require extra hardware support though, and thus would require more computational resources to implement. As a result, one should expect that the performance of the accelerator would not be quite so high as with the ReLU activation function.

Implement Streaming DDR Interface for FPGA Adding a streaming data interface for training data would reduce the amount of cycles during which the FPGA idles. This would be a strong improvement for performance. Adding a streaming data DDR interface for weights and activations would allow networks with larger footprints to be supported by the hardware model. Both of these modifications would be an overall improvement to the model.

Generated HDL for a Pre-Specified Network Architecture As one of the design goals was to be modular, if the streaming data interface were to be implemented, then it would be feasible to define a network architecture in a configuration file and create a program to generate HDL files for that network architecture. This would allow for a flexible, modular, FPGA-based framework that could implement any type of network, so long as the layer-types of that network were supported.

Conclusion

This thesis addressed the dearth of performance-optimized solutions for conducting online training of neural networks by proposing a novel hardware architecture. The proposed hardware accelerator achieves high performance by exploiting the fine-grain parallelism present at the neuron level during the training process.

After providing a background of neural networks, a software model was implemented to verify the chosen training algorithm for the hardware model. This provided the algorithmic foundation of the hardware model.

We then provided an explanation of the architecture and implementation of the hardware model. By carefully allocating computational kernel among the fully connected layers, a balanced computational scheme allowed for highly parallelized computation. The implemented design was a modular solution that resulted in a flexible design. Furthermore, full-scale and modular testbenches along with a convenient testing process resulted in the successful verification of the design.

The final design for the hardware accelerator was clocked at 50 MHz and satisfies all timing requirements. Furthermore, the implemented design uses low-power compared to GPUs. Experimental results of the hardware accelerator show that the proposed solution achieves a speedup of 17.35 compared to the next best online training model. At the same time, the accelerator is nearly as fast as the

GPU model performing training with a batch size of 50. The experiments also revealed that the bottleneck of the solution was from the MMIO communication between the PS and FPGA rather than the training on the FPGA. The results also showed that 18 bits of fixed-point precision is not enough to successfully converge to a local optimum during training, rather the training process will degrade in performance after a few epochs as precision error accumulates.

We then analyzed and discussed the experimental results, highlighting the need for more computational precision while showing the massive potential gains in performance from utilizing fine-grained parallelism.

Ultimately, this thesis has shown the feasibility of designing a hardware accelerator that uses neuron-level parallelism for the online training of neural networks. Furthermore, there are many potential future optimizations and improvements to the model that would increase both performance and functionality.

APPENDIX A

Stuff

APPENDIX B

Stuff2

Bibliography

- [AAB⁺15] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [C⁺15] François Chollet et al. Keras. <https://keras.io>, 2015.
- [cen01] ARM Info center. Arm developer suite axd and armsd debuggers guide, 2001. <http://infocenter.arm.com/help/index.jsp?topic=/com.arm.doc.dui0066d/CHDFAAEI.html>.
- [CES16] Yu-Hsin Chen, Joel Emer, and Vivienne Sze. Eyeriss: A spatial architecture for energy-efficient dataflow for convolutional neural networks. In *Proceedings of the 43rd International Symposium on Computer Architecture, ISCA '16*, pages 367–379, Piscataway, NJ, USA, 2016. IEEE Press.
- [Che16] Xinbo Chen. Energy efficiency analysis and optimization of convolutional neural networks for image recognition. Master’s thesis, Texas State University, may 2016.

- [GBB11] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 315–323, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR.
- [gma15] The mail you want, not the spam you don’t, Jul 2015. <https://gmail.googleblog.com/2015/07/the-mail-you-want-not-spam-you-dont.html>.
- [HZRS15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *CoRR*, abs/1502.01852, 2015.
- [JSD⁺14] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross B. Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. *CoRR*, abs/1408.5093, 2014.
- [JYP⁺17] Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre-luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike Daley, Matt Dau, Jeffrey Dean, Ben Gelb, Tara Vazir Ghaemmamghami, Rajendra Gottipati, William Gulland, Robert Hagmann, C. Richard Ho, Doug Hogberg, John Hu, Robert Hundt, Dan Hurt, Julian Ibarz, Aaron Jaffey, Alek Jaworski, Alexander Kaplan, Harshit Khaitan, Daniel Killebrew, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon, James Law, Diemthu Le, Chris Leary, Zhuyuan Liu, Kyle Lucke, Alan Lundin, Gordon MacKean, Adriana Maggiore, Maire Mahony, Kieran Miller, Rahul Nagarajan, Ravi Narayanaswami, Ray Ni, Kathy Nix, Thomas Norrie, Mark Omernick, Narayana Penukonda, Andy Phelps, Jonathan Ross, Matt Ross, Amir Salek, Emad Samadiani, Chris Severn, Gregory Sizikov, Matthew Snelham, Jed Souter, Dan Steinberg, Andy Swing, Mercedes Tan, Gregory Thorson, Bo Tian, Horia Toma, Erick Tuttle, Vijay Vasudevan, Richard Walter, Walter Wang, Eric Wilcox, and Doe Hyun Yoon. In-datacenter performance analysis of a tensor processing unit. In *Proceedings of the 44th Annual International Symposium on Computer Architecture*, ISCA ’17, pages 1–12, New York, NY, USA, 2017. ACM.
- [JZL⁺18] Wenbin Jiang, Yangsong Zhang, Pai Liu, Geyan Ye, and Hai Jin. *FiLayer: A Novel Fine-Grained Layer-Wise Parallelism Strategy for Deep Neural Networks: 27th International Conference on Artificial*

- Neural Networks, Rhodes, Greece, October 4–7, 2018, Proceedings, Part III*, pages 321–330. 10 2018.
- [Kar] Andrej Karpathy. Cs231n convolutional neural networks for visual recognition. <http://cs231n.github.io/neural-networks-3>.
- [KSH12] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*, NIPS’12, pages 1097–1105, USA, 2012. Curran Associates Inc.
- [LSSK19] Lu Lu, Yeonjong Shin, Yanhui Su, and George Em Karniadakis. Dying relu and initialization: Theory and numerical examples. *CoRR*, abs/1903.06733, 2019.
- [MPA⁺16] Divya Mahajan, Jongse Park, Emmanuel Amaro, Hardik Sharma, Amir Yazdanbakhsh, Joon Kyung Kim, and Hadi Esmaeilzadeh. Tabla: A unified template-based framework for accelerating statistical machine learning. pages 14–26, 03 2016.
- [MR17] Pavlo M. Radiuk. Impact of training set batch size on the performance of convolutional neural networks for diverse datasets. *Information Technology and Management Science*, 20, 12 2017.
- [PGC⁺17] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017.
- [QSX⁺16] Yuran Qiao, Junzhong Shen, Tao Xiao, Qianming Yang, Mei Wen, and Chunyuan Zhang. Fpga-accelerated deep convolutional neural networks for high throughput and energy efficiency: Fpga-accelerated deep convolutional neural networks. *Concurrency and Computation: Practice and Experience*, 05 2016.
- [sm-] The softmax function and its derivative. <https://eli.thegreenplace.net/2016/the-softmax-function-and-its-derivative/>.
- [TYRW14] Yaniv Taigman, Ming Yang, Marc’Aurelio Ranzato, and Lior Wolf. Deepface: Closing the gap to human-level performance in face verification. 09 2014.
- [XHQ⁺16] Yingce Xia, Di He, Tao Qin, Liwei Wang, Nenghai Yu, Tie-Yan Liu, and Wei-Ying Ma. Dual learning for machine translation. *CoRR*, abs/1611.00179, 2016.

- [Xila] Xilinx. *7 Series DSP48E1 Slice User Guide*. Xilinx.
- [Xilb] XilinxWiki. Zynq 2016.4 release. <https://xilinx-wiki.atlassian.net/wiki/spaces/A/pages/18842326/Zynq+2016.4+Release>.
- [ZFL⁺16] Wenlai Zhao, Haohuan Fu, Wayne Luk, Teng Yu, Shaojun Wang, Bo Feng, Yuchun Ma, and Guangwen Yang. F-cnn: An fpga-based framework for training convolutional neural networks. *2016 IEEE 27th International Conference on Application-specific Systems, Architectures and Processors (ASAP)*, pages 107–114, 2016.